# Apache Spark Usage in the Open Source Ecosystem
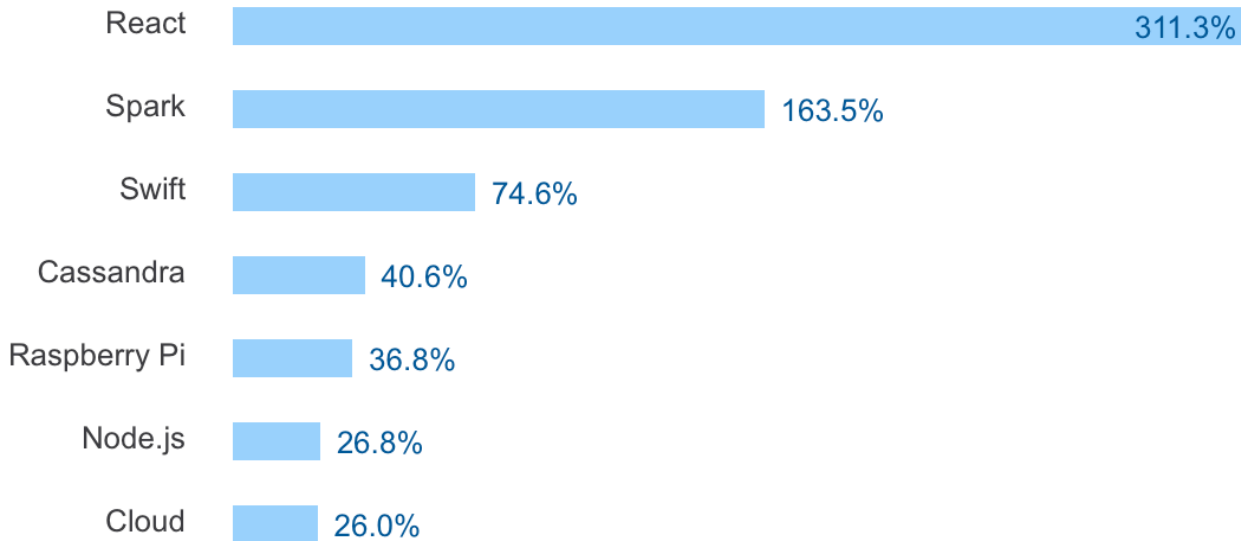
Hossein Falaki
@mhfalaki

databricks™

# About me

- Software Engineer / part-time Data Scientist at Databricks
- I started using Apache Spark since version 0.6
- Developed first version of Apache Spark CSV data source
- Worked on SparkR and R notebooks at Databricks

**databricks**™

# Stackoverflow 2016 trending tech

**Winners** | Losers

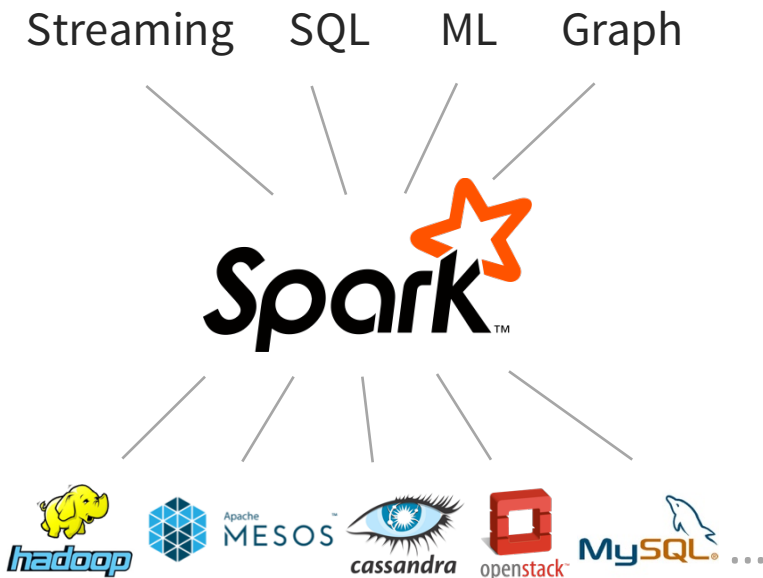| Technology | Growth |
|---|---|
| React | 311.3% |
| Spark | 163.5% |
| Swift | 74.6% |
| Cassandra | 40.6% |
| Raspberry Pi | 36.8% |
| Node.js | 26.8% |
| Cloud | 26.0% |



