

An introduction to Spark GraphFrame with examples analyzing the Wikipedia link graph

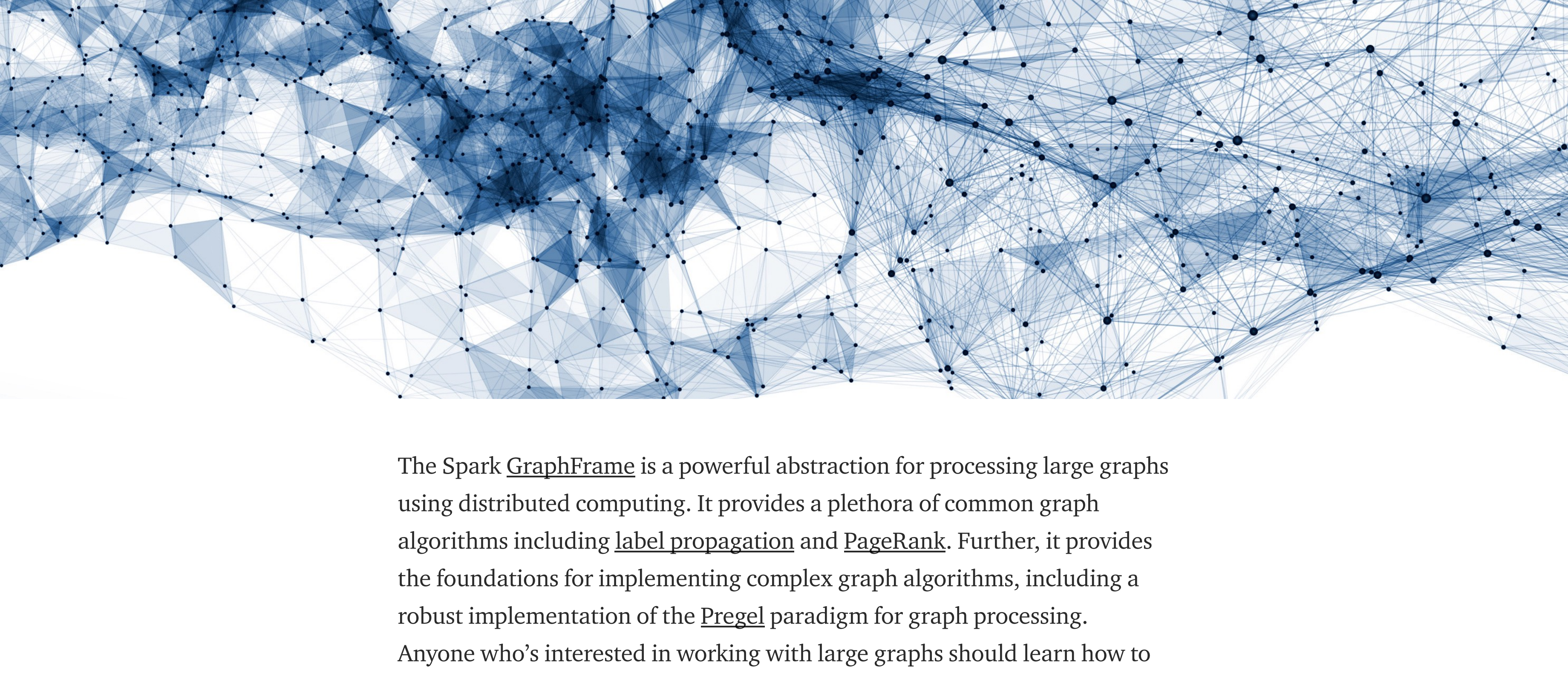
The Spark GraphFrame is an incredibly powerful tool for performing distributed computations with large graphical data. This article introduces the GraphFrame abstraction and shows how it can be leveraged to analyze the graph formed by the links between Wikipedia articles.



