DATA+AI
SUMMIT EUROPE
FORMERLY SPARK+AI SUMMIT

# Project Zen: Improving Apache Spark for Python Users

Hyukjin Kwon

Databricks Software Engineer

#DataTeams #DataAISummit

# Hyukjin Kwon



- Apache Spark Committer / PMC

- Major Koalas contributor

- Databricks Software Engineer

- @HyukjinKwon in GitHub

# Agenda

What is Project Zen?

Redesigned Documentation

PySpark Type Hints

Distribution Option for PyPI Users
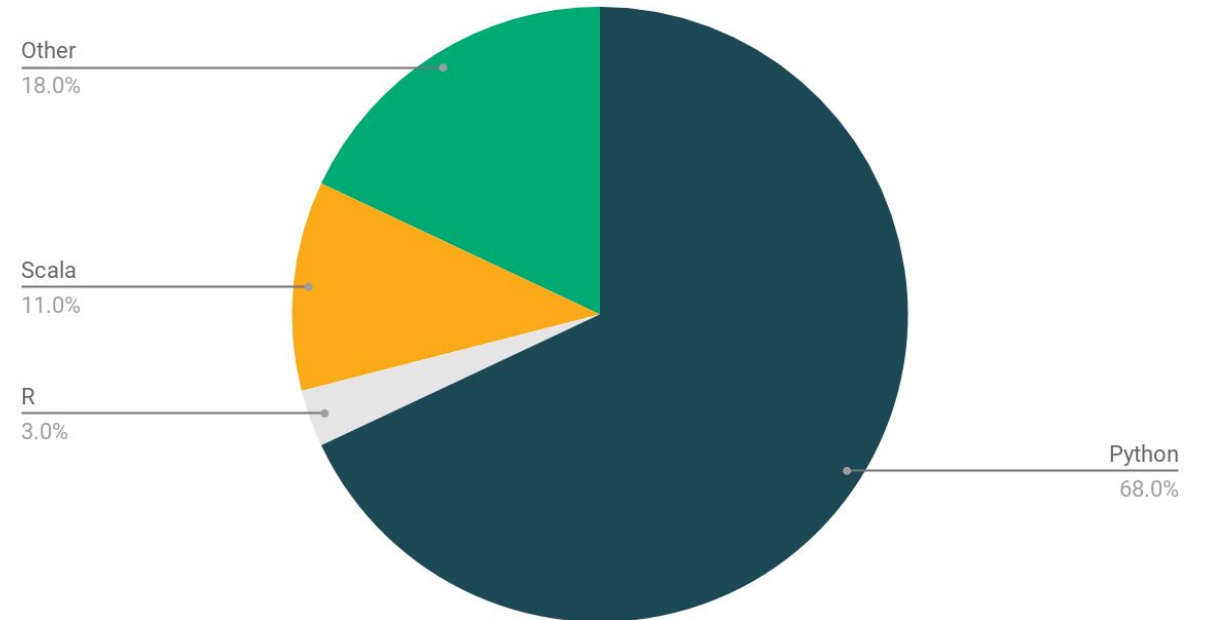
Roadmap

# What is Project Zen?

# Python Growth

## 68%

of notebook commands on Databricks are in Python

### Language Use in Notebooks



- Other 18.0%
- Scala 11.0%
- R 3.0%
- Python 68.0%

# PySpark Today

- ## Documentation difficult to navigate

  - All APIs under each module is listed in single page

  - No other information or classification

- ## Lack of information

  - No Quickstart page

  - No Installation page

  - No Introduction
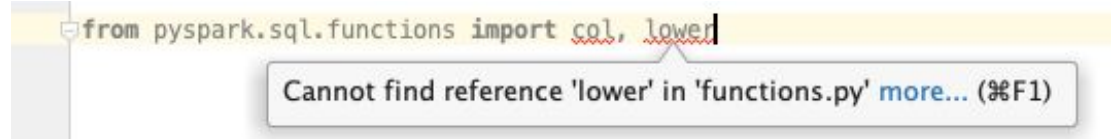


**PySpark documentation**

# PySpark Today

- ## IDE unfriendly
  - Dynamically defined functions reported as missing functions
  - Lack of autocompletion support
  - Lack of type checking support

- ## Notebook unfriendly
  - Lack of autocompletion support



```
from pyspark.sql.functions import col, lower
```