databricks

3

# Apache Spark Philosophy

**1** Unified engine
Support end-to-end applications

**2** High-level APIs
Easy to use, rich optimizations

**3** Integrate broadly
Storage systems, libraries, etc

Streaming    SQL    ML    Graph

**Spark**™

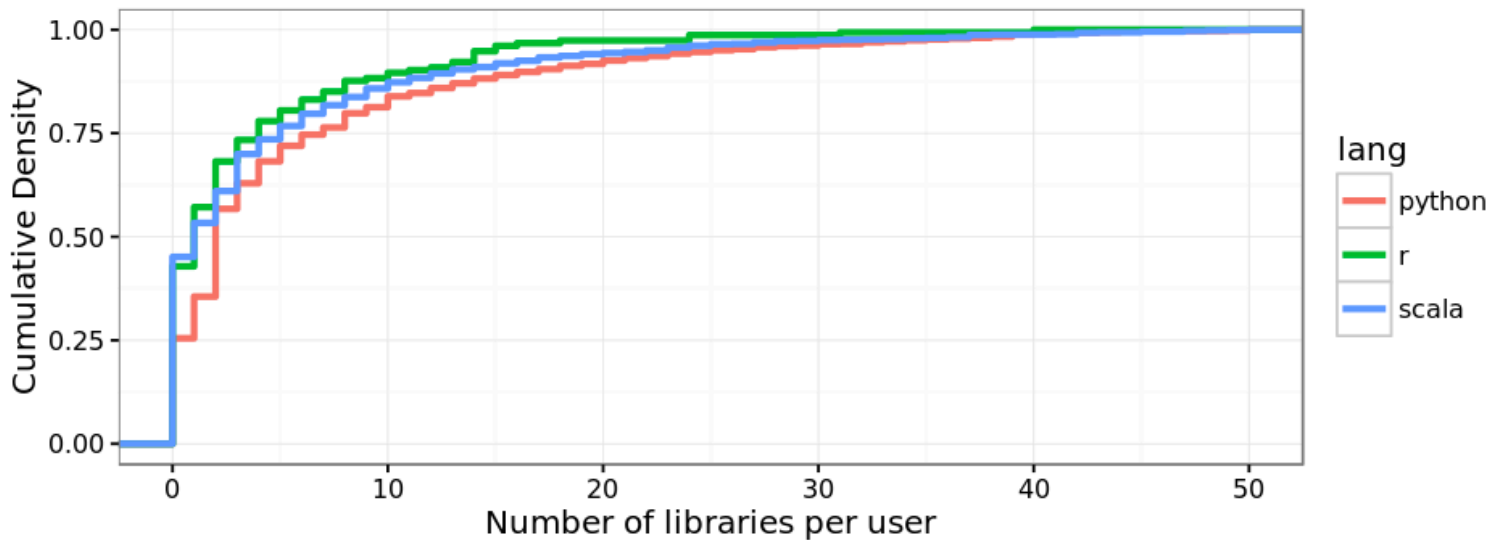hadoop    MESOS    cassandra    openstack    MySQL ...

databricks™

# Databricks Community Edition

- In February Databricks launched a free version of its cloud based platform in beta

- Since then more than 8,000 users registered

- Users created over 61,000 notebooks in different languages

- This is an analysis of third party libraries that our beta users imported to complement Apache Spark in Scala, Python, and R

# What % of users use other libraries

| Language | % users importing external libs | Average # libs | Median # libs |
|----------|--------------------------------|----------------|---------------|
| Python | 75 % | 9 | 2 |
| Scala | 55 % | 3 | 1 |
| R | 57 % | 6 | 1 |

