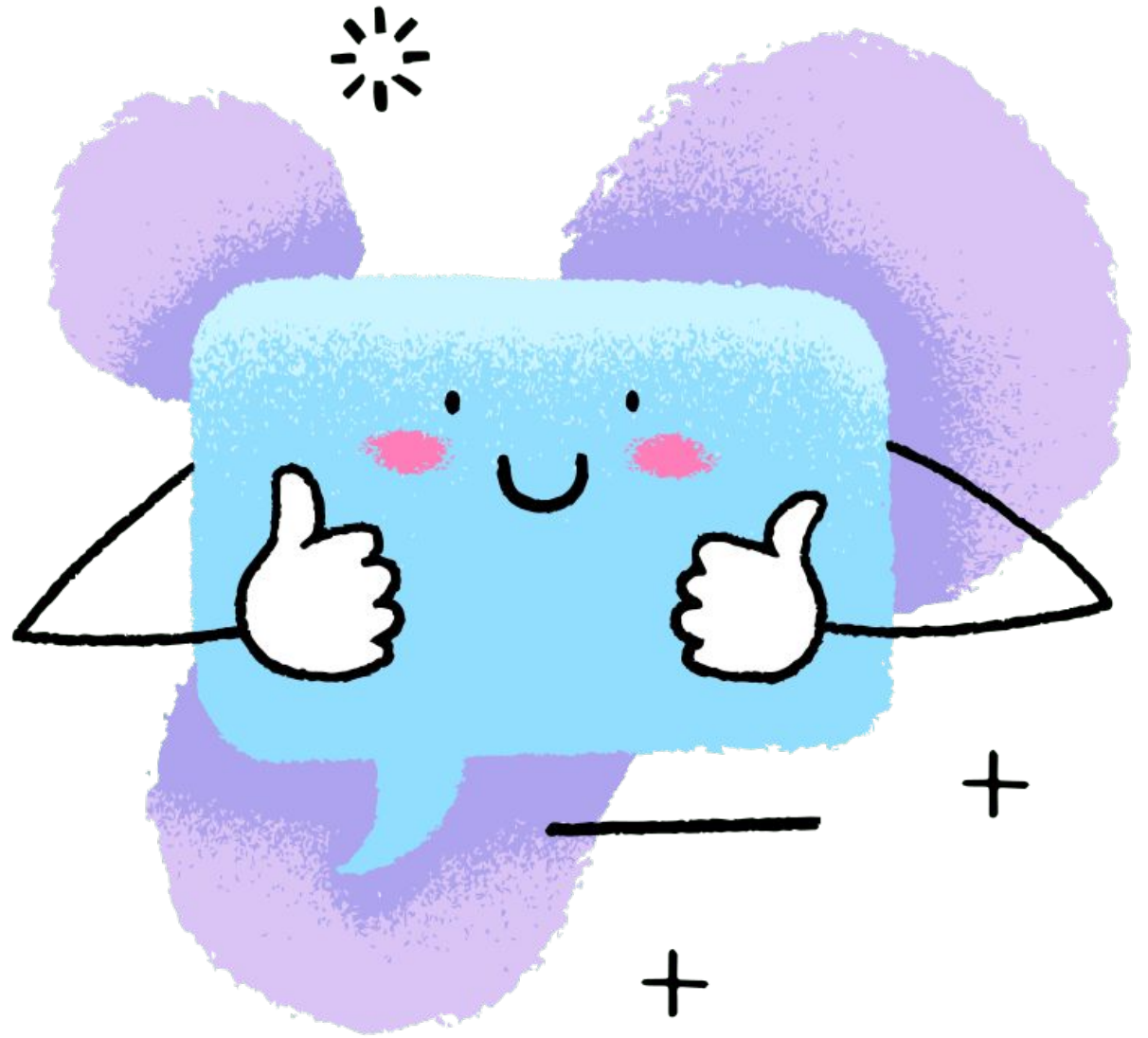


# PySpark



# API. Пользовательские функции

Apache Spark

# Что будет на уроке

1. Дополнение по API
2. Составные типы данных: Array, Map
3. User Defined Functions. Использование UDF в spark.sql
4. ML pipelines

# Union

## Нужно следить за порядком колонок

**union**(*other*)

[\[source\]](#)

Return a new **DataFrame** containing union of rows in this and another frame.

This is equivalent to *UNION ALL* in SQL. To do a SQL-style set union (that does deduplication of elements), use this function followed by a `distinct`.

Also as standard in SQL, this function resolves columns by position (not by name).

*New in version 2.0.*

**unionAll**(*other*)

[\[source\]](#)

Return a new **DataFrame** containing union of rows in this and another frame.

