точка

банк для предпринимателей
и предприятий

🐼 🐍 ▚ 🥷

# Джун против панды

*дано* ------------------------- *medium/habr/etc* ------------------------- *суровая реальность*

точка

# С чего все начиналось

**Дано:**
- 1 проект, написанный джуном
- 1 месяц до отпуска

**Найти:**
- 1 год данных, которые нужно
  просчитать и сохранить

до отпуска:  30 дней  ❤️❤️❤️

**Дано:**
- 1 проект, написанный джуном
- 1 месяц до отпуска

**Найти:**
- 1 год данных, которые нужно просчитать и сохранить

**Проблема #1**

Данные сохраняются со скоростью: 1 день данных за 1 день

**Проблема #2**

OOM

**Проблема #3**

pandas (медленно)

до отпуска:  30 дней  ❤️❤️❤️

Очень много способов ускорить **Pandas**

Stats&Data ninja · Follow
12 min read · Dec 7, 2022

PYTHON PROGRAMMING: MANIPULATING TABULAR DATA EFFICIENTLY

**How to Speed up Pandas by 100x**

With Great power comes great responsibility.

Pritish Jadhav · Follow

Как ускорить обработку данных в **Pandas** в **600** раз

Nick Komissarenko · Follow
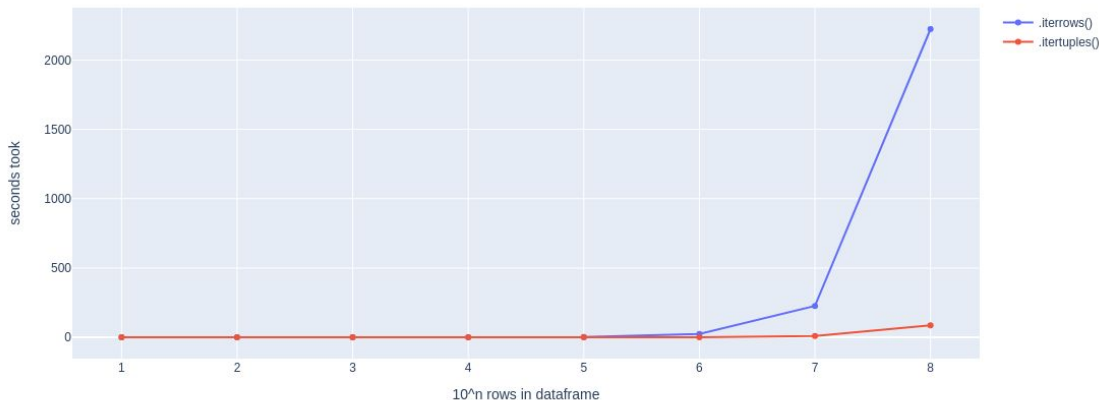5 min read · Sep 22, 2020

точка

# Итерирумся

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.

- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.

- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.
- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

Iteration speed on the number of rows in a dataframe

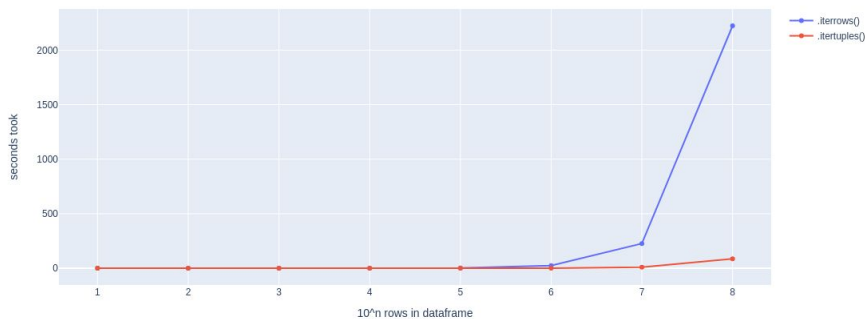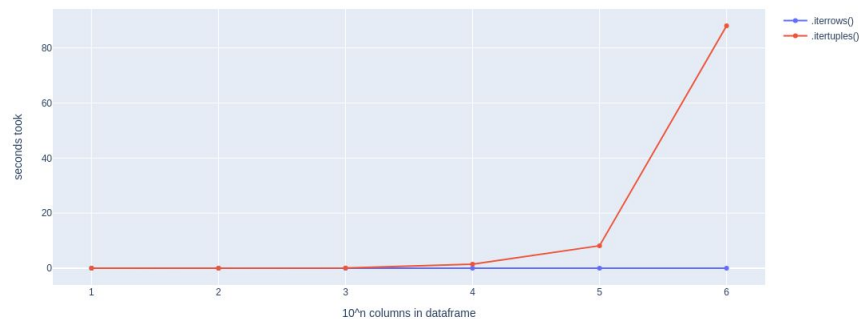To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.
- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

НО!

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.

- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

Iteration speed on the number of rows in a dataframe

Iteration speed on the number of columns in a dataframe

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.
- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

Iteration speed on the number of rows in a dataframe with 1000 columns



при 1000+ колонок - стоит
использовать iterrows

$10^n$ rows in dataframe

```
[5]: %memit pd.DataFrame(np.random.rand(10*6, 1000)).iterrows()

     peak memory: 139.55 MiB, increment: 0.86 MiB

[5]: %memit pd.DataFrame(np.random.rand(10*6, 1000)).itertuples()

     peak memory: 143.27 MiB, increment: 4.67 MiB
```
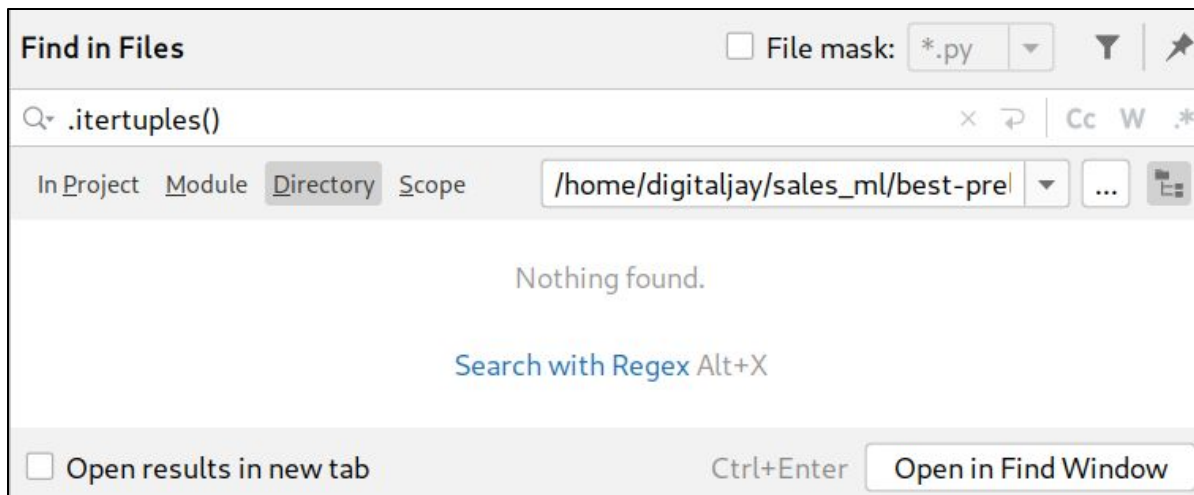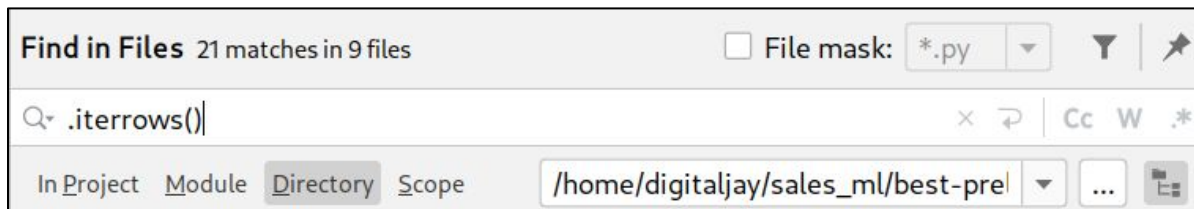
```
[5]:  %memit pd.DataFrame(np.random.rand(10*6, 1000)).iterrows()

      peak memory: 139.55 MiB, increment: 0.86 MiB

[5]:  %memit pd.DataFrame(np.random.rand(10*6, 1000)).itertuples()

      peak memory: 143.27 MiB, increment: 4.67 MiB
```

```
[5]: %memit pd.DataFrame(np.random.rand(10*6, 1000)).iterrows()

     peak memory: 139.55 MiB, increment: 0.86 MiB

[5]: %memit pd.DataFrame(np.random.rand(10*6, 1000)).itertuples()

     peak memory: 143.27 MiB, increment: 4.67 MiB
```