# Installing libraries is easy

## New Library

Language    Upload Python Egg or PyPI   ⇕

## Install PyPi Package

PyPi Name    PyPi Package

Install Library

## Upload Egg
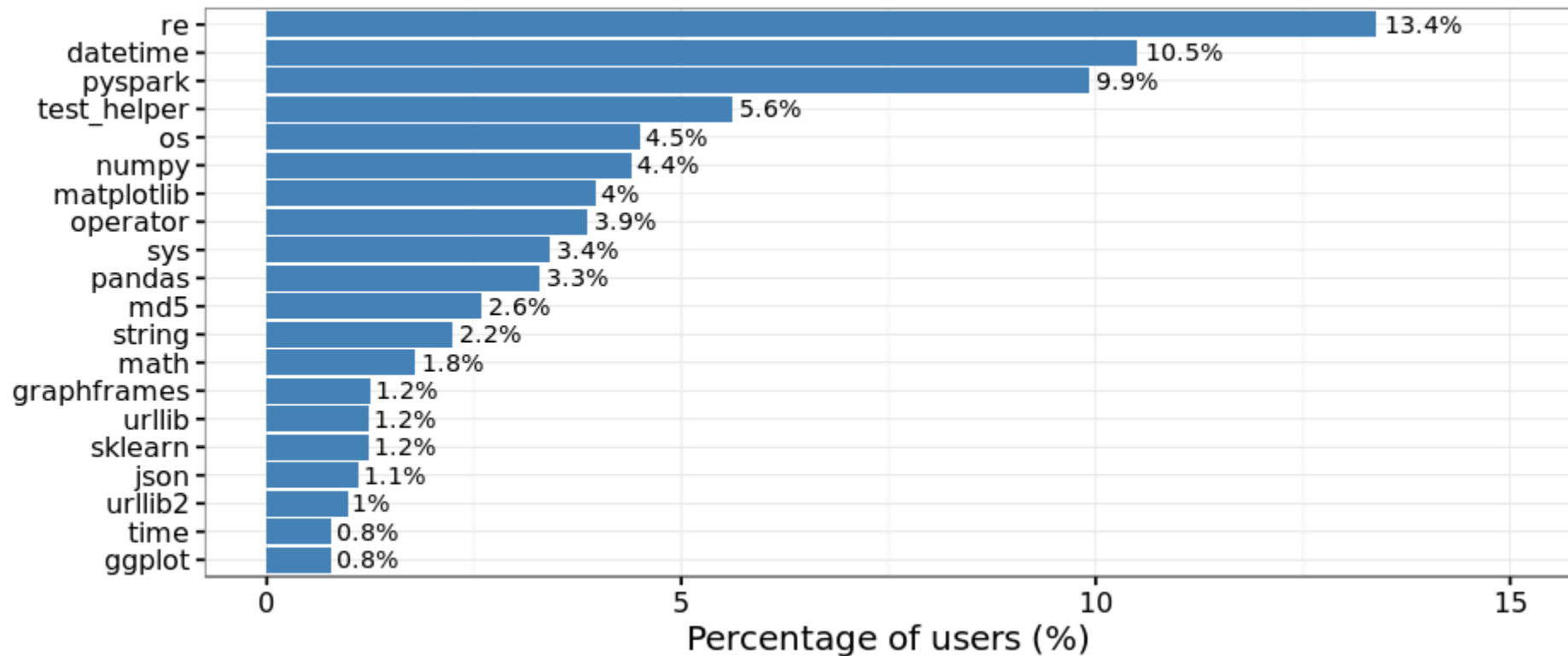
Library Name    Library Name

Egg File    Drop library egg here to upload

# Python Packages

# Most popular Python packages



A horizontal bar chart titled with the x-axis "Percentage of users (%)". The bars and their values from top to bottom:

| Package | Percentage of users (%) |
| --- | --- |
| re | 13.4% |
| datetime | 10.5% |
| pyspark | 9.9% |
| test_helper | 5.6% |
| os | 4.5% |
| numpy | 4.4% |
| matplotlib | 4% |
| operator | 3.9% |
| sys | 3.4% |
| pandas | 3.3% |
| md5 | 2.6% |
| string | 2.2% |
| math | 1.8% |
| graphframes | 1.2% |
| urllib | 1.2% |
| sklearn | 1.2% |
| json | 1.1% |
| urllib2 | 1% |
| time | 0.8% |
| ggplot | 0.8% |

databricks

# What is test_helper?

python™

search

» Package Index > test_helper > 0.2

## test_helper 0.2

*A testing helper for scalable machine learning mooc*

**Download**
test_helper-0.2.tar.gz

**Not Logged In**

Login
Register
Lost Login?
Use OpenID lp
Login with Google G

**Status**

Nothing to report

| File | Type | Py Version | Uploaded on | Size |
|------|------|------------|-------------|------|
| test_helper-0.2.tar.gz (md5) | Source | | 2015-05-11 | 1KB |

**Author:** Daniel Liu
**Home Page:** https://github.com/hpec/test_helper
**Download URL:** https://github.com/hpec/test_helper/tarball/0.1
**Keywords:** testing,autograder,mooc
**Package Index Owner:** hpec1
**DOAP record:** test_helper-0.2.xml

