# 10 things I learned
# about writing data pipelines
# in Python and Spark

#7 will blow your mind!

# In the next 30 minutes

- Patterns that you can apply today
- There will be <code>
- No machine learning - yet
- Use with caution*
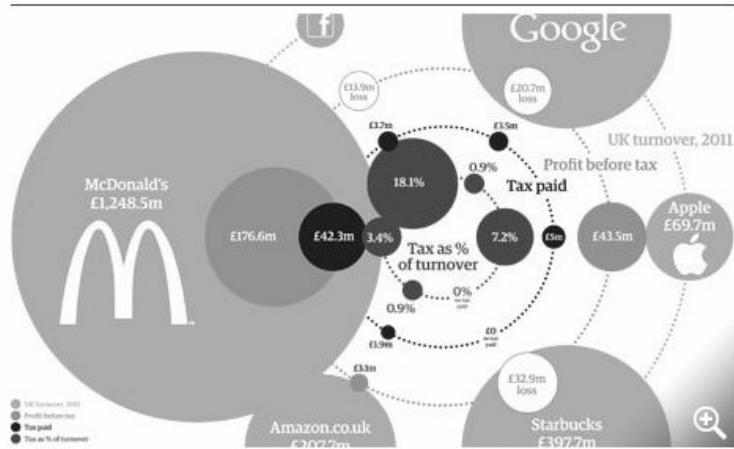- Mostly Python
- Who is this talk for?

# Background

# How much tax do Starbucks, Facebook and the biggest US companies pay in the UK

How much do the biggest US companies in the UK pay in tax?
• Download the data
• More data journalism and data visualisations from the Guardian



How much tax is paid by major US companies in the UK? Click image to embiggen

# Background



Enhanced Financials

Search for a company or director

Ali

**Ali Zaidi**
DUEDIL LIMITED

Companies / Directors in
## All Your Lists

View your lists

**Your Lists (10)**

| | |
|---|---|
| All Your Lists | |
| DD Competitors | 12 |
| Miscellaneous Companies | 1 |
| Customers | 0 |
| Suppliers | 0 |
| Prospects | 0 |
| Competitors | 0 |

View All

**DueDil Connect**

Browse your network of companies
Browse your team's network
Look at your company (duedil.com)
Add another account

**10** Lists
in your account

Companies / Directors
in the news

Companies / Directors
with opportunities

Companies / Directors
with risks

### In The News

Companies featured in the news lately. Great insights and conversation starters.

✖ No news matching the companies / directors in this list

### Opportunities

Sales triggers: companies with new accounts, new directorships, reviewing their budget...

| | | | | |
|---|---|---|---|---|
| DUEDIL LIMITED | 31 Oct 2015 | Exiting budget window | | +2 |
| DUEDIL LIMITED | 30 Sep 2015 | Still within budget window | | +2 |
| DUEDIL LIMITED | 31 Aug 2015 | Still within budget window | | +2 |
| DUEDIL LIMITED | 31 Jul 2015 | Entering budget window | | +2 |
| DUEDIL LIMITED | 04 May 2015 | Directorship details updated | | +2 |

Load more

### Risks

Risk factors: companies with new credit scores, CCJ updates, status change...
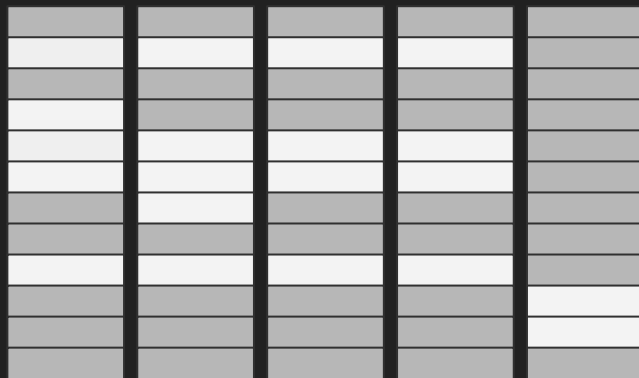
| | | | | |
|---|---|---|---|---|
| DUEDIL LIMITED | 22 Jun 2015 | New credit limit available | | +2 |
| DUEDIL LIMITED | 22 Jun 2015 | New accounts available | | +2 |
| DUEDIL LIMITED | 18 Jun 2015 | New credit score available | | +2 |
| DUEDIL LIMITED | 14 Oct 2014 | New credit score available | | +2 |
| DUEDIL LIMITED | 14 Oct 2014 | New credit limit available | | +2 |