This is equivalent to *UNION ALL* in SQL. To do a SQL-style set union (that does deduplication of elements), use this function followed by a `distinct`.

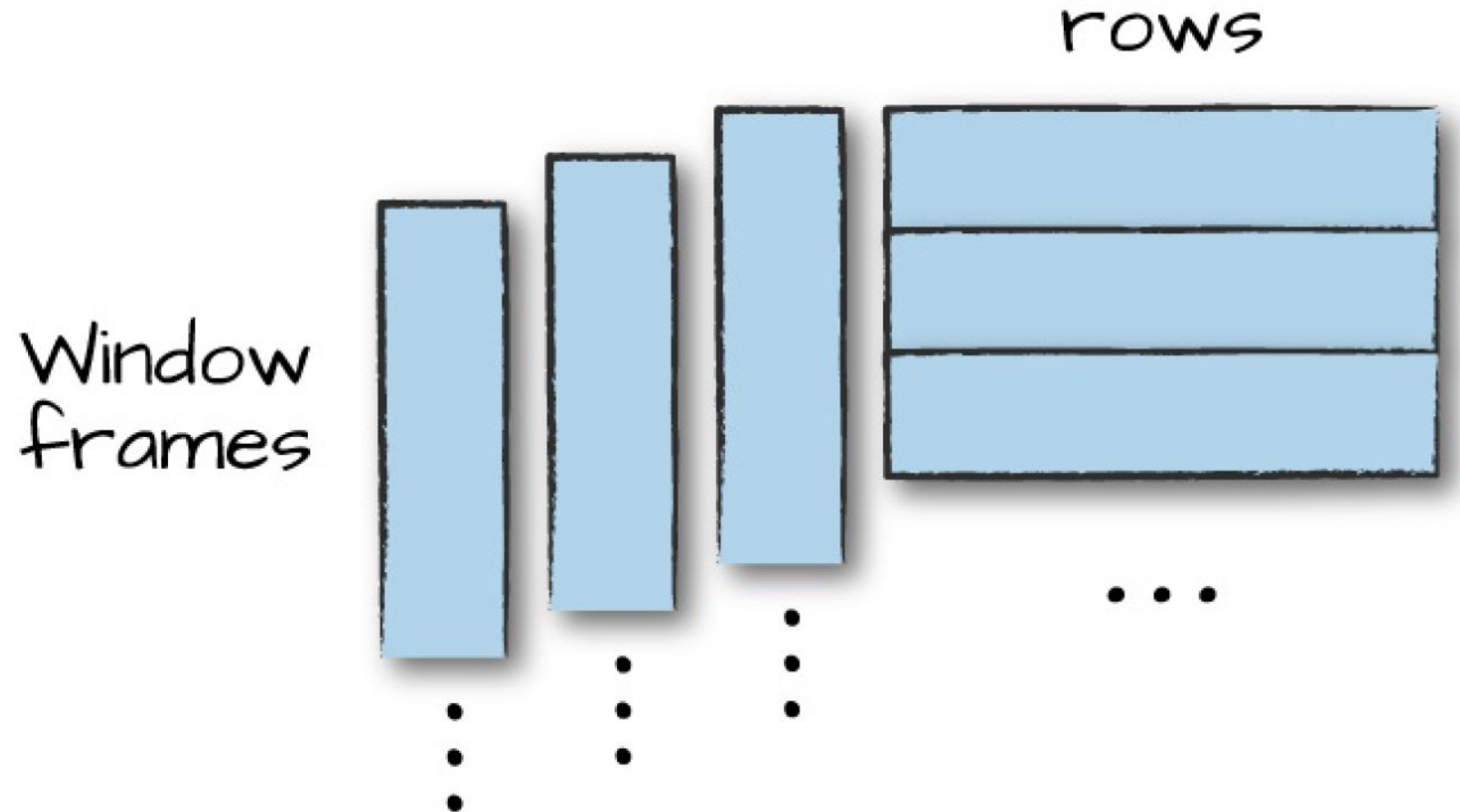
Also as standard in SQL, this function resolves columns by position (not by name).