databricks™

# What are these?

**ETL**
- re
- datetime
- pandas
- json
- csv
- string
- math / operator
- urllib / urllib2
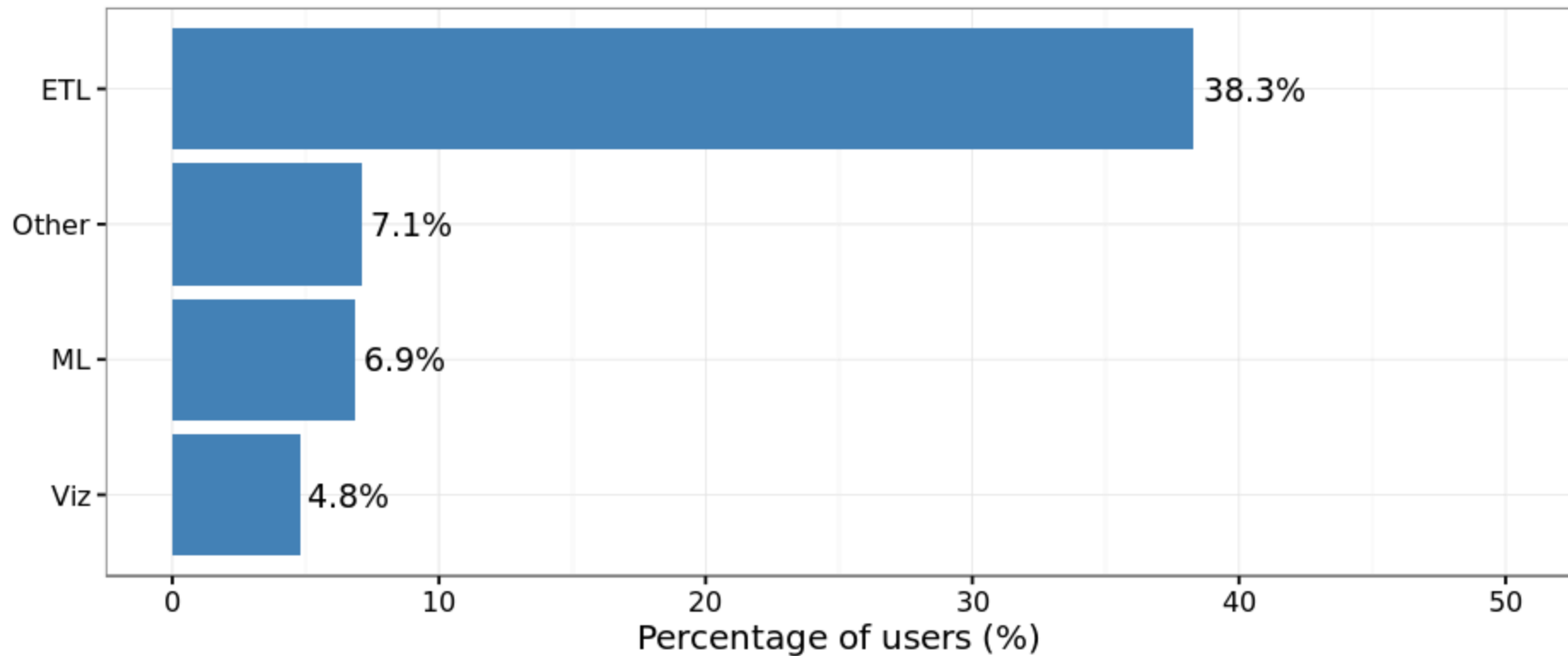
**Visualization**
- matplotlib
- ggplot
- seaborn