Load more
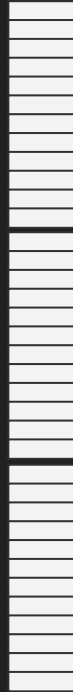
Press

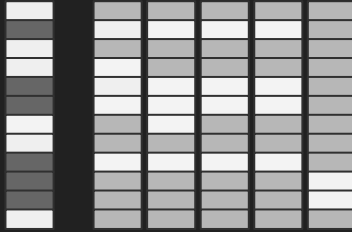# The Data



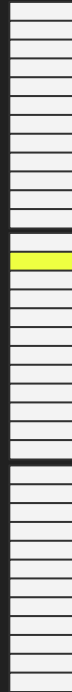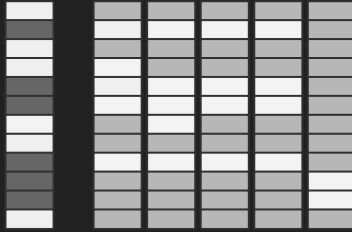Summary          Specialized          Standard Duedil

Financial data

Simple Metrics

Trend Metrics

Search
Segmenting
Benchmarking
Analysis

Financial data

Simple Metrics

Trend Metrics

Search
Segmenting
Benchmarking
Analysis

```python
def gearing(liabilities, shareholder_funds):
    return (liabilities / shareholder_funds) * 100
```

```python
def gearing(rec):
    return (rec['liabilities'] /
            rec['shareholder_funds']) * 100
```

```python
def gearing(rec):
    return (rec['liabilities'] * 1.0 /
            rec['shareholder_funds']) * 100
```

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
         return None
    return (rec['liabilities'] * 1.0 /
            rec['shareholder_funds']) * 100
```

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['shareholder_funds'] == 0:
        return None
    return (rec['liabilities'] * 1.0 /
            rec['shareholder_funds']) * 100
```

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['shareholder_funds'] == 0:
        return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
        return None
    return (rec['liabilities'] * 1.0 /
            rec['shareholder_funds']) * 100
```

# #1 Cleaning code will eat compute code and code quality

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['shareholder_funds'] == 0:
        return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
        return None
    return (rec['liabilities'] * 1.0 /
            rec['shareholder_funds']) * 100
```

```python
from __future__ import division

def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
         return None
    if rec['shareholder_funds'] == 0:
         return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
         return None
    return (rec['liabilities'] /
            rec['shareholder_funds']) * 100
```

# PEP 238 -- Changing the Division Operator

| PEP: | 238 |
|---|---|
| Title: | Changing the Division Operator |
| Author: | Moshe Zadka <moshez at zadka.site.co.il>, Guido van Rossum <guido at python.org> |
| Status: | Final |
| Type: | Standards Track |
| Created: | 11-Mar-2001 |
| Python-Version: | 2.2 |
| Post-History: | 16-Mar-2001, 26-Jul-2001, 27-Jul-2001 |

# #2 Disable Python's default Integer division*

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['shareholder_funds'] == 0:
        return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
        return None
    return (rec['liabilities'] /
            rec['shareholder_funds']) * 100
```

```python
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
        return None
    try:
        return (rec['liabilities'] /
                rec['shareholder_funds']) * 100
    except ZeroDivisionError:
        return None
```