Matt Hagy

Follow

Feb 22 · 4 min read



The Spark **GraphFrame** is a powerful abstraction for processing large graphs using distributed computing. It provides a plethora of common graph algorithms including **label propagation** and **PageRank**. Further, it provides the foundations for implementing complex graph algorithms, including a robust implementation of the **Pregel** paradigm for graph processing. Anyone who's interested in working with large graphs should learn how to apply this powerful tool.

In this article, I'll introduce you to the basic of GraphFrame and demonstrate how to use this tool through several examples. These examples consider the link graph between Wikipedia articles and I demonstrate how to analyze this graph by leveraging the GraphFrame abstraction.

Note that the examples in this post build off some more elementary Spark concepts such as DataFrames. Additionally, it uses basic Scala code to demonstrate algorithms. Only small and simple examples are shown so that one doesn't need to be well-familiar with these concepts to learn about the power of the GraphFrame abstraction.

Let's dive right in and consider how to create a GraphFrame. I'll start by introducing the Wikipedia link data in a basic Spark RDD. Each element of this RDD is a single Wikipedia article page represented with the following Scala class.

```
1 case class Page(title: String, links: Seq[String])
page.scala hosted with ❤ by GitHub view raw
```

Note that each link is the title of another page. I.e., each page knows all the other pages that it links to by title.

We begin with a single RDD `pages` that contains 19,093,693 Wikipedia article pages. From that, we generate two Spark DataFrames, one consisting of the vertices (i.e., page nodes) and the other consisting of the directed edges.

```
1 val vertices = pages.map(_._title).toDF("id")
2
3 case class Edge(src: String, dst: String)
4 val edges = (
5   pages
6   .flatMap(p => p.links.map(l => Edge(p.title, l)))
7   .toDF()
8 )
gf_vertices_and_edges.scala hosted with ❤ by GitHub view raw
```

Note that there are 206,181,091 directed edges in this graph.

Next, we create a GraphFrame using these two DataFrames.

```
1 val gf = GraphFrame(vertices, edges)
graphframe.scala hosted with ❤ by GitHub view raw
```

And that's all we have to do to access this powerful abstraction. We can now start using some of the builtin graph algorithms to analyze the Wikipedia link graph.

Let start by computing something simple: the Wikipedia pages with the largest number of outbound links. To this end, we can use the GraphFrame method `outDegrees`, which is a computed DataFame that corresponds to the number of outbound edges for each vertex. Since it's a DataFrame, we can use the `orderBy` method and `Limit` to select the top 10.

```
1 gf.outDegrees.orderBy($"outDegree".desc).limit(10)
gf_out_degrees.scala hosted with ❤ by GitHub view raw
```

This gives the following results.

id	outDegree
Wikipedia:WikiProject Medicine/Lists of pages/Articles	34140
Wikipedia:WikiProject Medicine/Popular pages En 2013b	29071
Draft:List of the prehistoric life of Germany	24747
Draft:List of the prehistoric life of France	24533
Wikipedia:0.8/Second half	24151
Wikipedia:0.8/First half	23156
Wikipedia:WikiProject Mountains/List of mountains	22767
Wikipedia:Database reports/Broken section anchors	22723
Wikipedia:WikiProject Curling/Curling (red links)	19765
Wikipedia:WikiProject Anatomy/Lists of pages/Articles	19661

out_degree.csv hosted with ❤ by GitHub view raw

Interestingly, we can see that many special "Wikipedia:"-prefix pages have the highest number of outbound links.

Next, let's consider the number of inbound edges for each page and find the top 10 linked-to pages using the corresponding GraphFrame method, `inDegrees`.

```
1 gf.inDegrees.orderBy($"inDegree".desc).limit(10)
in_degrees.scala hosted with ❤ by GitHub view raw
```

id	inDegree
United States	343728
France	151180
Animal	148740
India	138720
World War II	134905
Germany	129148
United Kingdom	116575
New York City	116016
Arthropod	112928
England	110900

in_degrees.csv hosted with ❤ by GitHub view raw

We can see that locations, including countries and cities, are among the most heavily linked to pages.