**Advanced analytics**
- numpy
- sklearn
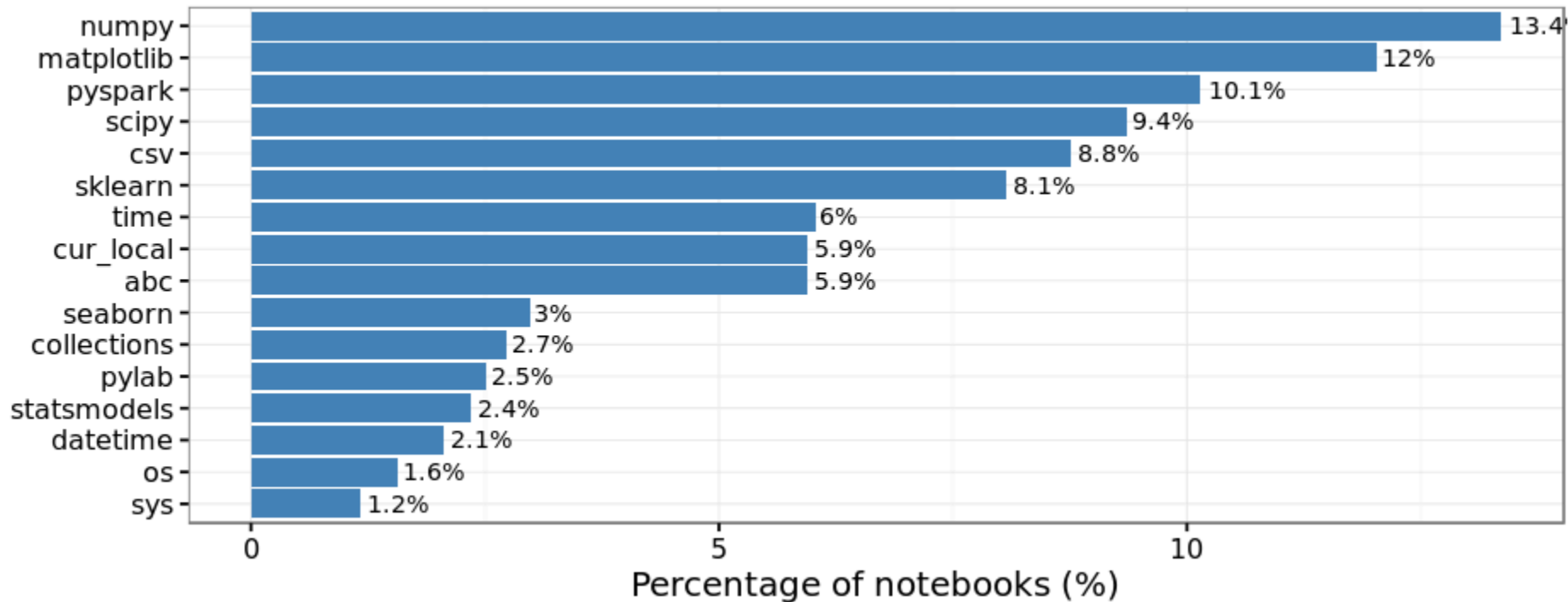- graphframes
- tensorflow
- scipy

**Other**
- test_helper
- os
- md5

# Python package categories

# What packages go together?



## Notebooks that used pandas also used

| Package | Percentage of notebooks (%) |
|---|---|
| numpy | 13.4 |
| matplotlib | 12% |
| pyspark | 10.1% |
| scipy | 9.4% |
| csv | 8.8% |
| sklearn | 8.1% |
| time | 6% |
| cur_local | 5.9% |
| abc | 5.9% |
| seaborn | 3% |
| collections | 2.7% |
| pylab | 2.5% |
| statsmodels | 2.4% |
| datetime | 2.1% |
| os | 1.6% |
| sys | 1.2% |

# Scala Packages

# Most popular Scala libraries

# What are these?

**ETL**
- java/scala util
- scala.collection
- scala.math
- java.{io, nio}
- java.text
- o.a.commons
- kafka
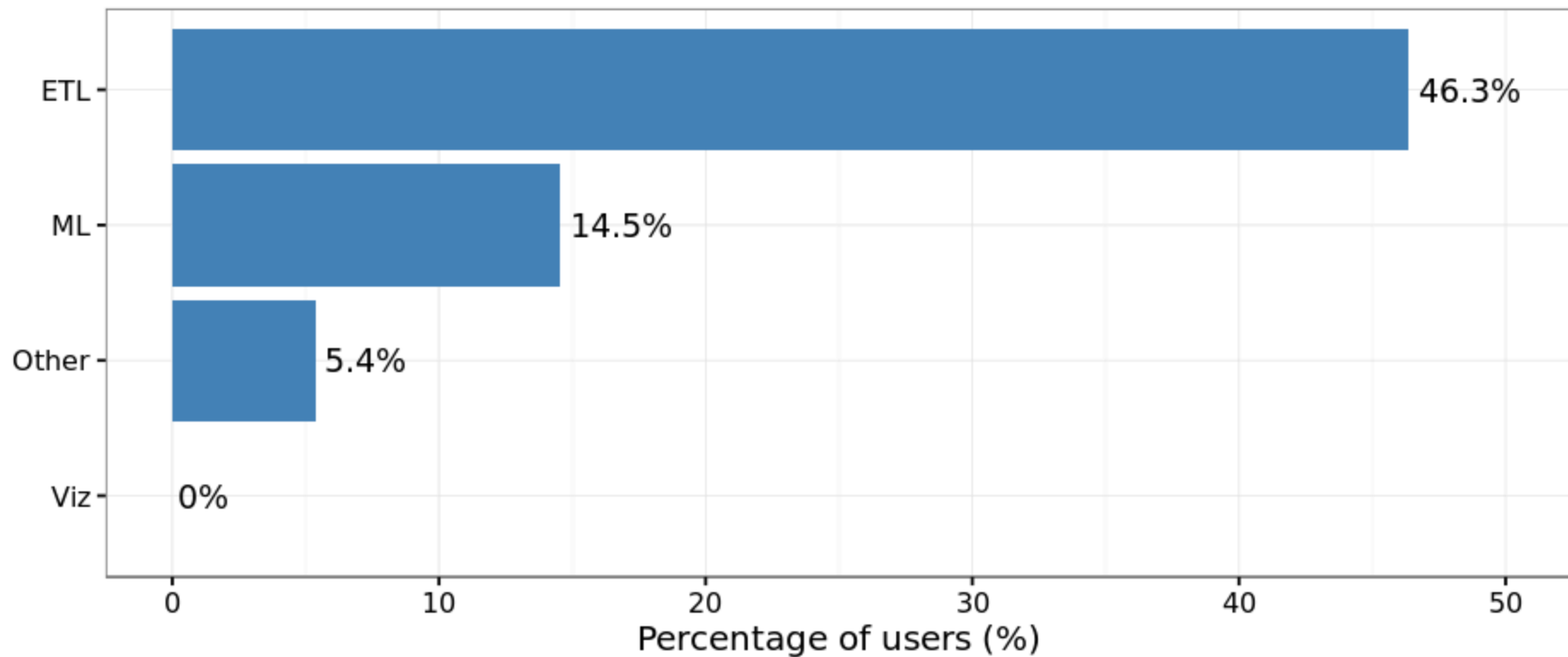- twitter4j