в **5** раз

меньше памяти потребляет
**iterrows** по сравнению с **itertuples**
при больших датафеймах

**1. Stop using iterrows() :**

- Data manipulation often requires iterating over dataframe rows.

- `iterrows()` is often the go-to option for such use cases. However, it is notoriously slow and can be easily swapped by `itertuples()` .

Например, итерация по строкам с помощью метода **.iterrows**(). Это наиболее медленный способ, к тому же не сохраняет типы данных. Другие варианты — использовать .itertuples(), где на каждой итерации строка рассматривается как именованный tupple. Это во много раз быстрее, чем .iterrows(). Еще один аналог — **.iteritems**().

точка

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.

- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

To iterate over the rows of a DataFrame, you can use the following methods:

- `iterrows()` : Iterate over the rows of a DataFrame as (index, Series) pairs. This converts the rows to Series objects, which can change the dtypes and has some performance implications.
- `itertuples()` : Iterate over the rows of a DataFrame as namedtuples of the values. This is a lot faster than `iterrows()` , and is in most cases preferable to use to iterate over the values of a DataFrame.

> ⚠️ **Warning**
>
> Iterating through pandas objects is generally **slow**. In many cases, iterating manually over the rows is not needed and can be avoided with one of the following approaches:
>
> - Look for a *vectorized* solution: many operations can be performed using built-in methods or NumPy functions, (boolean) indexing, …
> - When you have a function that cannot work on the full DataFrame/Series at once, it is better to use `apply()` instead of iterating over the values. See the docs on function application.
> - If you need to do iterative manipulations on the values but performance is important, consider writing the inner loop with cython or numba. See the enhancing performance section for some examples of this approach.
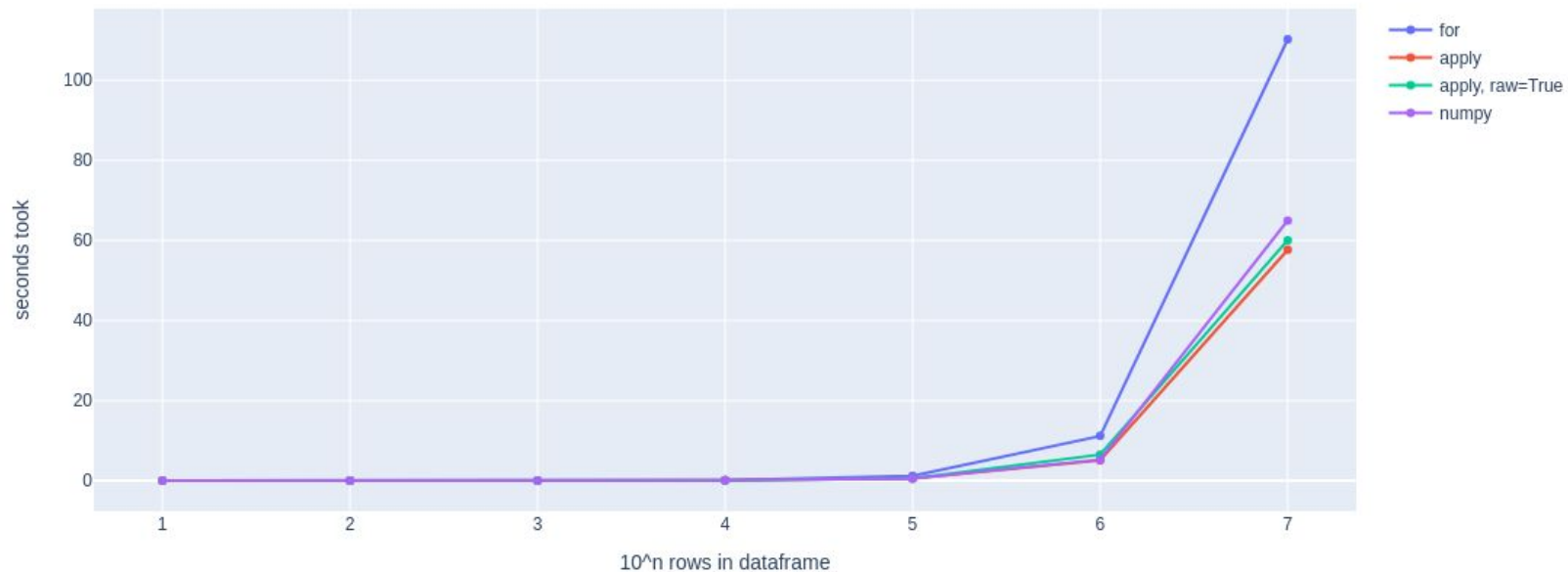
точка

apply?

## Apply

The 'apply' function effectively does the same thing as the loop. It will create a new column titled
reward and apply the calculation function... ...defined by in it. The 'apply' function is
a faster way to run a loop to yo...

**3. apply() is just a glorified for loop:**

- A more traditional way of applying a function to dataframe rows
  involves using `apply()` method.

- Under the hood, `apply()` uses a loop with an added overhead. It can
  often be avoided by leveraging vectorized operations.