# explode & flatten (spark 2.4)

## Массив значений в строки

```
arrayArrayData = [  
    ("James",[["Java", "Scala", "C++"], ["Spark", "Java"]]),  
    ("Michael",[["Spark", "Java", "C++"], ["Spark", "Java"]]),  
    ("Robert",[["CSharp", "VB"], ["Spark", "Python"]])  
]  
  
df = spark.createDataFrame(data=arrayArrayData, schema = ['name', 'subjects'])  
  
from pyspark.sql.functions import explode  
df.select(df.name, explode(df.subjects)).show(truncate=False)  
  
from pyspark.sql.functions import flatten  
df.select(df.name, flatten(df.subjects)).show(truncate=False)
```

# Window functions



# Window functions

WINDOW FUNCTIONS USAGE & SYNTAX	PYSPARK WINDOW FUNCTIONS DESCRIPTION
<code>row_number(): Column</code>	Returns a sequential number starting from 1 within a window partition
<code>rank(): Column</code>	Returns the rank of rows within a window partition, with gaps.
<code>percent_rank(): Column</code>	Returns the percentile rank of rows within a window partition.
<code>dense_rank(): Column</code>	Returns the rank of rows within a window partition without any gaps. Where as Rank() returns rank with gaps.
<code>ntile(n: Int): Column</code>	Returns the ntile id in a window partition
<code>cume_dist(): Column</code>	Returns the cumulative distribution of values within a window partition
<code>lag(e: Column, offset: Int): Column</code> <code>lag(columnName: String, offset: Int): Column</code> <code>lag(columnName: String, offset: Int, defaultValue: Any): Column</code>	returns the value that is `offset` rows before the current row, and `null` if there is less than `offset` rows before the current row.
<code>lead(columnName: String, offset: Int): Column</code> <code>lead(columnName: String, offset: Int): Column</code> <code>lead(columnName: String, offset: Int, defaultValue: Any): Column</code>	returns the value that is `offset` rows after the current row, and `null` if there is less than `offset` rows after the current row.



# Booleans

## Фильтрация по значению в колонке

```
1 # in Python
2 from pyspark.sql.functions import col
3 df.where(col("InvoiceNo") != 536365)\
4 .select("InvoiceNo", "Description")\
5 .show(5, False)|
```

```
1 from pyspark.sql.functions import instr
2
3 priceFilter = col("UnitPrice") > 600
4 descripFilter = instr(df.Description, "POSTAGE") >= 1
5
6 df.where(df.StockCode.isin("DOT"))\
7 .where(priceFilter | descripFilter).show()
```



# Регулярные выражения

```
1 from pyspark.sql.functions import regexp_replace
2 regex_string = "BLACK|WHITE|RED|GREEN|BLUE"
3 df.select(
4     regexp_replace(col("Description"), regex_string, "COLOR")
5     .alias("color_clean"),
6     col("Description")
7 ).show(2)
```

# Nulls

первое не-null значение из списка столбцов

```
1 from pyspark.sql.functions import coalesce
2
3 df.select(coalesce(col("Description"), col("CustomerId")))\
4     .show()
```