**Visualization**
- ?

**Advanced_analytics**
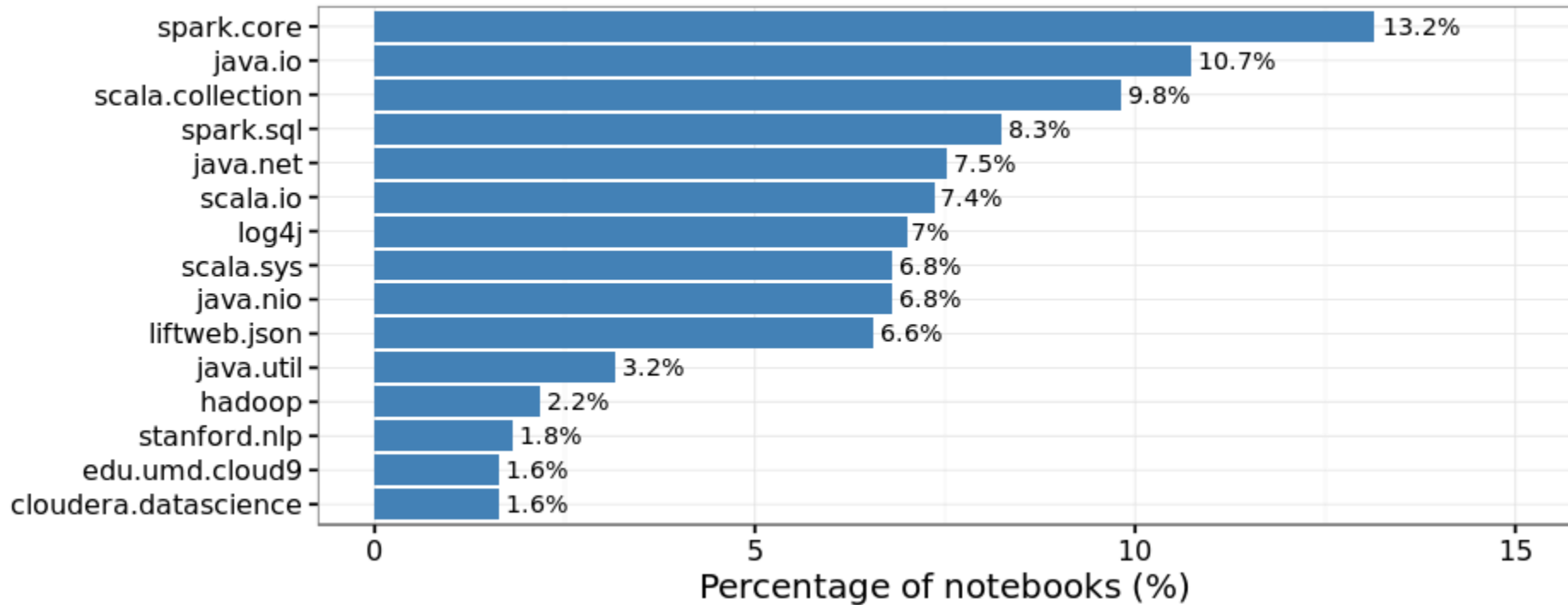- spark.ml
- graphframes

**Other**
- java.net
- scala.sys

databricks

# Scala package categories

# What libraries go together?

## Notebooks that used spark.ml also used

| Library | Percentage |
|---|---|
| spark.core | 13.2% |
| java.io | 10.7% |
| scala.collection | 9.8% |
| spark.sql | 8.3% |
| java.net | 7.5% |
| scala.io | 7.4% |
| log4j | 7% |
| scala.sys | 6.8% |
| java.nio | 6.8% |
| liftweb.json | 6.6% |
| java.util | 3.2% |
| hadoop | 2.2% |
| stanford.nlp | 1.8% |
| edu.umd.cloud9 | 1.6% |
| cloudera.datascience | 1.6% |