**apply?**

Function application speed on the number of rows in a dataframe with 10 columns

# apply?

точка

Function application speed on the number of columns in a dataframe with 10 rows

```
[6]:  # for
      %memit [sum_with_changing_sign(i) for i in df.to_numpy()]
```

peak memory: 219.48 MiB, increment: 2.40 MiB

```
[6]:  # apply
      %memit df.apply(sum_with_changing_sign)
```

peak memory: 218.34 MiB, increment: 3.28 MiB

```
[6]:  # apply, raw=True
      %memit df.apply(sum_with_changing_sign, raw=True)
```

peak memory: 217.59 MiB, increment: 0.72 MiB

```
[6]:  # numpy
      %memit np.sum(df.to_numpy() * ((np.indices(df.shape).sum(axis=0) % 2) + ((np.indices(df.shape).sum(axis=0) % 2)-1))*(-1), axi
```

peak memory: 520.99 MiB, increment: 303.49 MiB

# Function application

To apply your own or another library's functions to pandas objects, you should be aware of the three methods below. The appropriate method to use depends on whether your function expects to operate on an entire `DataFrame` or `Series`, row- or column-wise, or elementwise.

1. Tablewise Function Application: `pipe()`
2. Row or Column-wise Function Application: `apply()`
3. Aggregation API: `agg()` and `transform()`
4. Applying Elementwise Functions: `map()`

точка

# 1.8x speed up

apply -> apply(raw=True)

Ура! Теперь до конца просчета данных - всего **100** дней!

```
df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
```

# Database-style DataFrame or named Series joining/merging

pandas has full-featured, **high performance** in-memory join operations idiomatically very similar to relational databases like SQL. These methods perform 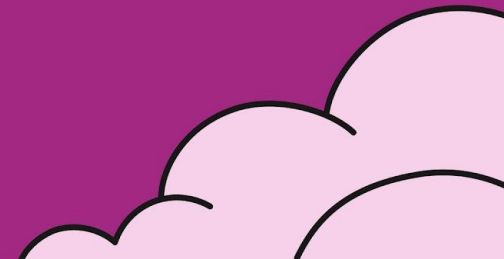significantly better (in some cases well over an order of magnitude better) than other open source implementations (like `base::merge.data.frame` in R). The reason for this is careful algorithmic design and the internal layout of the data in `DataFrame`.

See the cookbook for some advanced strategies.

Users who are familiar with SQL but new to pandas might be interested in a comparison with SQL.

pandas provides a single function, `merge()`, as the entry point for all standard database join operations between `DataFrame` or named `Series` objects:

**точка**

```
178   -              df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
   179 +              df = df.set_index(["id", "ad_login"]).loc[no_phone_key[["id", "ad_login"]], :].reset_index()
```

```
178  -           df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
179  +           df = df.set_index(["id", "ad_login"]).loc[no_phone_key[["id", "ad_login"]], :].reset_index()
```

Merge speed on the number of columns in a dataframe with 10 columns

```
[5]: # merge
     %memit no_phone_key[['id', 'ad_login']].merge(df, on=['id', 'ad_login'], how='left')

     peak memory: 2494.50 MiB, increment: 1046.69 MiB
```

```
[5]: # loc
     %memit df.set_index(['id', 'ad_login']).loc[no_phone_key[['id', 'ad_login']].values.tolist(), :].reset_index()

     peak memory: 2269.16 MiB, increment: 818.97 MiB
```

```
178  -                      df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
   179  +                      df = df.set_index(["id", "ad_login"]).loc[no_phone_key[["id", "ad_login"]], :].reset_index()
```
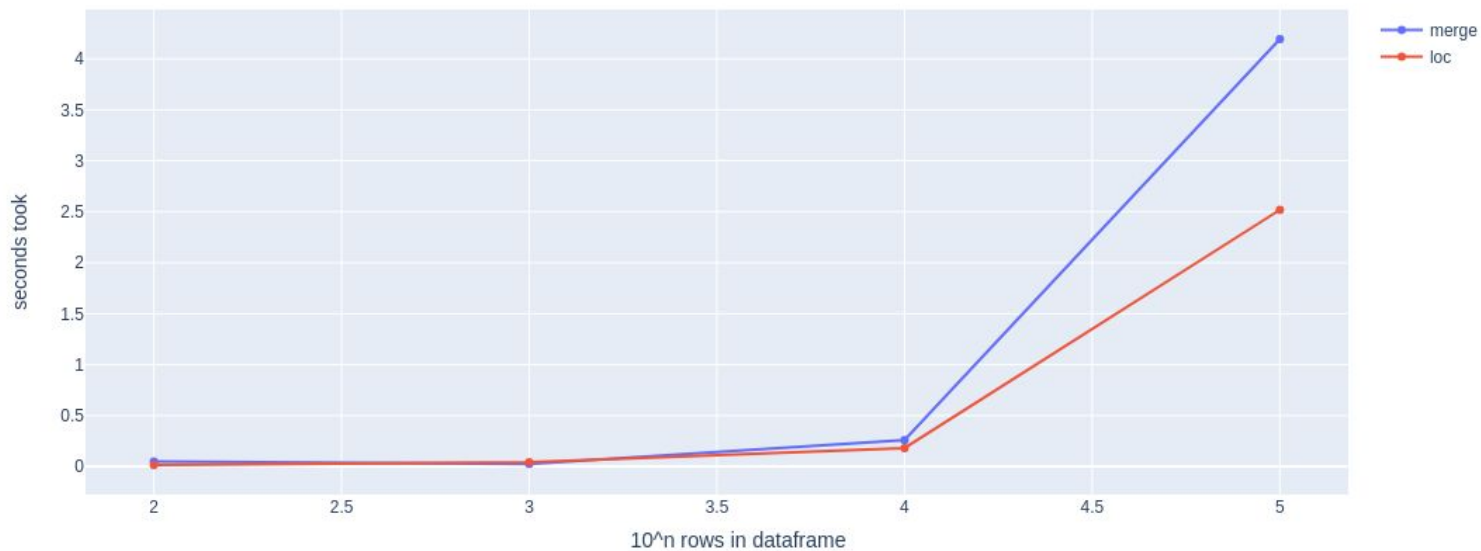
```
178   -                    df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
      179   +              df = df.set_index(["id", "ad_login"]).loc[no_phone_key[["id", "ad_login"]], :].reset_index()
```
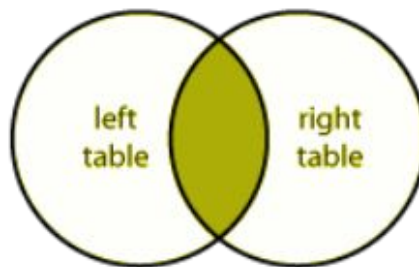
INNER JOIN

```
181  -                    df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
182  -
     180  +                    df = loc_with_fill(
     181  +                        fields_to_loc=["id", "ad_login"],
     182  +                        frame_to_take_as_index=no_phone_key,
     183  +                        frame_to_loc_from=df,
     184  +                        drop_duplicates=False,
     185  +                    )
```

```
134  + def loc_with_fill(
135  +     fields_to_loc: List[str],
136  +     frame_to_take_as_index: DataFrame,
137  +     frame_to_loc_from: DataFrame,
138  +     drop_duplicates=False,
139  + ) -> DataFrame:
140  +     ind_frame = frame_to_take_as_index[fields_to_loc]
141  +     if drop_duplicates:
142  +         ind_frame = ind_frame.drop_duplicates()
143  +     if len(fields_to_loc) > 1:
144  +         reindex = MultiIndex.from_frame(ind_frame)
145  +     else:
146  +         reindex = Index(ind_frame[fields_to_loc[0]].values)
147  +     return frame_to_loc_from.set_index(fields_to_loc).reindex(reindex).reset_index(names=fields_to_loc)
```

```
134  + def loc_with_fill(
135  +     fields_to_loc: List[str],
136  +     frame_to_take_as_index: DataFrame,
137  +     frame_to_loc_from: DataFrame,
138  +     drop_duplicates=False,
139  + ) -> DataFrame:
140  +     ind_frame = frame_to_take_as_index[fields_to_loc]
141  +     if drop_duplicates:
142  +         ind_frame = ind_frame.drop_duplicates()
143  +     if len(fields_to_loc) > 1:
144  +         reindex = MultiIndex.from_frame(ind_frame)
145  +     else:
146  +         reindex = Index(ind_frame[fields_to_loc[0]].values)
147  +     return frame_to_loc_from.set_index(fields_to_loc).reindex(reindex).reset_index(names=fields_to_loc)
```