удаление строк, содержащих null

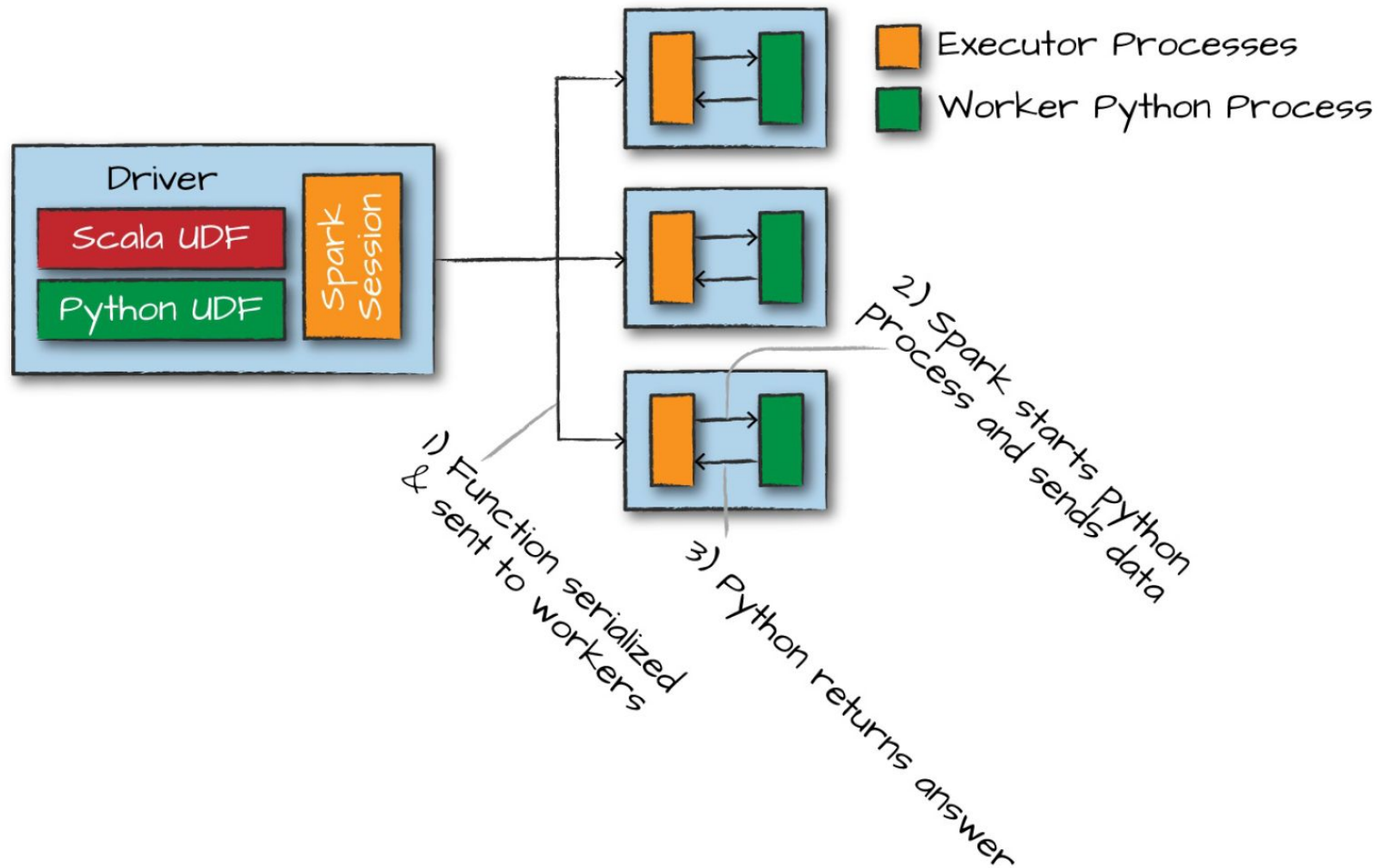
```
1 df.na.drop()
2 df.na.drop("any") # drop if ANY column is null
3
4 df.na.drop("all") # drop if ALL columns are null
5
6 df.na.drop("all", subset=["StockCode", "InvoiceNo"])|
```

# Nulls

## заполнение пустых значений

```
1 df.na.fill("it was null value")
2
3 # specify values with dict
4 fill_cols_vals = {"StockCode": 5, "Description" : "No Value"}
5 df.na.fill(fill_cols_vals)
```

# User-Defined Functions



# Performance concerns with UDFs

- UDFs are black-box to Spark optimizations.
- UDFs block many spark optimizations like
  - WholeStageCodegen
  - Null Optimizations
  - Predicate Pushdown
  - More optimizations from Catalyst Optimizer
- **String Handling within UDFs**
  - UTF-8 to UTF-16 conversion. Spark maintains string in UTF-8 encoding versus Java runtime encodes in UTF-16.
  - Any String input to UDF requires UTF-8 to UTF-16 conversion.
  - Conversely, a String output requires a UTF-16 to UTF-8 conversion.