Percentage of notebooks (%)

# R Packages

# Most popular R packages



Bar chart — Percentage of users (%):
- ggplot2: 20.6%
- magrittr: 6.3%
- dplyr: 5.6%
- reshape2: 2.4%
- plyr: 2.1%
- beanplot: 2.1%
- jsonlite: 1.4%
- h2o: 1.4%
- tidyr: 1%
- lubridate: 1%
- devtools: 1%
- caret: 1%
- reshape: 0.7%
- plotly: 0.7%
- MASS: 0.7%
- httr: 0.7%
- htmltools: 0.7%
- ggplot: 0.7%
- e1071: 0.7%
- data.table: 0.7%

# What are these?

**ETL**
- dplyr
- plyr
- reshape2
- jsonlite
- tidyr
- lubridate
- httr
- data.table

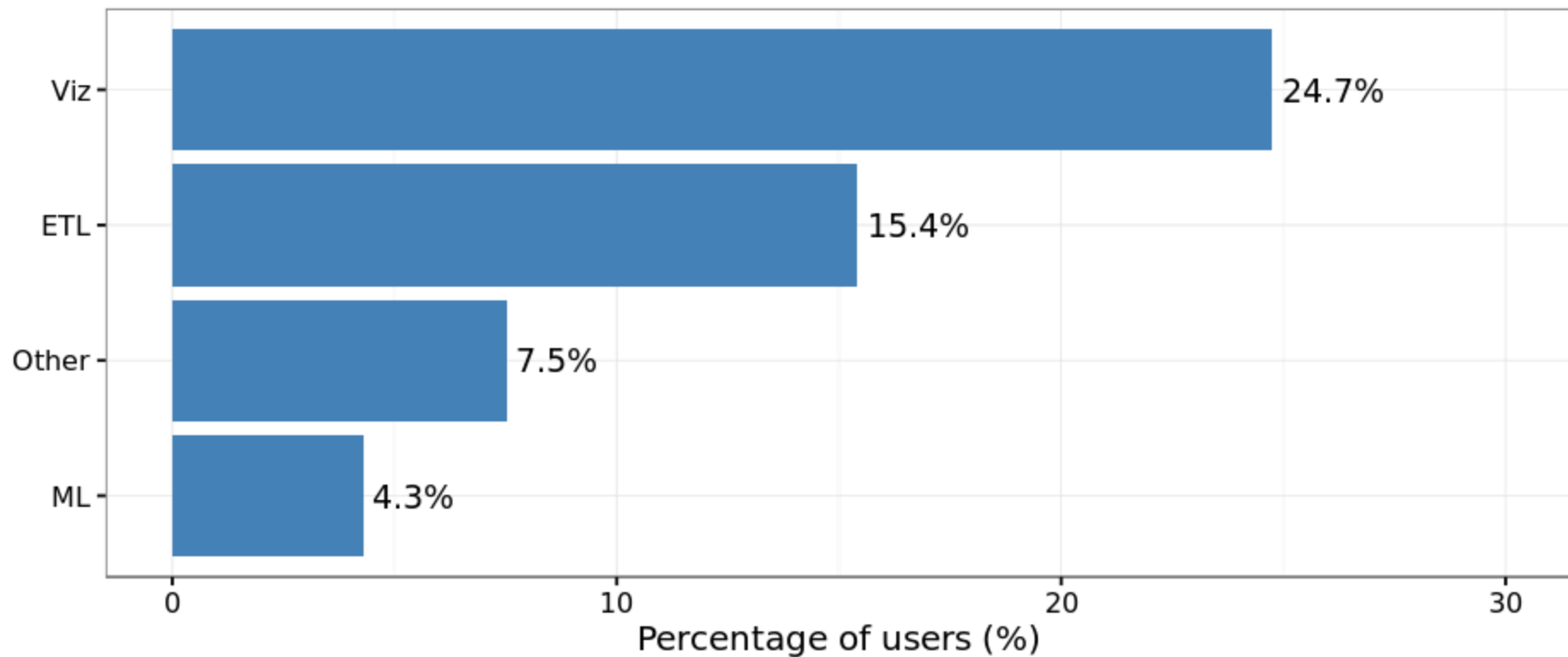**Visualization**
- ggplot2
- beanplot
- plotly
- …