Now let's explore some more complex graph computations. Let's consider the most heavily-linked to article, "United States", and find the other articles in the link graph that are furthest away from this article in terms of the number of links you have to follow to arrive at "United States".

To this end, we can use the GraphFrame method, `shortestPaths`, which takes a collection of landmark vertices and returns the path length for every vertex in the graph to every landmark vertex.

```
1 val us_distance = gf.shortestPaths.landmarks(List("United States")).run()
2
3 // Exclude articles that have a colon in their titles
4 val stringContainsColon = udf((s: String) => s.contains(":"))
5
6 // Unpack the map of landmark -> path length
7 val getUsDistance = udf((m: Map[String, Int]) => m.getOrElse("United States", -1))
8
9 (us_distance
10  .filter(!stringContainsColon($"id"))
11  .withColumn("us_distance", getUsDistance($"distances"))
12  .filter($"us_distance" > 0)
13  .orderBy($"us_distance".desc)
14  .limit(10))
wiki_us_distance.scala hosted with ❤ by GitHub view raw
```

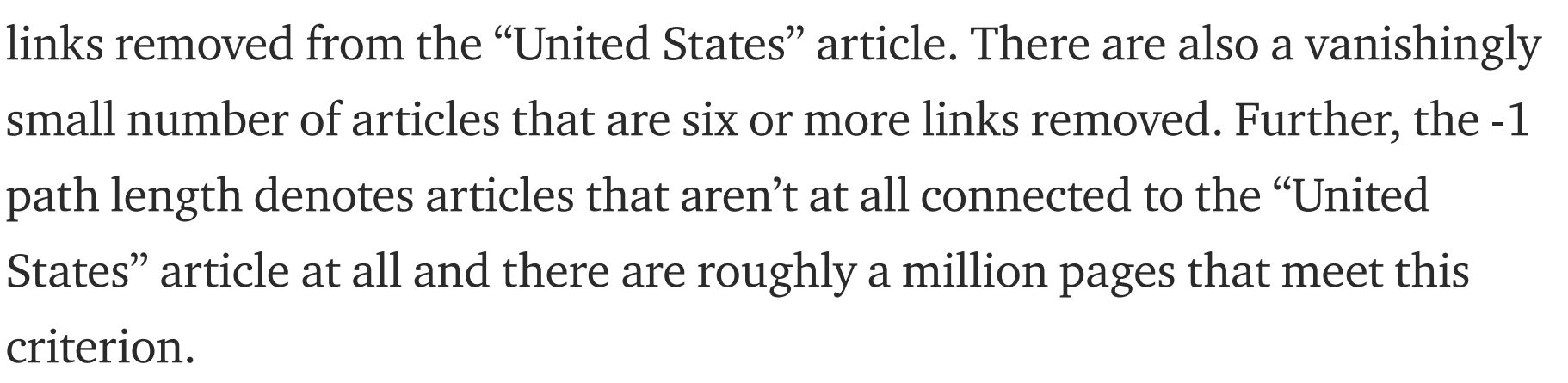
id	us_distance
Lumbocostal arches	10
Arcus lumbocostalis	10
Lumbocostal arch (disambiguation)	10
Dorsal scapular vessels	9
Lacrimal fossa (disambiguation)	9
Lumbocostal arch	9
Dorsal scapular (disambiguation)	9
Successive pairs	9
Lone chooser	9
The Art of Kissing Properly (album)	9

us_distance.csv hosted with ❤ by GitHub view raw

It is interesting that there are only three articles that are a full 10 links removed from the "United States" article in terms of shortest path length.

Let's go one step further and compute the number of pages at each path length. I.e., how many Wikipedia articles are 1 link removed from "United States", how many are 2 links removed, etc.

```
1 display(
2   us_distance
3   .filter(!stringContainsColon($"id"))
4   .withColumn("us_distance", getUsDistance($"distances"))
5   .groupBy($"us_distance")
6   .count()
7   .orderBy($"us_distance")
8 )
us_dist_distance.scala hosted with ❤ by GitHub view raw
```



(Note, that I'm using the built-in `display` function of the **Databricks** notebook. I'll talk more about the joy of using Databricks in a future post.)