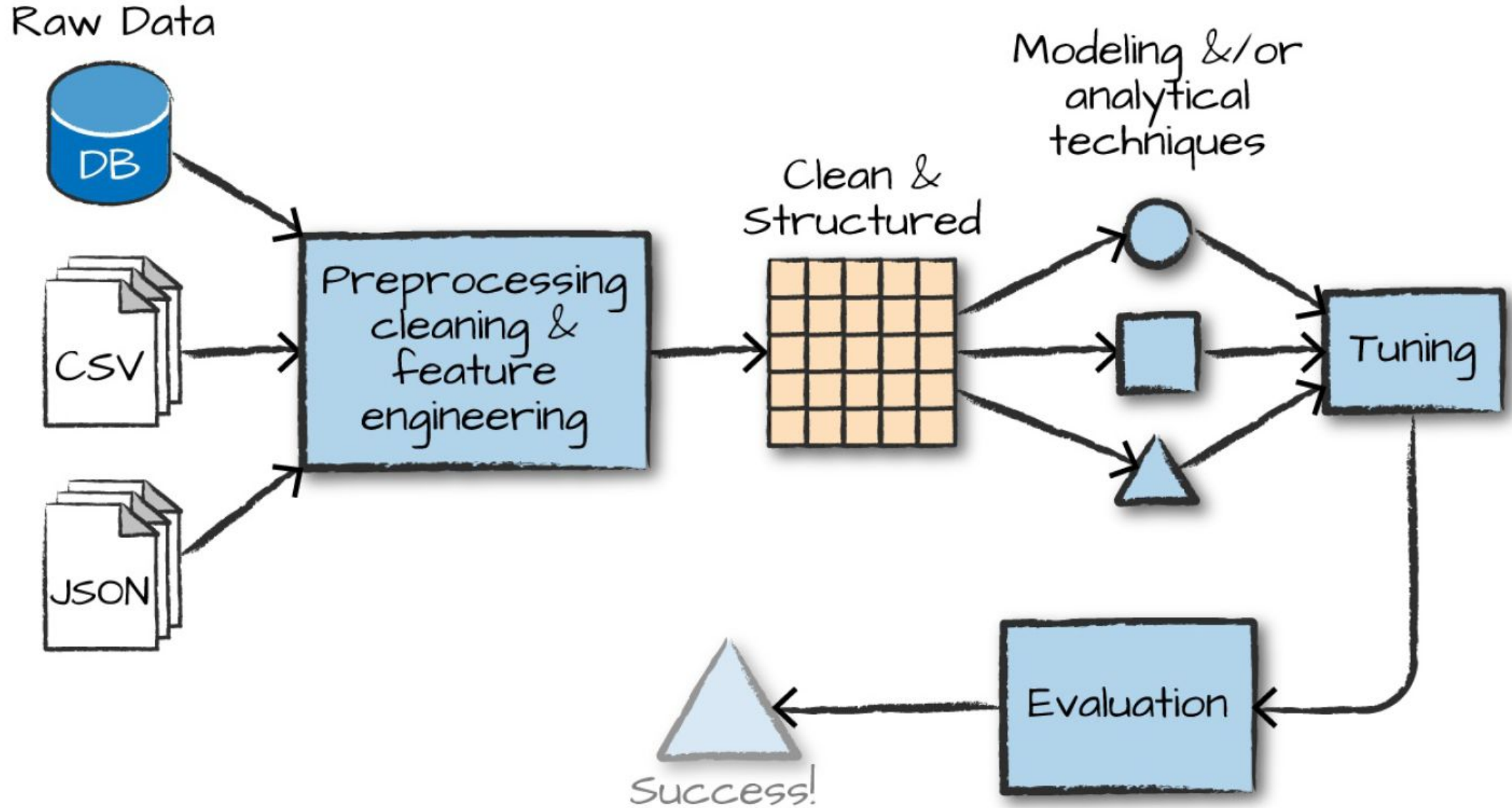
# Ссылки

[Introducing Window Functions in Spark SQL - The Databricks Blog](https://www.youtube.com/watch?v=2QNk-bcjN-I)  
<https://www.youtube.com/watch?v=2QNk-bcjN-I>