**Advanced analytics**
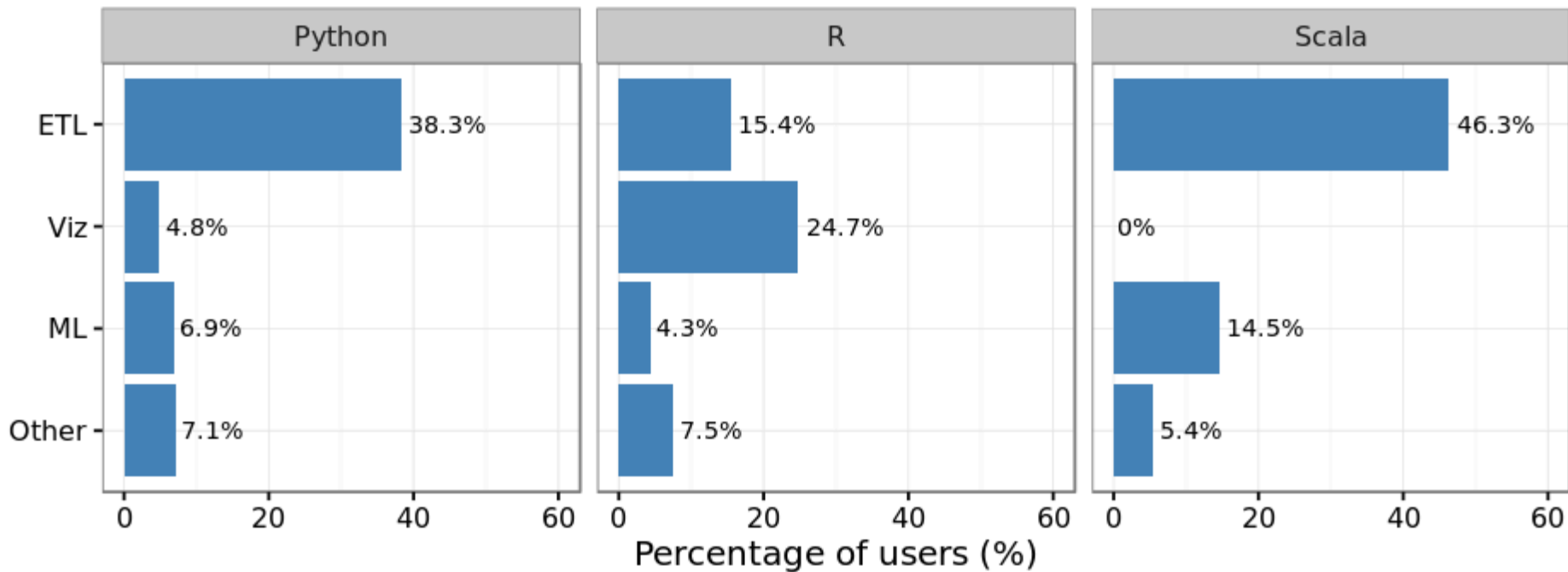- sparkr
- h2o
- caret
- e1071

**Other**
- devtools
- magrittr

# R package categories

# Comparing Python, Scala & R

# Languages have unique features



| 1. ETL | → | 2. Explore | → | 3. Model |

Scala**/** Python / R          R / Python          Scala / Python**/** R

- 25 % of users, use multiple languages
- 3% of notebooks mix different languages

# Summary

- Spark users extensively mix it with other packages in different languages
  - One of goals of Spark project is working well with other projects
- ETL related libraries are the most popular category
  - Opportunities for new data sources
- Notebooks are being used for "small data" as well as "big data."
- Languages and their ecosystems have diverse capabilities. Users seem to be mixing languages to their advantage
  - Scala is missing visualization libraries

# Try your favorite library in Databricks

Try latest version of Apache Spark and preview of Spark 2.0

## http://databricks.com/ce

Create Cluster

**New Cluster** | Cancel | **Create Cluster**

**Cluster Name**

New Cluster

**Spark Version**

✓ Spark 2.0 (apache/branch-2.0 preview)
Spark 1.3.0 (Hadoop 1)
Spark 1.4.1 (Hadoop 1)
Spark 1.5.2 (Hadoop 1)
Spark 1.6.0 (Hadoop 1)
Spark 1.6.1 (Hadoop 1)
Spark 1.6.1 (Hadoop 2)

l automatically t
de your Databri

databricks™

Thank you!

# What packages are used together?



Notebooks that used sparkr also used

| Package | Percentage |
|---|---|
| ggplot2 | 12% |
| dplyr | 11.4% |
| plyr | 10.2% |
| rpart.plot | 9.9% |
| rpart | 9.9% |
| randomforest | 9.9% |
| metrics | 9.9% |
| e1071 | 9.9% |
| caret | 9.9% |
| magrittr | 1.5% |

Percentage of notebooks (%)