

Introduction to Data Analytics & Python

Hans van der Heijden



BUSINESS
SCHOOL

Learning outcomes

- To understand and appreciate the tools of data science, in an accounting context
- To be able to detect unusual transactions in large business data sets
- To be able to articulate why they are unusual

Agenda

1. Why data science
2. Introduction to Python
3. Example: Three-way match in procure-to-pay
4. Three-way match in Python

Why data science?

Business Data

Tesco

- 70 transactions per second
- 6 million transactions per day

Amazon

- 306 transactions per second
- 26.4 million transactions per day

- Twitter

- 6000 tweets per second
- 500 million tweets per day

Data science

The role of data science is to take raw business data and turn it into something meaningful, ie. an insight, a trend, an answer to a specific question

Input: Raw Data

- Large tabular datasets

Output: Excel file

- File contains human-readable table or chart

Data science in accounting

Data scientists are in demand given the reliance on financial data in accounting, auditing, banking & finance.

Data Scientist salaries in United Kingdom

£56,494 per year

Based on 48,655 salaries



[Data Scientist salaries by company in United Kingdom](#)

Data science, statistics, and programming

Data science is not the same as statistics but statistical techniques are often used in data science

- Summarising (using descriptive statistics)
- Forecasting (using statistics such as linear regression)

Data science is not programming (or coding) either, but programming is often used in data science as well

- E.g., creating a for-loop to go through all business transactions

Data science shortcomings

1. General expectations gap between what's possible and what's believed to be possible
2. “Garbage in / Garbage out”
 - Poor quality data can still lead to an “exact” final outcome
3. “Black box” manipulations may limit trust in final outcomes

Data science in auditing

- Past practice – sample based auditing
 - Auditor takes a sample of all transactions, audits just these, and forms opinion.
- Sampling not always adequate
 - 25 out of 80,000 transactions (insurer)
 - 10 out of 6 million mortgages (retail bank)
- New practice – population-based auditing
 - Auditor takes *all* transactions
 - Runs them through software to see which ones are suspect
 - Audits suspect transactions & forms opinion.

Examples of checks (1)

- Any data entry made in weekend or late at night
- Any suspiciously round number
- Any payment just under particular threshold
- Any payment with a description containing the word 'bonus' or similar
- Any missing invoice number

Examples of checks (2)

- Any invoice with an outgoing payment higher than stated on purchase order or goods receipt

We will proceed with this check, and implement it in Python.

Introduction to Python

Python

- General purpose computer language
 - Created 1991 in Amsterdam
- Popular in part because of “human readability”
 - (it is easy to see what the code does)
- Named after Monty Python
- No 1 language for data science

Setting up Python

<https://colab.research.google.com>

The screenshot shows the Google Colab interface. The title bar says "Welcome To Colaboratory - Colab". The left sidebar has a "Table of contents" section with links to "Getting started", "Data science", "Machine learning", "More Resources", "Machine Learning Examples", and "Section". The main content area is titled "What is Colaboratory?". It explains that Colaboratory allows you to write and execute Python in your browser with zero configuration, free access to GPUs, and easy sharing. It also links to "Introduction to Colab". Below this, there's a section titled "Getting started" with a note about Colab notebooks being interactive environments. A code cell is shown with the following Python code:

```
[ ] seconds_in_a_day = 24 * 60 * 60
seconds_in_a_day
```

The output of the cell is "86400". A note says to execute the code by clicking and pressing Command/Ctrl+Enter. It also mentions that variables defined in one cell can be used in others. Another code cell is shown with the following Python code:

```
[ ] seconds_in_a_week = 7 * seconds_in_a_day
seconds_in_a_week
```