# Spark ML

# The machine learning workflow



# The machine learning workflow in Spark

Structured  
APIs

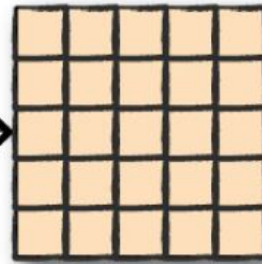
Raw Data



Preprocessing  
cleaning &  
feature  
engineering

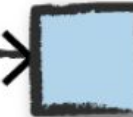
Transformers  
& Estimators

Clean &  
Structured



Estimators  
& Models

Modeling &/or  
analytical  
techniques



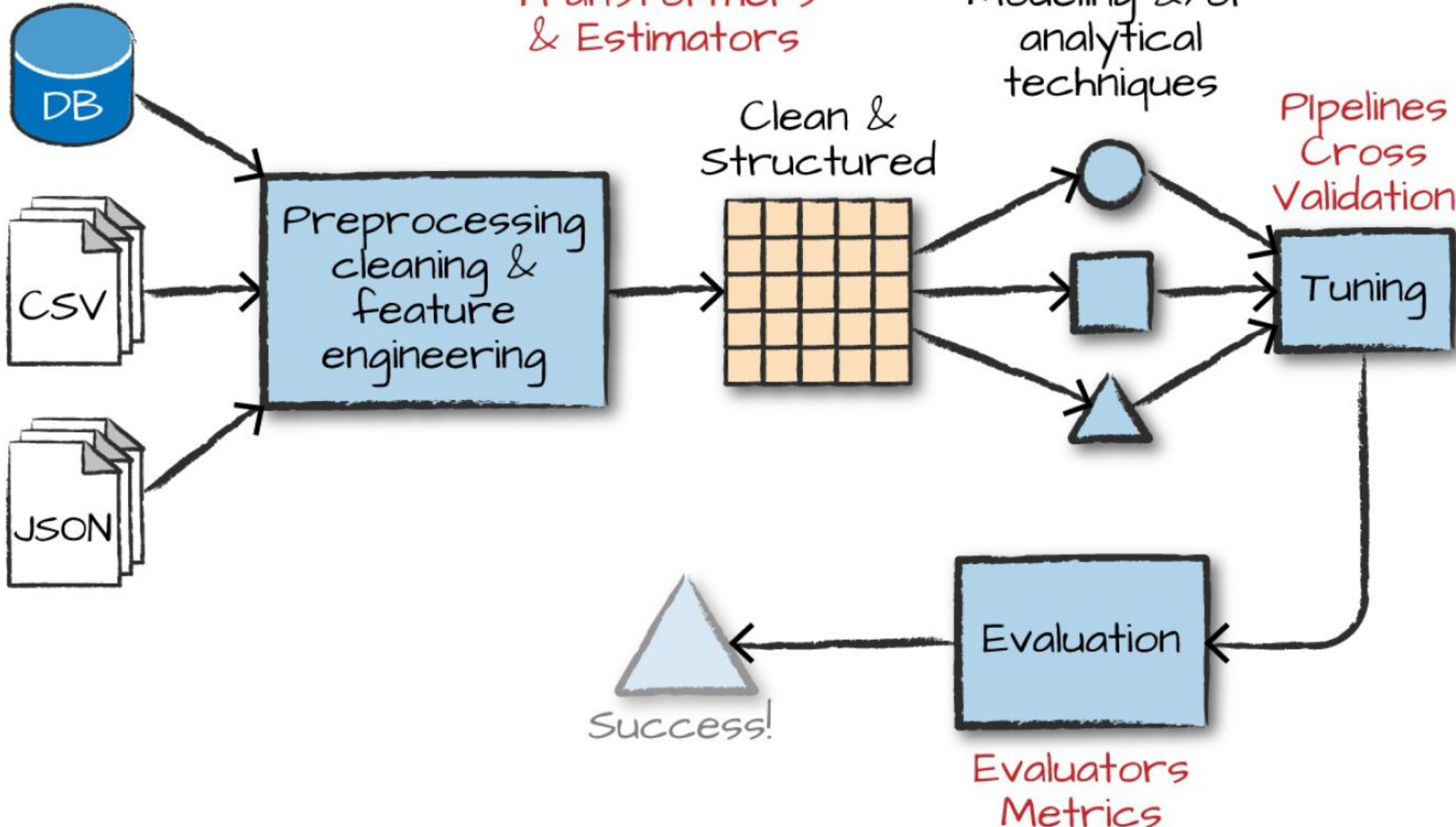
Pipelines  
Cross  
Validation

Tuning

Evaluation

Evaluators  
Metrics

Success!