# pandas.DataFrame.reindex

`DataFrame.reindex(labels=None, *, index=None, columns=None, axis=None, method=None, copy=None, level=None, fill_value=nan, limit=None, tolerance=None)`

Conform DataFrame to new index with optional filling

Places NA/NaN in locations having no value in the pr unless the new index is equivalent to the current one

```
>>> index = ['Firefox', 'Chrome', 'Safari', 'IE10', 'Konqueror']
>>> df = pd.DataFrame({'http_status': [200, 200, 404, 404, 301],
...                    'response_time': [0.04, 0.02, 0.07, 0.08, 1.0]},
...                   index=index)
>>> df
           http_status  response_time
Firefox            200           0.04
Chrome             200           0.02
Safari             404           0.07
IE10               404           0.08
Konqueror          301           1.00
```

Create a new index and reindex the dataframe. By default values in the new index that do not have corresponding records in the dataframe are assigned NaN.

```
>>> new_index = ['Safari', 'Iceweasel', 'Comodo Dragon', 'IE10',
...              'Chrome']
>>> df.reindex(new_index)
                http_status  response_time
Safari                404.0           0.07
Iceweasel               NaN            NaN
Comodo Dragon           NaN            NaN
IE10                  404.0           0.08
Chrome                200.0           0.02
```
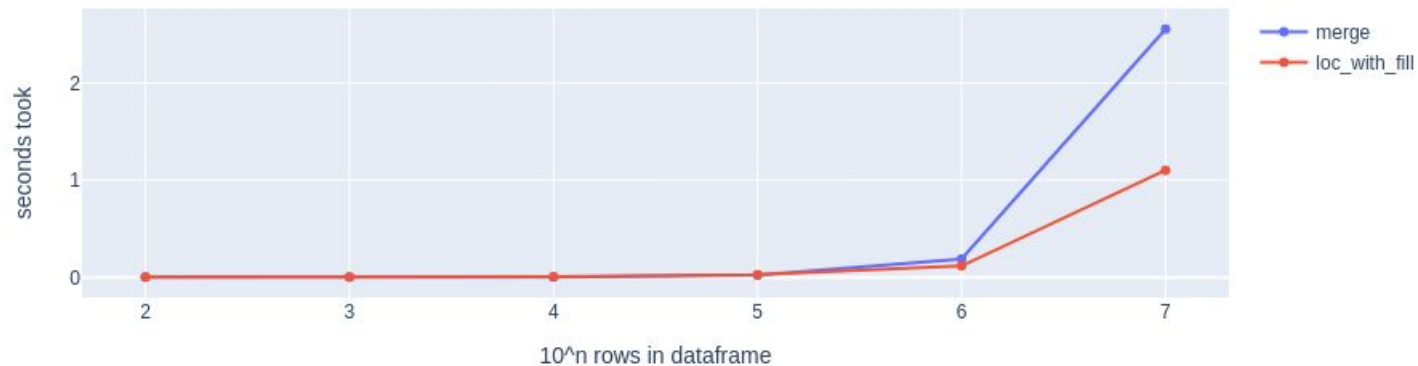
43

```
181  -                    df = no_phone_key[["id", "ad_login"]].merge(df, on=["id", "ad_login"], how="left")
182  -
     180  +                df = loc_with_fill(
     181  +                    fields_to_loc=["id", "ad_login"],
     182  +                    frame_to_take_as_index=no_phone_key,
     183  +                    frame_to_loc_from=df,
     184  +                    drop_duplicates=False,
     185  +                )
```

Merge speed on the number of columns in a dataframe

```
[5]:  # merge
      %memit no_phone_key[['id', 'ad_login']].merge(df, on=['id', 'ad_login'], how='left')

      peak memory: 7195.32 MiB, increment: 96.28 MiB
```

```
[5]:  # loc_with_fill
      %memit loc_with_fill(fields_to_loc=['id', 'ad_login'], frame_to_take_as_index=no_phone_key, frame_to_loc_from=df, drop_duplicat

      peak memory: 7206.66 MiB, increment: 106.60 MiB
```
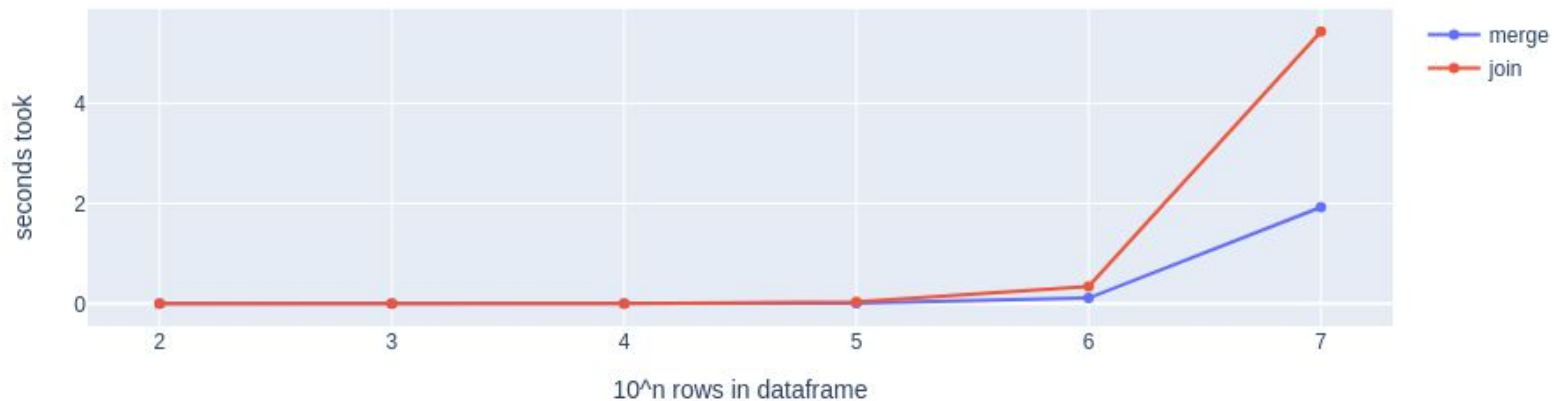
```
193    -                    df = no_phone_key[["id", "ad_login"]].merge(df, how="left", on=["id"])
194    +                    df = no_phone_key[["id", "ad_login"]].set_index("id").join(df.set_index("id"), how="left").reset_index()
```

```
193   -                df = no_phone_key[["id", "ad_login"]].merge(df, how="left", on=["id"])
194   +                df = no_phone_key[["id", "ad_login"]].set_index("id").join(df.set_index("id"), how="left").reset_index()
```