Colab basics

Notebook

- Basic structure of Colab

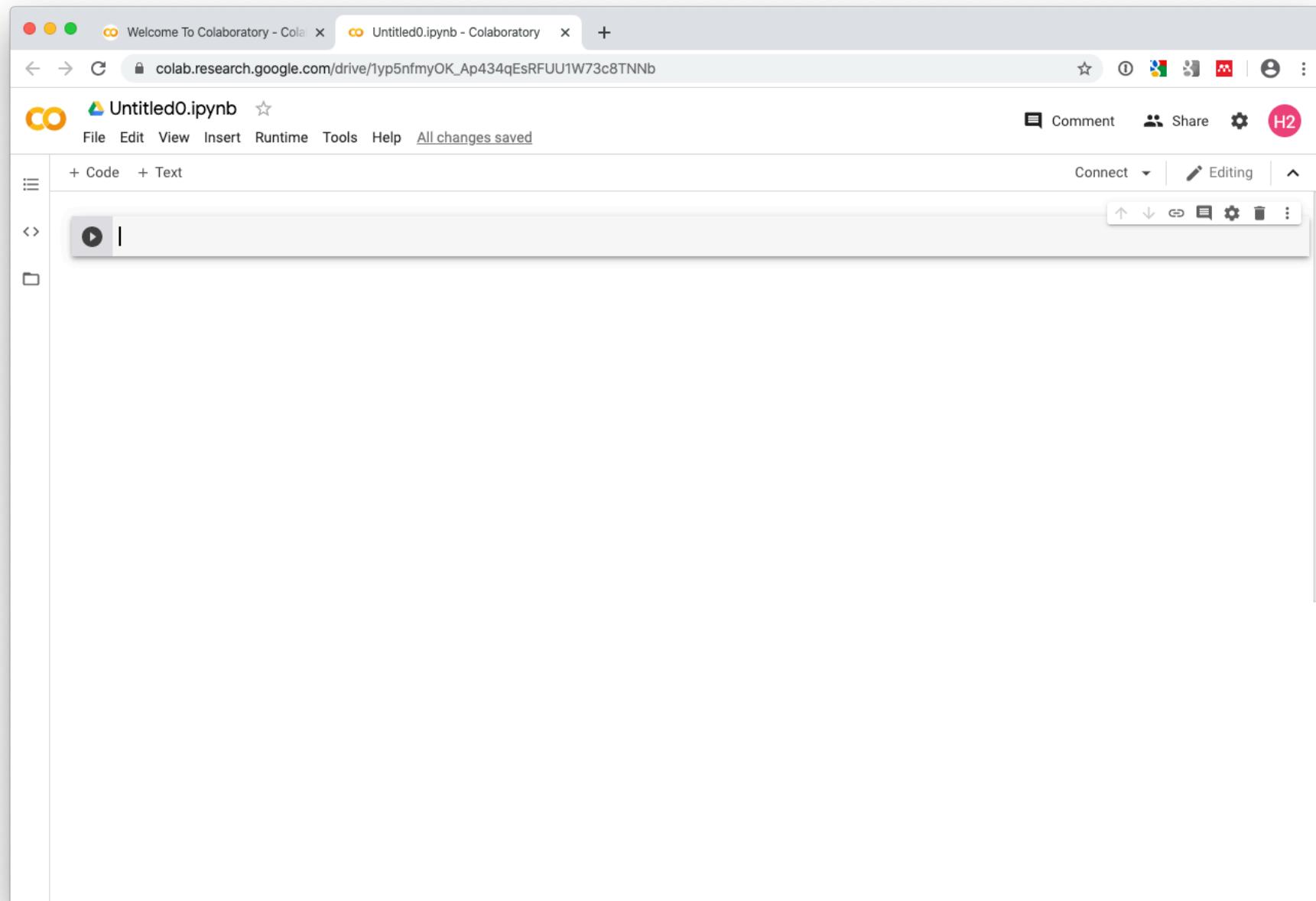
Python run-time

- A run-time is the "engine" running your code, in our case Python

Cells

- Notebooks have cells, cells contain the code
- To 'run' a cell, put cursor in cell and type Shift-Enter

New notebook



Core features

- Input-Output
 - Type `print('hello world')` and press shift-enter
- Variables
- Boolean conditions (true and false)
- Functions
- Libraries

Functions

- `min(1, 2)`
- `list = (1, 2, 3, 4, 5)`
`sum(list)`
- `round(3.2223, 2)`

These are examples of functions: they take input and produce output.