Here we can see that there are millions of pages that are between 2 and 4 links removed from the "United States" article. There are also a vanishingly small number of articles that are six or more links removed. Further, the -1 path length denotes articles that aren't at all connected to the "United States" article at all and there are roughly a million pages that meet this criterion.

Lastly, no demonstration of GraphFrame would be complete without showing how easily it can be used to perform the PageRank algorithm.

```
1 val pr = gf.pageRank.resetProbability(0.15).tol(0.01).run()
2 pr.vertices.orderBy($"pagerank".desc).limit(10)
wiki_page_rank.scala hosted with ❤ by GitHub view raw
```

The page rank of the top 10 articles is:

id	pagerank
United States	31393.4
World War II	14776.6
France	13613.1
United Kingdom	13282.6
The New York Times	12575.4
Germany	12129.1
Race and ethnicity in the United States Census	11639.6
List of sovereign states	10892.6
New York City	10423.3
India	10250.5

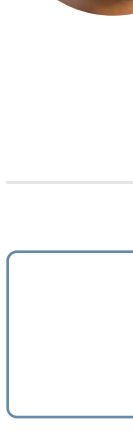
wiki_page_rank.csv hosted with ❤ by GitHub view raw

It may not be surprising that "United States" has the highest vertex rank as we know it's the most heavily linked-to article.

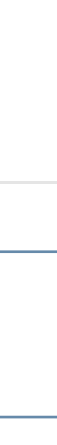
I hope these examples have helped to convince you that GraphFrame is a powerful abstraction for performing distributed computation on large graphs. There are additional more advanced concepts that I hope to share in future articles, including how to implement custom graph algorithms using the Pregel compute paradigm with GraphFrame.


Update: If you've found the Scala examples in this article interesting and would like to learn more about this powerful programming language you can check out my article, [Quickly learning the basics of Scala through Structure and Interpretation of Computer Programs examples \(Part 1\)](#).

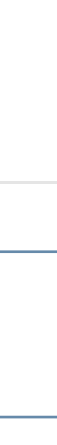
Spark Graphframe Data Science Analysis Graph




334 claps










The best of Towards Data Science, in your inbox.

Sign up for a free account to follow Towards Data Science. You'll get their best stories delivered to your email, personalized for you.

Follow



WRITTEN BY

Matt Hagy


Software Engineer and fmr. Data Scientist and Manager, Ph.D. in Computational Statistical Chemistry. (matthagy.com)

Follow

See responses (1)

More From Medium

More from Towards Data Science



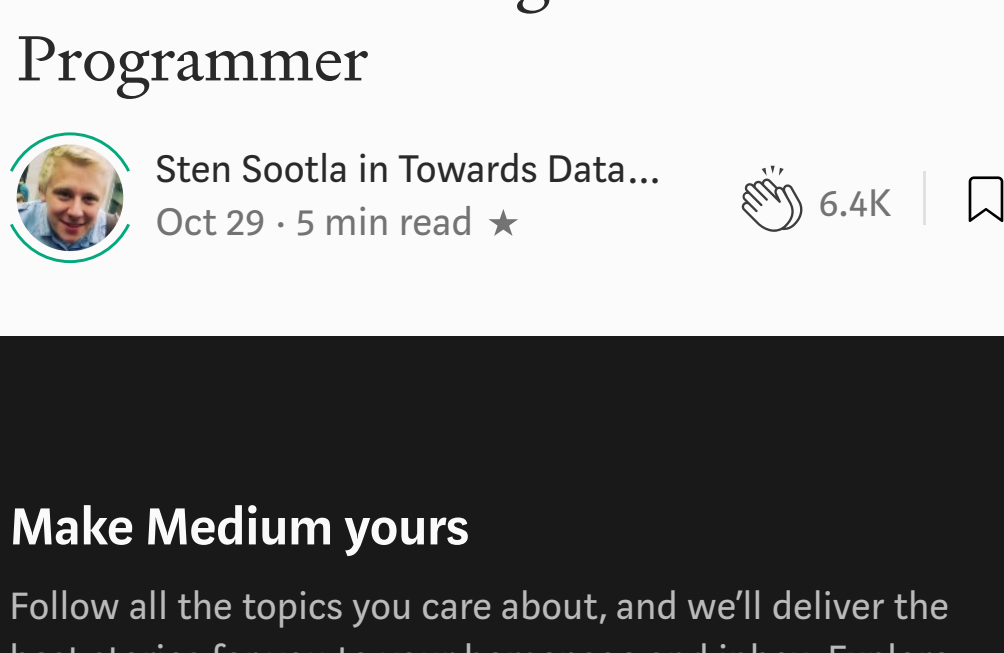
Want a data science job? Use the weekend project principle to get it