Merge speed on the number of columns in a dataframe

**merge: а правда ли merge - это join?**

точка

```
193  -          df = no_phone_key[["id", "ad_login"]].merge(df, how="left", on=["id"])
194  +          df = no_phone_key[["id", "ad_login"]].set_index("id").join(df.set_index("id"), how="left").reset_index()
```

Merge speed on the number of columns in a dataframe

2696 строк сурсов для
оптимизации

```
[5]:    # merge
        %memit df1.merge(df2, left_on=['df1_join1'], right_on=['df2_join1'], how='inner')

        peak memory: 5105.04 MiB, increment: 3129.31 MiB

[5]:    # join
        %memit df1.set_index('df1_join1').join(df2.set_index('df2_join1'), how='inner').reset_index()

        peak memory: 8457.75 MiB, increment: 6480.49 MiB
```
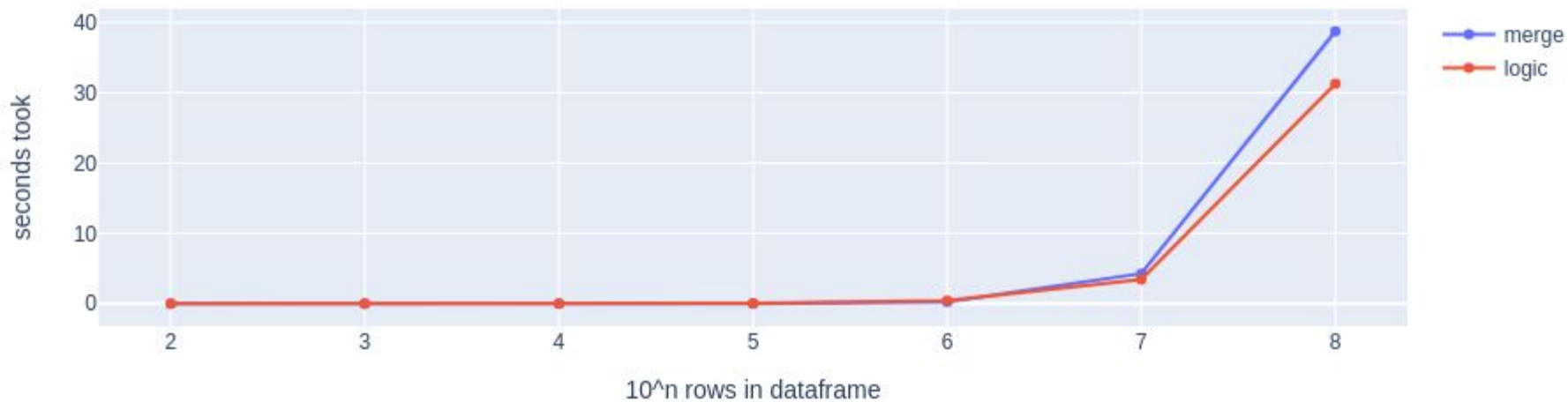
```
37    -        cached_features = validate_prelead_features(cached_features, feature_category)
38    -        cache_presence_indicator = key.merge(
39    -            cached_features[index_cols].drop_duplicates(),
40    -            indicator=True,
41    -            how="left",
42    -            on=index_cols,
43    -        )
      38  +        cached_features = validate_prelead_features(cached_features, feature_category).set_index(index_cols)
      39  +
      40  +        key = key.set_index(index_cols)
      41  +        if len(cached_features) > 0:
      42  +            keys_in_cache = (
      43  +                key.copy().loc[key.index.intersection(cached_features.index), :].reset_index(names=index_cols)
      44  +            )
      45  +            keys_not_in_cache = (
      46  +                key.copy().loc[key.index.difference(cached_features.index), :].reset_index(names=index_cols)
      47  +            )
      48  +        else:
      49  +            keys_in_cache = pd_DataFrame(columns=key.columns)
      50  +            keys_not_in_cache = key.copy().reset_index(names=index_cols)
44    51
45    -        keys_in_cache = key[(cache_presence_indicator["_merge"] == "both").values]
46    -        keys_not_in_cache = key[(cache_presence_indicator["_merge"] == "left_only").values]
      52  +        key.reset_index(names=index_cols, inplace=True)
      53  +        cached_features.reset_index(names=index_cols, inplace=True)
```

**точка**

```
37   -        cached_features = validate_prelead_features(cached_features, feature_category)
38   -        cache_presence_indicator = key.merge(
39   -            cached_features[index_cols].drop_duplicates(),
40   -            indicator=True,
41   -            how="left",
42   -            on=index_cols,
43   -        )
     38  +  cached_features = validate_prelead_features(cached_features, feature_category).set_index(index_cols)
     39  +
     40  +  key = key.set_index(index_cols)
     41  +  if len(cached_features) > 0:
     42  +      keys_in_cache = (
     43  +          key.copy().loc[key.index.intersection(cached_features.index), :].reset_index(names=index_cols)
     44  +      )
     45  +      keys_not_in_cache = (
     46  +          key.copy().loc[key.index.difference(cached_features.index), :].reset_index(names=index_cols)
     47  +      )
     48  +  else:
     49  +      keys_in_cache = pd_DataFrame(columns=key.columns)
     50  +      keys_not_in_cache = key.copy().reset_index(names=index_cols)
44   51
45   -        keys_in_cache = key[(cache_presence_indicator["_merge"] == "both").values]
46   -        keys_not_in_cache = key[(cache_presence_indicator["_merge"] == "left_only").values]
     52  +  key.reset_index(names=index_cols, inplace=True)
     53  +  cached_features.reset_index(names=index_cols, inplace=True)
```

## Merge speed on the number of columns in a dataframe

```
[6]:  # merge
      %memit merge_logic(key, cached_features, index_cols)
```

peak memory: 7243.54 MiB, increment: 83.61 MiB

```
[6]:  # logic with intersection
      %memit intersect_logic(key, cached_features, index_cols)
```

peak memory: 7234.76 MiB, increment: 79.11 MiB

**Find in Files** 56 matches in 24 files  ☐ File mask: *.py ▾ ▼ ▼ | 📌

🔍▾ .merge( ✕ ↵ | Cc W .*

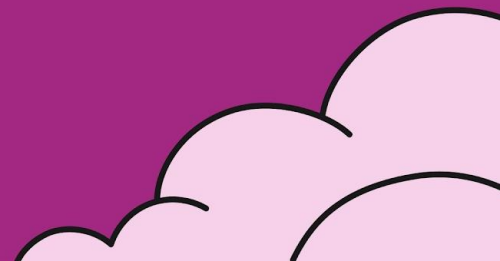In Project   Module   Directory   Scope   /home/digitaljay/sales_ml/best-pre| ▼   ...   🗁