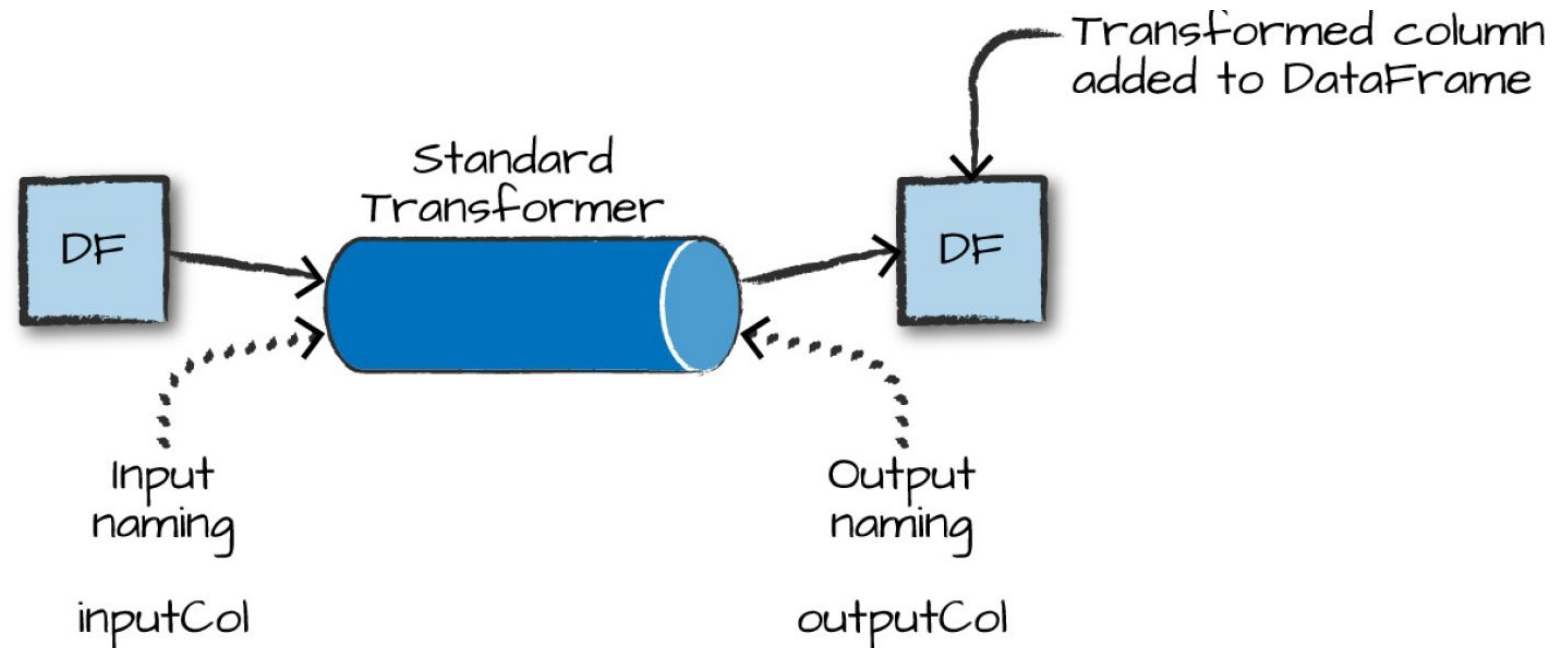


# Pipeline components

## Transformer and Estimator

A **Transformer** can transform one DataFrame into another DataFrame.

An **Estimator** can be fit on a DataFrame to produce a Transformer



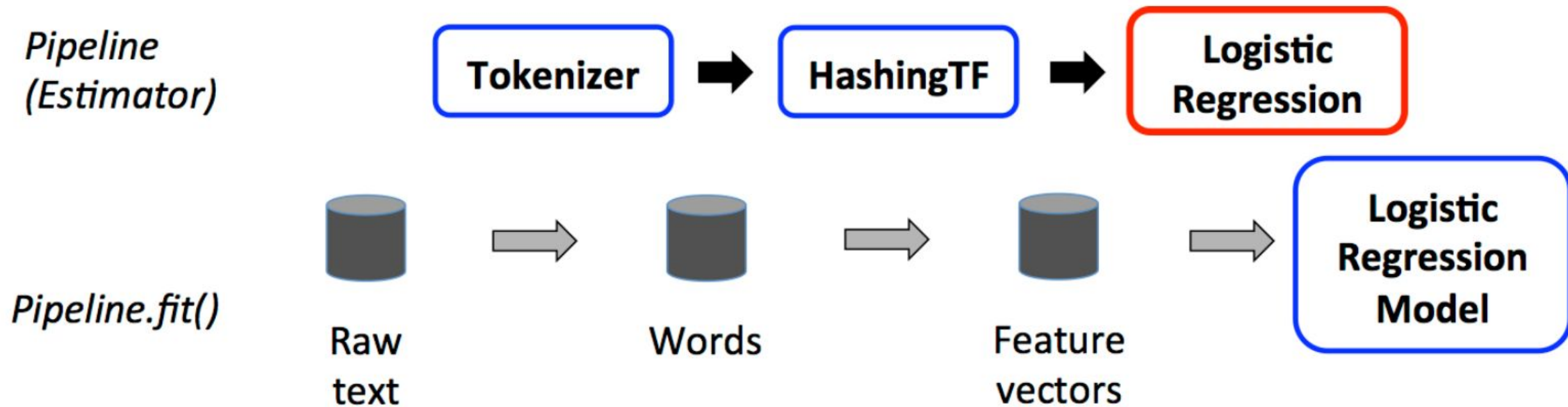
# Low-level data types

## Vector

```
1 from pyspark.ml.linalg import Vectors
2 denseVec = Vectors.dense(1.0, 2.0, 3.0)
3 size = 3
4 idx = [1, 2] # locations of non-zero elements in vector
5 values = [2.0, 3.0]
6 sparseVec = Vectors.sparse(size, idx, values)
```