Daniel Bourke in Towards Data...

Nov 2 · 4 min read

★

More from Towards Data Science



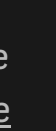
How To Fake Being a Good Programmer

Sten Sootla in Towards Data...

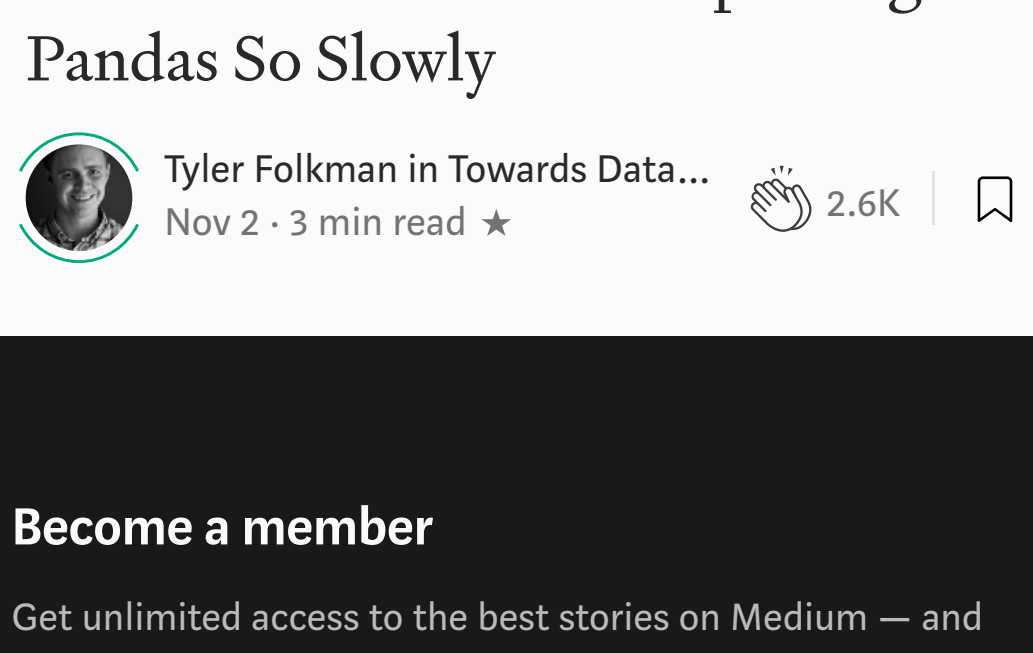
Oct 29 · 5 min read

★

6.4K



More from Towards Data Science



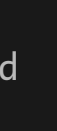
One Word of Code to Stop Using Pandas So Slowly

Tyler Folkman in Towards Data...

Nov 2 · 3 min read

★

2.6K



Discover Medium

Welcome to a place where words matter. On Medium, smart voices and original ideas take center stage - with no ads in sight. [Watch](#)

Make Medium yours

Follow all the topics you care about, and we'll deliver the best stories for you to your homepage and inbox. [Explore](#)

Become a member

Get unlimited access to the best stories on Medium — and support writers while you're at it. Just \$5/month. [Upgrade](#)

Medium

About Help Legal