точка

типы данных

| Kind of Data | Data Type | Scalar | Array | String Aliases |
|---|---|---|---|---|
| tz-aware datetime | `DatetimeTZDtype` | `Timestamp` | `arrays.DatetimeArray` | `'datetime64[ns, <tz>]` |
| Categorical | `CategoricalDtype` | (none) | `Categorical` | `'category'` |
| period (time spans) | `PeriodDtype` | `Period` | `arrays.PeriodArray` `'Period[<freq>]'` | `'period[<freq>]'`, |
| sparse | `SparseDtype` | (none) | `arrays.SparseArray` | `'Sparse'`, `'Sparse[in` `'Sparse[float]'` |
| intervals | `IntervalDtype` | `Interval` | `arrays.IntervalArray` | `'interval'`, `'Interva` `'Interval[<numpy_dtyp` `'Interval[datetime64` `<tz>]]'`, `'Interval[timedelta6` |
| nullable integer | `Int64Dtype`, … | (none) | `arrays.IntegerArray` | `'Int8'`, `'Int16'`, `'In` `'Int64'`, `'UInt8'`, `'U` `'UInt32'`, `'UInt64'` |
| `nullable` `float` | `Float64Dtype`, … | (none) | `arrays.FloatingArray` | `'Float32'`, `'Float64'` |
| Strings | `StringDtype` | `str` | `arrays.StringArray` | `'string'` |
| Boolean (with NA) | `BooleanDtype` | `bool` | `arrays.BooleanArray` | `'boolean'` |

Calculation speed on the number of rows in a dataframe

Calculation speed on the number of rows in a dataframe

Addressing a specific part of your quote:

> For example, when reading a 16-bit value on a 64-bit machine, a full 64 bits worth of data must still be read from memory. The desired 16-bit field then has to be masked off and possibly shifted into place within the destination register.

This is true for some architectures, but not all of them. Particularly, the x86-64 architecture has instructions for working directly with 16-bit values, so doing so on that architecture will not require mask and shift operations. Futhermore, while 64 bits[1] will still be fetched from memory, if subsequent instructions need to use the data in the other 48 bits they will be able to do so without needing to access main memory again due to the processor's cache. For this architecture specifically, therefore, this advice is wrong.

[1]: or, more likely, 128 bits, as this is the width of the cache on most modern processors

Share   Improve this answer   Follow

answered Jun 15, 2020 at 23:00

occipita

**209** ● 2 ● 5

```
[3]:  # int64
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int64')>50).sum(axis=1)

      peak memory: 997.96 MiB, increment: 852.95 MiB
```

```
[3]:  # uint64
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint64')>50).sum(axis=1)

      peak memory: 1755.40 MiB, increment: 1611.11 MiB
```

```
[3]:  # int32
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int32')>50).sum(axis=1)

      peak memory: 1378.75 MiB, increment: 1234.32 MiB
```

```
[3]:  # uint32
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint32')>50).sum(axis=1)

      peak memory: 1377.25 MiB, increment: 1234.32 MiB
```

```
[3]:  # int16
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int16')>50).sum(axis=1)

      peak memory: 1181.45 MiB, increment: 1038.37 MiB
```

```
[3]:  # uint16
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint16')>50).sum(axis=1)

      peak memory: 1173.63 MiB, increment: 1030.44 MiB
```

```
[3]:  # int8
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int8')>50).sum(axis=1)

      peak memory: 1092.50 MiB, increment: 950.10 MiB
```

```
[3]:  # uint8
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint8')>50).sum(axis=1)

      peak memory: 1095.40 MiB, increment: 951.97 MiB
```

```
[3]:  # int64
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int64')>50).sum(axis=1)

      peak memory: 997.96 MiB, increment: 852.95 MiB

[3]:  # uint64
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint64')>50).sum(axis=1)

      peak memory: 1755.40 MiB, increment: 1611.11 MiB

[3]:  # int32
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int32')>50).sum(axis=1)

      peak memory: 1378.75 MiB, increment: 1234.32 MiB

[3]:  # uint32
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint32')>50).sum(axis=1)

      peak memory: 1377.25 MiB, increment: 1234.32 MiB

[3]:  # int16
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int16')>50).sum(axis=1)

      peak memory: 1181.45 MiB, increment: 1038.37 MiB

[3]:  # uint16
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint16')>50).sum(axis=1)

      peak memory: 1173.63 MiB, increment: 1030.44 MiB

[3]:  # int8
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int8')>50).sum(axis=1)

      peak memory: 1092.50 MiB, increment: 950.10 MiB

[3]:  # uint8
      %memit (pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint8')>50).sum(axis=1)

      peak memory: 1095.40 MiB, increment: 951.97 MiB
```
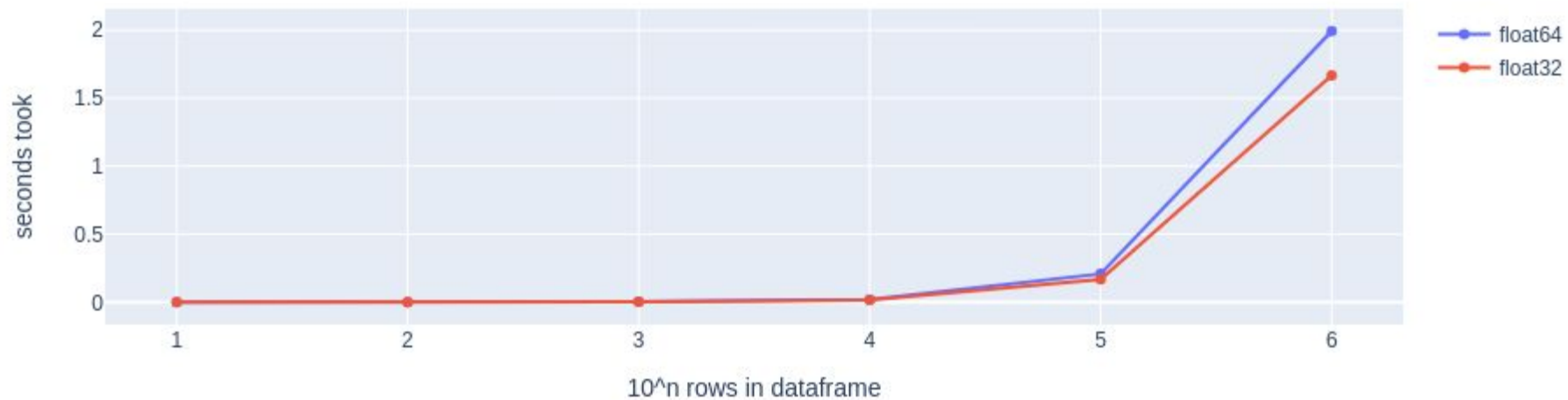
```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 100)), dtype='int64')
```

```python
# int64
%memit df
```
peak memory: 7774.91 MiB, increment: 0.31 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint64')
```

```python
# uint64
%memit df
```
peak memory: 904.18 MiB, increment: 0.06 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int32')
```

```python
# int32
%memit df
```
peak memory: 537.11 MiB, increment: 0.04 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint32')
```

```python
# uint32
%memit df
```
peak memory: 521.03 MiB, increment: 0.04 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int16')
```

```python
# int16
%memit df
```
peak memory: 332.25 MiB, increment: 0.04 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint16')
```