# A Pipeline

blue - transformations, red - estimator

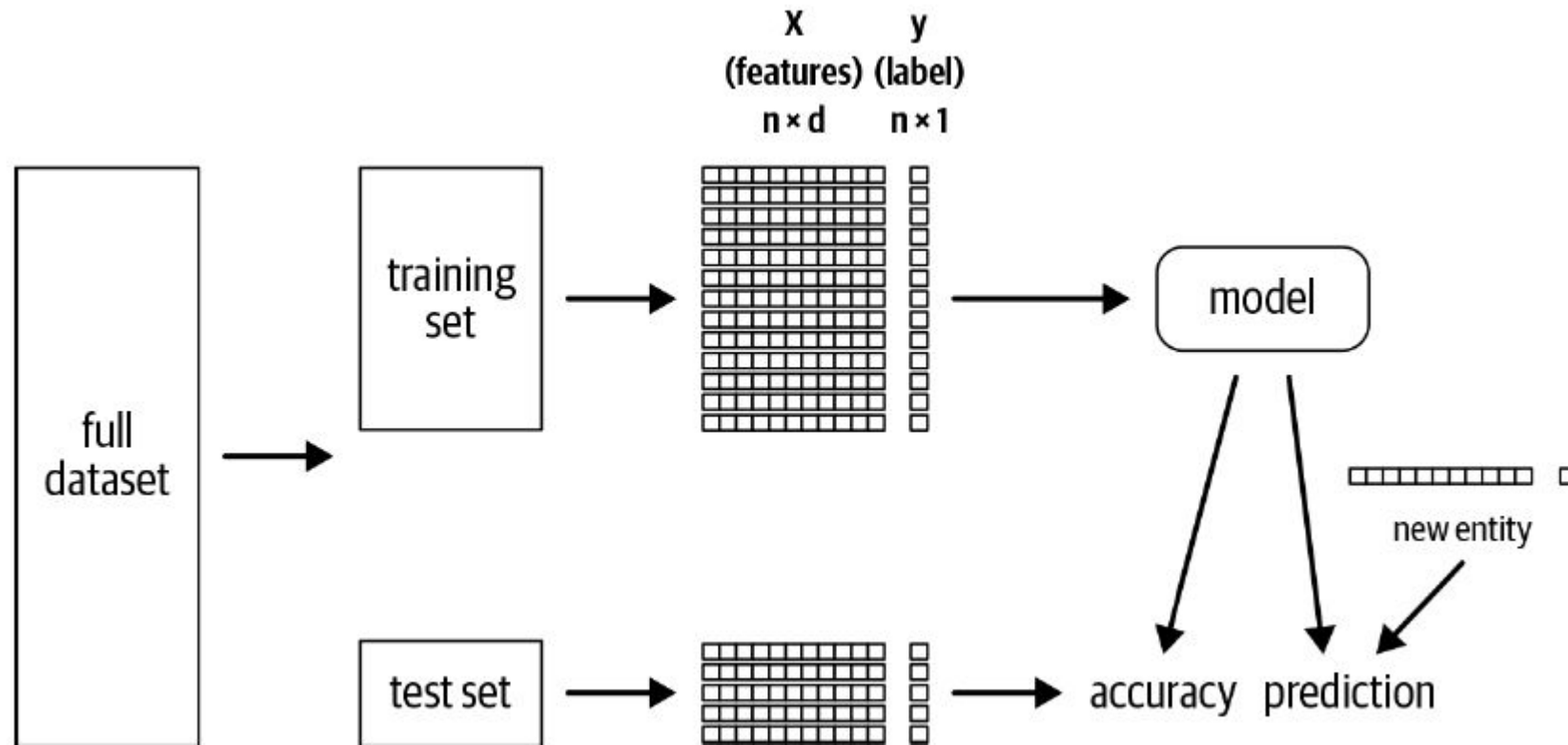




# Logistic Regression

```
1 from pyspark.mllib.linalg import Vectors
2 from pyspark.ml.classification import LogisticRegression
3 from pyspark.ml.param import Param, Params
4
5 # Prepare training data from a list of (label, features) tuples.
6 training = sqlContext.createDataFrame([
7     (1.0, Vectors.dense([0.0, 1.1, 0.1])),
8     (0.0, Vectors.dense([2.0, 1.0, -1.0])),
9     (0.0, Vectors.dense([2.0, 1.3, 1.0])),
10    (1.0, Vectors.dense([0.0, 1.2, -0.5]))], ["label", "features"])
11
12 # Create a LogisticRegression instance. This instance is an Estimator.
13 lr = LogisticRegression(maxIter=10, regParam=0.01)
14 # Print out the parameters, documentation, and any default values.
15 print "LogisticRegression parameters:\n" + lr.explainParams() + "\n"
16
17 # Learn a LogisticRegression model. This uses the parameters stored in
18    lr.
19 model1 = lr.fit(training)
```

# Creating Training and Test Data Sets



# Ссылки

<https://spark.apache.org/docs/latest/ml-pipeline.html>

**Спасибо!**  
**Каждый день**  
**вы становитесь**  
**лучше :)**