Cannot find reference 'lower' in 'functions.py' more... (⌘F1)

**Missing import in IDE**

```
rdd = spark.sparkContext.
```

| main | if __name__ == '__main__': exp |
| not | not expr |
| par | (expr) |

Press ^. to choose the selected (or first) suggestion and insert a dot afterwards ≥ π

**Autocompletion in IDE**

```
In [1]: rdd = spark.sparkContext.
```

**Autocompletion in Jupyter**

# PySpark Today

- ## Less Pythonic

  - Deprecated Python built-in instance support when creating a DataFrame

```
>>> spark.createDataFrame([{'a': 1}])
```

```
/.../session.py:378: UserWarning: inferring schema from dict is deprecated, please use
pyspark.sql.Row instead
```

# PySpark Today

- ## Missing distributions with other Hadoop versions in PyPI

  - Missing Hadoop 3 distribution

  - Missing Hive 1.2 distribution



| Filename, size | File type | Python version | Upload date | Hashes |
|---|---|---|---|---|
| pyspark-3.0.1.tar.gz (204.2 MB) | Source | None | Sep 8, 2020 | View |

spark-3.0.1-bin-hadoop2.7-hive1.2.tgz  2020-08-28 18:25  209M
spark-3.0.1-bin-hadoop2.7.tgz          2020-08-28 18:25  210M
spark-3.0.1-bin-hadoop3.2.tgz          2020-08-28 18:25  214M
spark-3.0.1-bin-without-hadoop.tgz     2020-08-28 18:25  149M

**Apache Mirror**

**PyPI Distribution**

# PySpark Today

- ## Inconsistent exceptions and warnings

  - Unclassified exceptions and warnings

```
>>> spark.range(10).explain(1, 2)
```

```
Traceback (most recent call last):
    ...
Exception: extended and mode should not be set together.
```

# The Zen of Python

## The Zen of Python

```
Beautiful is better than ugly.

Explicit is better than implicit.

Simple is better than complex.

Complex is better than complicated.

Flat is better than nested.

Sparse is better than dense.

Readability counts.

Special cases aren't special enough to break the rules.
```

**PEP 20 — The Zen of Python**

# Project Zen (SPARK-32082)

- ## Be Pythonic

  - The Zen of Python

  - Python friendly

- ## Better and easier use of PySpark

  - Better documentation

  - Clear exceptions and warnings

  - Python type hints: autocompletion, static type checking and error detection

  - More options for `pip` installation

- ## Better interoperability with other Python libraries

  - pandas, pyarrow, NumPy, Koalas, etc.

  - Visualization

# Redesigned Documentation

# Problems in PySpark Documentation

- ## Everything in few pages

  - Whole module in single page w/o classification

- ## Difficult to navigate

  - Very long to stroll down

  - Virtually no structure

- ## No other useful pages

  - How to start?

  - How to ship 3rd party packages together?

  - How to install?

  - How to debug / setup an IDE?



**(Old) PySpark documentation**

# New PySpark Documentation



**New PySpark documentation**

# New PySpark Documentation



**New user guide page**

# New PySpark Documentation



Top menu

Getting Started   **User Guide**   API Reference   Development   Migration Guide

Search

Sub-titles in the current page

### Series to Scalar

The type hint can be expressed as `pandas.Series`, ... -> Any.

By using `pandas_udf()` with the function having such type hints above, it creates a Pandas UDF similar to PySpark's aggregate functions. The given function takes *pandas.Series* and returns a scalar value. The return type should be a primitive data type, and the returned scalar can be either a python primitive type, e.g., `int` or `float` or a numpy data type, e.g., `numpy.int64` or `numpy.float64`. `Any` should ideally be a specific scalar type accordingly.

This UDF can be also used with `GroupedData.agg()` and *Window*. It defines an aggregation from one or more *pandas.Series* to a scalar value, where each *pandas.Series* represents a column within the group or window.