```python
# uint16
%memit df
```
peak memory: 334.51 MiB, increment: 0.03 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='int8')
```

```python
# int8
%memit df
```
peak memory: 234.75 MiB, increment: 0.04 MiB

```python
df = pd.DataFrame(np.random.randint(100, size=(10**7, 10)), dtype='uint8')
```

```python
# uint8
%memit df
```
peak memory: 240.97 MiB, increment: 0.04 MiB

Calculation speed on the number of rows in a dataframe

```
[3]:  # float64
      %memit pd.DataFrame(np.random.rand(10**7, 10), dtype='float64')

      peak memory: 880.34 MiB, increment: 739.81 MiB
```
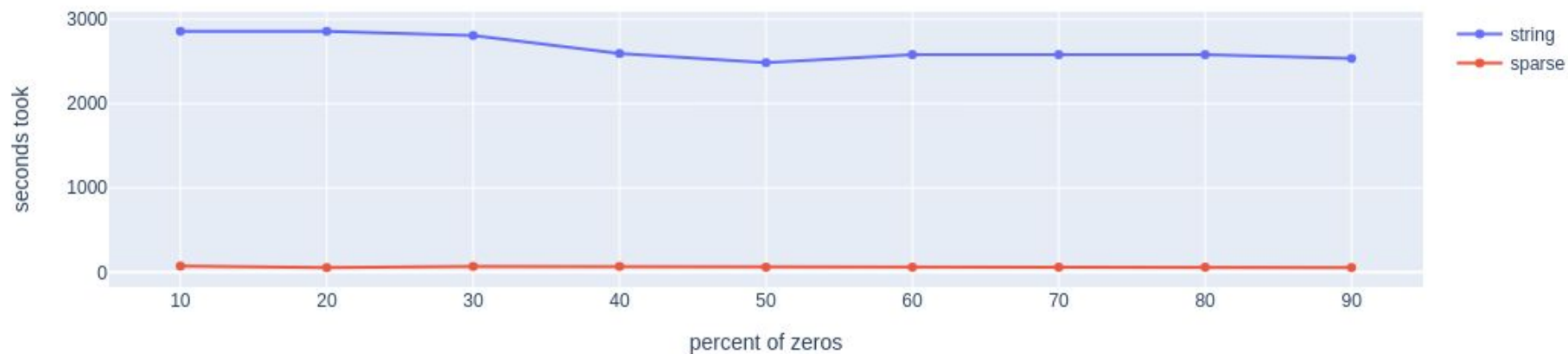
```
[3]:  # float32
      %memit pd.DataFrame(np.random.rand(10**7, 10), dtype='float32')

      peak memory: 1284.04 MiB, increment: 1139.30 MiB
```

```
38    # без этой эвристики начинает конвертировать всякие id'шники, а это жрёт оперативку и время
39    # TODO: Сам поправь, критично, никогда не бывает больше 1
40    categorical_can_be_used = (len(col.unique()) / (len(col) + 1) >= 3) and colname in categ_cols
```

group&count speed depending on the number of unique values on dataframe with shape (10**8, 10)

sum speed depending on the percent of zeros in dataframe with shape (10**8, 10)

```
[3]:  # string
      %memit pd.DataFrame(np.random.randint(0, 10**6, (10**6, 10)), dtype='str')

      peak memory: 903.05 MiB, increment: 764.17 MiB
```

```
[3]:  # category
      %memit pd.DataFrame(np.random.randint(0, 10**6, (10**6, 10)), dtype='category')

      peak memory: 517.06 MiB, increment: 376.34 MiB
```

```
[3]:  # string
      %memit pd.DataFrame(np.random.randint(0, 10, (10**6, 10)), dtype='str')

      peak memory: 900.20 MiB, increment: 757.54 MiB
```

```
[3]:  # sparse
      %memit pd.DataFrame(np.random.randint(0, 10, (10**6, 10)), dtype='Sparse[string]')

      peak memory: 344.79 MiB, increment: 203.62 MiB
```
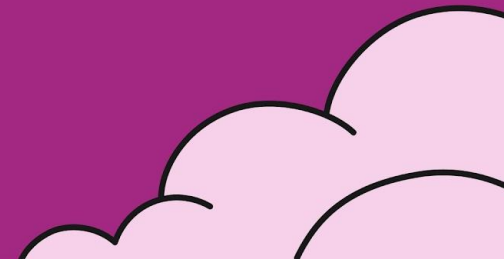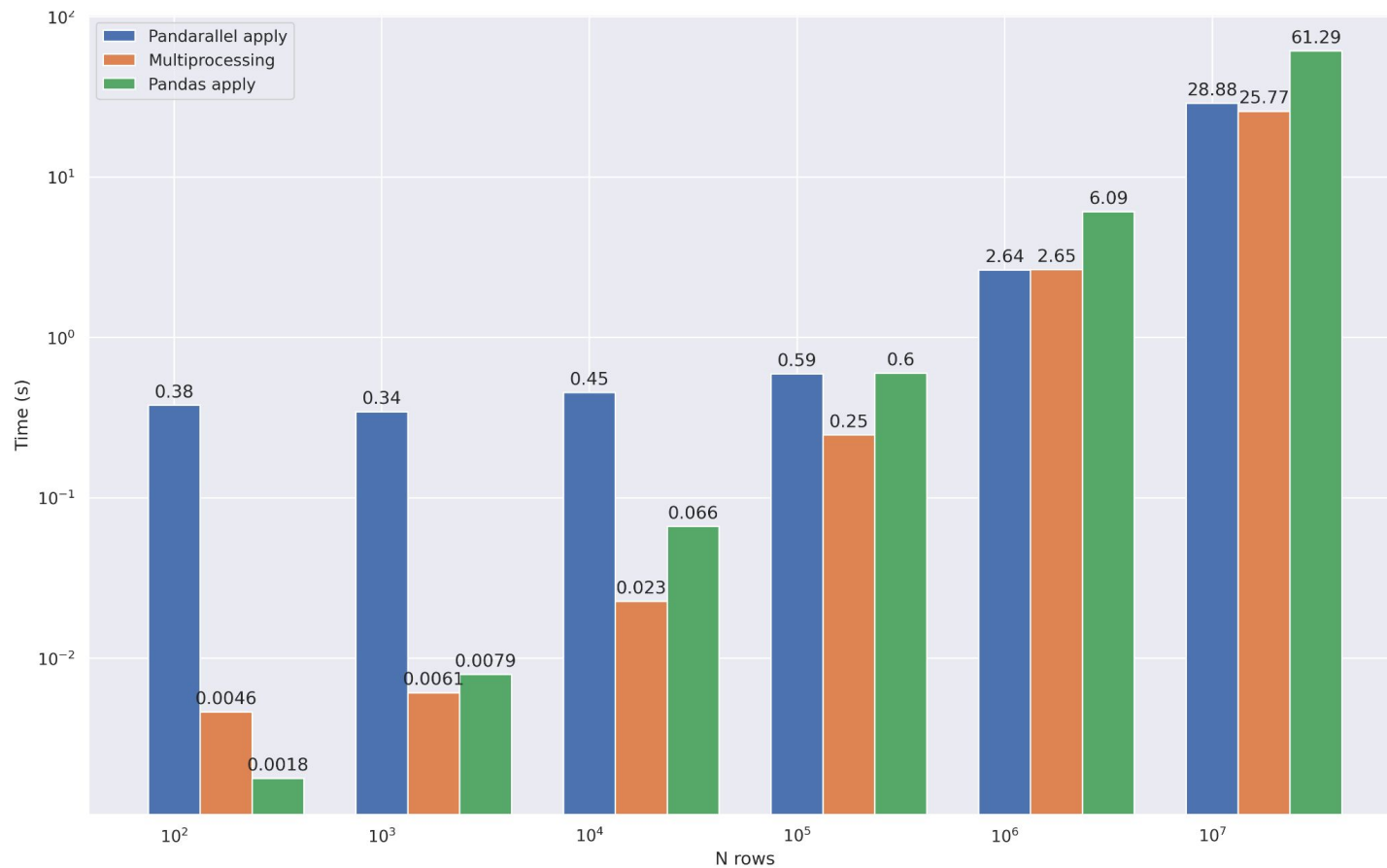
# 2.5x speed up