```python
@handle_div_by_0
def gearing(rec):
    if 'shareholder_funds' not in rec or
        'liabilities' not in rec:
        return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
        return None
    return (rec['liabilities'] /
            rec['shareholder_funds']) * 100
```

```python
from functools import wraps

def handle_div_by_0(f):
    @wraps(f)
    def wrapper(*arg, **kwargs):
        try:
            return f(*arg, **kwargs)
        except ZeroDivisionError:
            return None
    return wrapper
```

# #3 Handle ZeroDivisionError in decorators**

```python
def gearing(rec):
    if 'shareholder_funds' not in  rec or
        'liabilities' not in  rec:
         return None
    if rec['liabilities'] is None or
        rec['shareholder_funds'] is None:
         return None
    return (rec['liabilities'] /
            rec['shareholder_funds']) * 100
```

```python
def gearing(shareholder_funds, liabilities, **kwargs):
    if liabilities is None or shareholder_funds is None:
        return None
    return (liabilities / shareholder_funds) * 100
```

```python
def gearing(shareholder_funds, liabilities, **kwargs):
    '''Replace gearing(rec) with gearing(**rec)'''
    if liabilities is None or shareholder_funds is None:
        return None
    return (liabilities / shareholder_funds) * 100
```

```python
def gearing(shareholder_funds, liabilities, **kwargs):
    '''Replace gearing(rec) with gearing(**rec)
    And this would still not work :P
    '''

    if liabilities is None or shareholder_funds is None:
        return None
    return (liabilities / shareholder_funds) * 100
```

```python
def gearing(shareholder_funds=None, liabilities=None, **kwargs):
    if liabilities is None or shareholder_funds is None:
        return None
    return (liabilities / shareholder_funds) * 100
```

# #4 Use **kwargs to unpack records into local variables

```python
def gearing(shareholder_funds=None, liabilities=None, **kwargs):
    if liabilities is None or shareholder_funds is None :
        return None
    return (liabilities / shareholder_funds) * 100
```

```python
def gearing(shareholder_funds=None, liabilities=None, **kwargs):
    if None in [shareholder_funds, liabilities] :
        return None
    return (liabilities / shareholder_funds) * 100
```

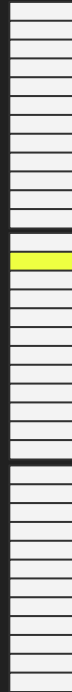# #5 Refactor checks for many variables having a value to "in"

```python
def gearing(shareholder_funds=None, liabilities=None, **kwargs):
    if None in [shareholder_funds, liabilities]:
        return None
    return (liabilities / shareholder_funds) * 100
```
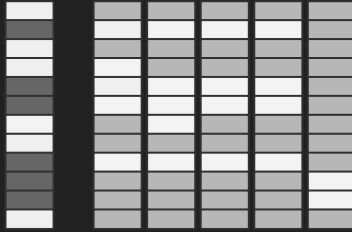
```python
@handle_none_args
def gearing(shareholder_funds, liabilities, **kwargs):
    return (liabilities / shareholder_funds) * 100
```

```python
from functools import wraps
from inspect import getargspec

def handle_none_args(f):
    @wraps(f)
    def wrapper(**kwargs):
        for arg in getargspec(f)[0]:
            if kwargs.get(arg) is None:
                return None
        return f(**kwargs)
    return wrapper
```
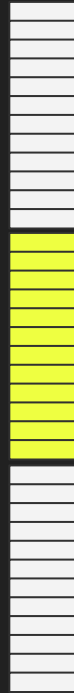
# #6 Be a craftsman but don't be a perfectionist

Financial data

Simple Metrics

Trend Metrics

Search
Segmenting
Benchmarking
Analysis

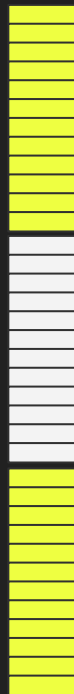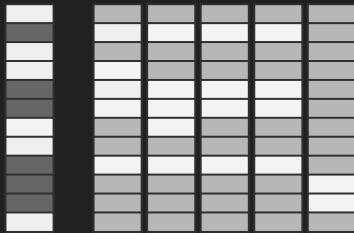Financial data

Simple Metrics

Trend Metrics

Search
Segmenting
Benchmarking
Analysis