Note that this type of UDF does not support partial aggregation and all data for a group or window will be loaded into memory. Also, only unbounded window is supported with Grouped aggregate Pandas UDFs currently. The following example shows how to use this type of UDF to compute mean with a group-by and window operations:

```
import pandas as pd
```

**Apache Arrow in PySpark**

3rd Party Python Packages

On this page

Ensure PyArrow Installed

Enabling for Conversion to/from Pandas

**Pandas UDFs (a.k.a. Vectorized UDFs)**

Series to Series

Iterator of Series to Iterator of Series

Iterator of Multiple Series to Iterator of Series

**Series to Scalar**

Pandas Function APIs

Usage Notes

Other pages

Contents in the current page

**DATA+AI** SUMMIT EUROPE

**#DataTeams  #DataAISummit**

# New API Reference



**New API reference page**

# New API Reference



**New API reference page**

# New API Reference



**New API reference page**

# New API Reference



**New API reference page**

Table for each classification

# New API Reference



**New API reference page**

# New API Reference



**New API reference page**

Each page for each API

# Quickstart



**Quickstart page**

# Live Notebook



**Move to live notebook (Binder integration)**

# Live Notebook



**Live notebook (Binder integration)**

# Other New Pages

## Development

- Contributing to PySpark
  - Contributing by Testing Releases
  - Contributing Documentation Changes
  - Preparing to Contribute Code Changes
  - Code Style Guide
- Testing PySpark
  - Running Individual PySpark Tests
  - Running tests using GitHub Actions
- Debugging PySpark
  - Remote Debugging (PyCharm Professional)
  - Checking Resource Usage (`top` and `ps`)
  - Profiling Memory Usage (Memory Profiler)
  - Identifying Hot Loops (Python Profilers)
- Setting up IDEs
  - PyCharm

## Getting Started

This page summarizes the basic steps required

- Installation
  - Python Version Supported
  - Using PyPI
  - Using Conda
  - Manually Downloading
  - Installing from Source
  - Dependencies
- Quickstart
  - DataFrame Creation
  - Viewing Data
  - Selecting and Accessing Data
  - Applying a Function
  - Grouping Data
  - Getting Data in/out
  - Working with SQL

## User Guide

- Apache Arrow in PySpark
  - Ensure PyArrow Installed
  - Enabling for Conversion to/from Pandas
  - Pandas UDFs (a.k.a. Vectorized UDFs)
  - Pandas Function APIs
  - Usage Notes
- 3rd Party Python Packages
  - Using PySpark Native Features
  - Using Zipped Virtual Environment
  - Using PEX

## Migration Guide

This page describes the migration guide specific to PySpark. Many items of applied when migrating PySpark to higher versions because PySpark interna Please also refer other migration guides such as Migration Guide: SQL, Data

- Upgrading from PySpark 2.4 to 3.0
- Upgrading from PySpark 2.3 to 2.4
- Upgrading from PySpark 2.3.0 to 2.3.1 and above
- Upgrading from PySpark 2.2 to 2.3
- Upgrading from PySpark 1.4 to 1.5
- Upgrading from PySpark 1.0-1.2 to 1.3

**New useful pages**

# PySpark Type Hints

# What are Python Type Hints?

```python
def greeting(name):
    return 'Hello ' + name
```

**Typical Python codes**

```python
def greeting(name: str) -> str:
    return 'Hello ' + name
```

**Python codes with type hints**

```python
def greeting(name: str) -> str: ...
```

**Stub syntax (.pyi file)**

# Why are Python Type Hints good?

- IDE Support

- Notebook Support

- Documentation

- Static error detection



**Before type hints**



**After type hints**

# Why are Python Type Hints good?

- IDE Support

- Notebook Support

- Documentation

- Static error detection

```
In [1]: rdd = spark.sparkContext.
```

**Before type hints**

```
In [1]: rdd = spark.sparkContext.
                accumulator
                addFile
                addPyFile
                applicationId
                appName
                binaryFiles
                binaryRecords
                broadcast
                cancelAllJobs
                cancelJobGroup
```

**After type hints**

# Why are Python Type Hints good?

- IDE Support

pyspark.sql.functions.corr(*col1*, *col2*)                                                    [source]

**Before type hints**

- Notebook Support

- Documentation

pyspark.sql.functions.corr(*col1: pyspark.sql.column.Column, col2: pyspark.sql.column.Column*) →
**pyspark.sql.column.Column**                                                                  [source]