int64, float64 -> int8, float32

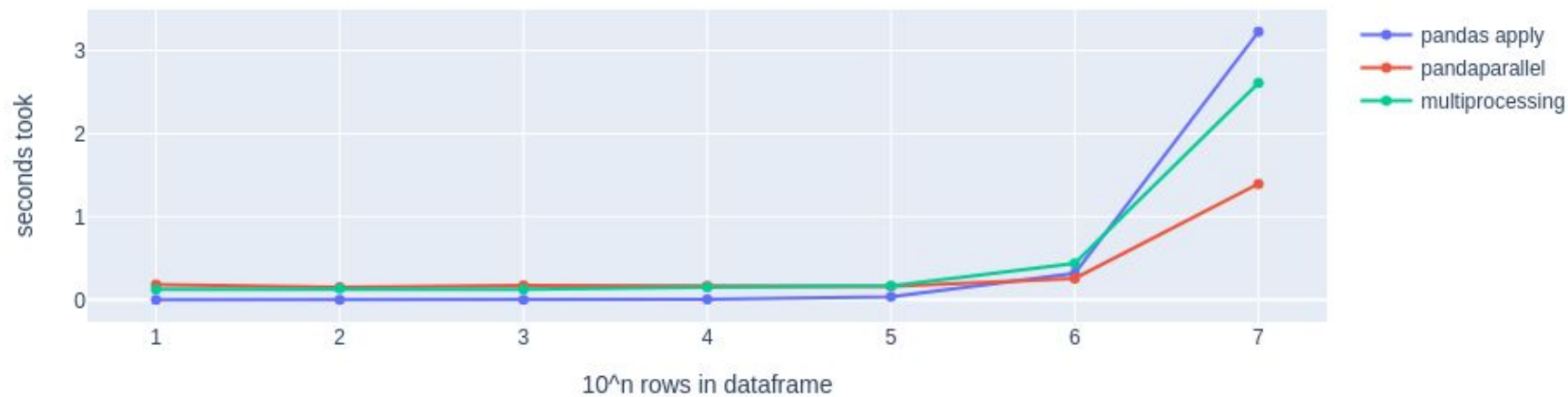Ура! Теперь до конца просчета данных - всего **20** дней!

точка

parallel

Parallelisation

```
[7]: def apply_func(i):
         return i**2
```

```
[8]: df = pd.DataFrame(np.random.rand(10**7, 1), columns=["sample_column"])
```

```
[9]: # pandas apply
     %memit df['sample_column'].apply(apply_func)
```

peak memory: 1195.19 MiB, increment: 978.40 MiB

```
[9]: # pandaparallel
     %memit df['sample_column'].parallel_apply(apply_func)
```

peak memory: 352.48 MiB, increment: 137.52 MiB

```
[9]: # multiprocessing
     %memit mp.Pool().map(apply_func, df['sample_column'])
```

peak memory: 599.39 MiB, increment: 382.56 MiB

**точка**

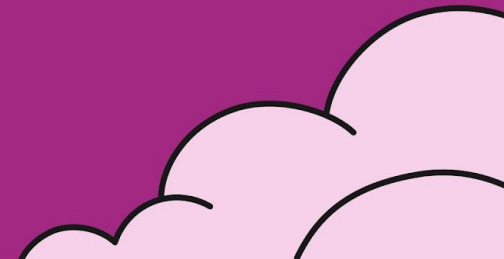| Without parallelization | With parallelization |
|---|---|
| `df.apply(func)` | `df.parallel_apply(func)` |
| `df.applymap(func)` | `df.parallel_applymap(func)` |
| `df.groupby(args).apply(func)` | `df.groupby(args).parallel_apply(func)` |
| `df.groupby(args1).col_name.rolling(args2).apply(func)` | `df.groupby(args1).col_name.rolling(args2).parallel_` |
| `series.map(func)` | `series.parallel_map(func)` |
| `series.apply(func)` | `series.parallel_apply(func)` |
| `series.rolling(args).apply(func)` | `series.rolling(args).parallel_apply(func)` |

# 3x speed up

pandaparallel
Ура! Теперь до конца просчета данных - всего **6** дней!

ПОБЕДА!

точка

отпуск

1.  iterrows ----------------> itertuples (если кол-во колонок <1000)

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

3. merge ------------------> loc/reindex

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

3. merge ------------------> loc/reindex

4. int64, float64 ---------> int8, float32

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

3. merge -----------------> loc/reindex

4. int64, float64 ---------> int8, float32

5. string ------------------> sparse/categorical

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

3. merge -----------------> loc/reindex

4. int64, float64 ---------> int8, float32

5. string ------------------> sparse/categorical

6. ------------------------> pandaparallel

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ------------------> apply(raw=True)

3. merge -----------------> loc/reindex

4. int64, float64 ---------> int8, float32

5. string ------------------> sparse/categorical

6. --------------------------> pandaparallel

1. iterrows ----------------> itertuples (если кол-во колонок <1000)

2. apply ----------------⚡> apply(raw=True)

3. merge ------------------> loc/reindex

4. int64, float64 ----⚡-⚡> int8, float32

5. string -------------------> sparse/categorical

6. ---------------------⚡-⚡> pandaparallel

⚡  быстрее всего добавить

1. iterrows ----------------> itertuples (если кол-во колонок <1000)
2. apply ----------------⚡> apply(raw=True)
3. merge ------------------> loc/reindex
4. int64, float64 ----⚡-⚡> int8, float32
5. string -------------------> sparse/categorical
6. ----------------------⚡-⚡> pandaparallel

⚡ быстрее всего добавить

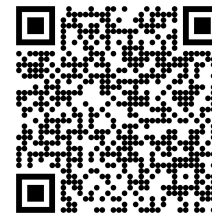🗄 сокращают/не меняют потребление памяти

точка

1. Документация != последняя инстанция

2. Не верь статьям без бенчмарков

3. Знание своих данных позволяет писать более оптимальный код

# Елизавета Пушкарева

Data Scientist в Точке

tg: @digitaljay

бенчмарки

JET BRAINS

Blog

# When should you stick with pandas?

All of this sounds so amazing that you're probably wondering why you would even bother with pandas anymore. Not so fast! While Polars is superb for doing extremely efficient data transformations, it is currently not the optimal choice for data exploration or for use as part of machine learning pipelines. These are areas where pandas continues to shine.