```python
def calc_metrics(rec):
    '''Without kwarg unpacking'''
    rec['gearing'] = gearing(
        rec['shareholder_funds'],
        rec['liabilities']
    )
    # 200 more...
```

```python
METRICS_MAP = [
    ('gearing', gearing),
    # 100 more
]

def calc_metrics(rec):
    for field, func in METRICS_MAP:
        rec[field] = func(**rec)
```

# #7 Mapping dictionaries are awesome

Financial data

Simple Metrics

Trend Metrics

Search
Segmenting
Benchmarking
Analysis

```
TREND_METRICS_MAP = {
    ...
    prev_3: {
        'turnover': (cagr3, 'turnover'),
        # 50 more
    }
    ...
}
```

```python
def sliding_window(iterable, size=2):
    iters = tee(iterable, size)
    for i in xrange(1, size):
        for each in iters[i:]:
            next(each, None)
    return izip(*iters)
```

```python
def sliding_window(iterable,  size=2):
    iters = tee(iterable, size)
    for i in xrange(1, size):
        for each in iters[i:]:
            next(each, None)
    return izip(*iters)

sliding_window(accounts)
sliding_window(accounts, 3)
```

```python
def sliding_window(iterable, size):
    ...

prev_2 = partial(sliding_window, size=2)
prev_3 = partial(sliding_window, size=3)

prev_2(accounts)
prev_3(accounts)
```

# #8 Use partials to make APIs more accessible

```
>>> fns = []
>>> for i in range(3):
...     fns.append(lambda: i)
...
>>> for fn in fns:
>>>     print fn()
```

```
>>> fns = []
>>> for i in range(3):
...     fns.append(lambda: i)
...
>>> for fn in fns:
>>>     print fn()
2
2
2
```

```
>>> fns = []
>>> for i in range(3):
...   def f(i=i):
...     return i
...   fns.append(f)
...
>>> for fn in fns:
...   print fn()
0
1
2
```

# #9 Be aware of common Python gotchas

```
>>> type(1)
int
>>> type(1.0)
float
```

```
>>> data = [{'x': 353463435324.0}, {'x': 4534123897.0}]

>>> rdd = sc.parallelize(data)

>>> fields = [types.StructField('a', types.FloatType())]

>>> schema = types.StructType(fields)

>>> sql.jsonRDD(rdd, schema=schema).collect()
```

```
>>> data = [{'x': 35346343 5324.0}, {'x': 453412 3897.0}]

>>> rdd = sc.parallelize(data)

>>> fields = [types.StructField('a', types.FloatType())]

>>> schema = types.StructType(fields)

>>> sql.jsonRDD(rdd, schema=schema).collect()
[Row(a=3534634352 64.0), Row(a=45341240 32.0)]
```

# #10 Be aware of common coding gotchas

# WhatsApp ups its group chat limit to 256 users

by MIX — 3 months ago in APPS

# Identifying the previous set of accounts

history

subject accounts
date = d
period = m months

previous accounts are the most recent accounts in the range:
d - (m - 1) months
to
d - (m + 12) months

d - (m - 1) months                      d - (m +12) months

**This allows:**
- Correct calculation of deltas and percentage deltas even if there are multiple accounts filed for largely overlapping time periods;
- Calculation of average balance sheet values over an accounting period for use in certain ratios and KPIs.

# #Bonus
Nothing beats awesome specs

www.duedil.com/careers

# Questions

alixedi commented on Nov 11, 2015

http://jenkins.duedil.net/job/snapshot-run-generic/396/console

Took 7 hours, 14 mins

Hello: @alixedi, GitHub: alixedi, Essays: alixedi.github.io