**After type hints**

- Static error detection

# Why are Python Type Hints good?

- IDE Support

- Notebook Support

- Documentation

- Static error detection



**Static error detection**

https://github.com/zero323/pyspark-stubs#motivation

# Python Type Hints in PySpark



[SPARK-32714][PYTHON] Initial pyspark-stubs port. #29591

## Built-in in the upcoming Apache Spark 3.1!

Community support: zero323/pyspark-stubs

User facing APIs only

Stub (.pyi) files

# Installation Option for PyPI Users

# PyPI Distribution



**PySpark on PyPI**

# PyPI Distribution

- ## Multiple distributions available

    - Hadoop 2.7 and Hive 1.2

    - Hadoop 2.7 and Hive 2.3

    - Hadoop 3.2 and Hive 2.3

    - Hive 2.3 without Hadoop



| | |
|---|---|
| spark-3.0.1-bin-hadoop2.7-hive1.2.tgz | 2020-08-28 18:25 209M |
| spark-3.0.1-bin-hadoop2.7.tgz | 2020-08-28 18:25 210M |
| spark-3.0.1-bin-hadoop3.2.tgz | 2020-08-28 18:25 214M |
| spark-3.0.1-bin-without-hadoop.tgz | 2020-08-28 18:25 149M |

**Multiple distributions
in Apache Mirror**

- ## PySpark distribution in PyPI

    - Hadoop 2.7 and Hive 1.3

| Filename, size | File type | Python version | Upload date | Hashes |
|---|---|---|---|---|
| pyspark-3.0.1.tar.gz (204.2 MB) | Source | None | Sep 8, 2020 | View |

**One distribution in PyPI**

# New Installation Options

```
HADOOP_VERSION=3.2 pip install pyspark
```

**Spark with Hadoop 3.2**

```
HADOOP_VERSION=2.7 pip install pyspark
```

**Spark with Hadoop 2.7**

```
HADOOP_VERSION=without pip install pyspark
```

**Spark without Hadoop**

```
PYSPARK_RELEASE_MIRROR=http://mirror.apache-kr.org HADOOP_VERSION=2.7 pip install
```

**Spark downloading from the
specified mirror**

# Why not `pip --install-options`?



**Ongoing issues in pip**

# Roadmap

# Roadmap

- ## Migrate to NumPy documentation style

  - Better classification

  - Better readability

  - Widely used

```
"""Specifies some hint on the current
:class:`DataFrame`.

:param name: A name of the hint.
:param parameters: Optional parameters.
:return: :class:`DataFrame`
```

**reST style**

```
"""Specifies some hint on the
current :class:`DataFrame`.

Parameters
----------
name : str
    A name of the hint.
parameters : dict, optional
    Optional parameters

Returns
-------
DataFrame
```

**Numpydoc style**

# Roadmap

- ## Standardize warnings and exceptions

  - Classify the exception and warning types

  - Python friendly messages instead of JVM stack trace

```
>>> spark.range(10).explain(1, 2)
```
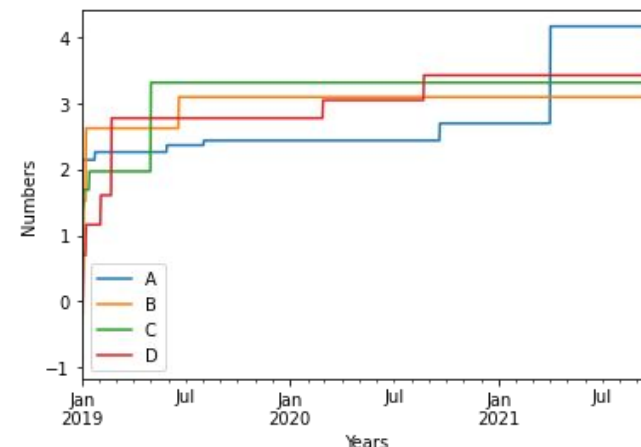
```
Traceback (most recent call last):
    ...
Exception: extended and mode should not be set together.
```

**Plain `Exception` being thrown**

# Roadmap

- **Interoperability between NumPy, Koalas, other libraries**

  - Common features in DataFrames

  - NumPy universe function

- **Visualization and plotting**

  - Make a chart from Spark DataFrame

# Re-cap

# Re-cap

- Python and PySpark are becoming more and more popular

- PySpark documentation is redesigned with many new pages

- Auto-completion and type checking in IDE and notebooks

- PySpark download options in PyPI

# Re-cap: What's next?

- Migrate to NumPy documentation style

- Standardize warnings and exceptions

- Visualization

- Interoperability between NumPy, Koalas, other libraries

# Question?