Create your own functions

```
def value(price, quantity):  
    return price * quantity
```

- Starts with def
- Note the indent on the second line

```
value(10, 100) (shift-enter)  
    • 1000
```

Example: Three-way match in procure-to-pay

Three-way match

Three-way match is a type of invoice verification

Three documents must match before payment goes out to a supplier

- Purchase Order (price & quantity)
- Goods Received Document (quantity)
- Invoice Received (price & quantity)

Anything unusual requires further investigation
Potentially human error, misunderstanding, or in the worst case fraud

Three business documents

- Purchase order
 - Created by purchase department in *our* company
 - Specifies quantity & price
- Goods received
 - Created by warehouse in *our* company
 - Specified quantity, refers to purchase order
- Invoice received
 - Created by *supplier*, and then sent to *our* accounting department
 - Specifies quantity, price, total value to pay
 - Refers to purchase order

Purchase Order

Miss Red

PURCHASE ORDER

Sumigi Support's business

purchase order no
date

430864
24 Dec 2018

Product code	Quantity	Description
85206	300	Apple Trees, less trade discount @ 20%

AUTHORISED:

Signature:

Miss Red

Date: *24 Dec 2018*

Invoice Received

Sumigi Support's business

INVOICE

Miss Blue

invoice no **83357**
account **MISS13**
your reference **769031**
date **4 Jan 2019**

Product code	Description	Quantity	Price	Unit	Total	Discount @20%	Net
85437	T-shirts	100	1.25	each	125.00	25.00	100.00

Terms

Net monthly
Carriage paid
E & OE

Goods total	100.00
VAT @21%	21.00
Grand total	121.00

An example

- Our company wishes to purchase 51 boxes of cheese at £72.07 each
 - Our purchasing department creates the purchase order and sends to supplier.
- 51 boxes of cheese arrive
 - Our warehouse creates the goods receipt document. Refers back to the purchase order.
- Cheese supplier sends invoice of £3675.57
 - Our accounting department receives the invoice. It contains the purchase order reference.
- Company pays £3675.57 to cheese supplier
 - Accounting department carries out three-way match, and if OK pays out.

Invoice verification

Invoices must always be verified, because what can happen?

1. Invoice is fake (no purchase order)
2. Invoice is premature (no goods received yet)
3. Invoice amount is incorrect (higher than price * quantity)
4. Price is different on invoice than on purchase order
5. Quantity is different on purchase order
6. Quantity is different on goods receipt

Three-way match in Python

Spreadsheets

<https://github.com/accountinglab/audit-analytics>

The screenshot shows a GitHub repository page for 'accountinglab/audit-analytics'. The repository has 1 branch and 0 tags. The main commit is 'LICENSE' (Initial commit) by 'accountinglab' at 3 minutes ago. Following this are 14 other commits, all titled 'student-X.xlsx' where X is a number from 0 to 8, each added via upload 17 minutes ago. The repository has 1 watch, 0 stars, and 0 forks. The 'About' section notes 'No description, website, or topics provided.' and lists the 'MIT License'. The 'Releases' section indicates 'No releases published' and 'Create a new release'. The 'Packages' section shows 'No packages published' and 'Publish your first package'.

File	Description	Time Ago
LICENSE	Initial commit	20 minutes ago
assessment-form.docx	Add files via upload	19 minutes ago
practice-set-1.pdf	Add files via upload	13 minutes ago
practice-set-2.pdf	Add files via upload	3 minutes ago
student-0.xlsx	Add files via upload	17 minutes ago
student-1.xlsx	Add files via upload	17 minutes ago
student-2.xlsx	Add files via upload	17 minutes ago
student-3.xlsx	Add files via upload	17 minutes ago
student-4.xlsx	Add files via upload	17 minutes ago
student-5.xlsx	Add files via upload	17 minutes ago
student-6.xlsx	Add files via upload	17 minutes ago
student-7.xlsx	Add files via upload	17 minutes ago
student-8.xlsx	Add files via upload	17 minutes ago

Examine your spreadsheet

Click on student-X.xlsx where X is the last digit of your student ID number

Click on 'Download'

The screenshot shows a GitHub repository page for 'audit-analytics / student-0.xlsx'. At the top left, there's a dropdown menu labeled 'main' and a link to 'audit-analytics / student-0.xlsx'. On the right, there are buttons for 'Go to file' and '...'. Below the header, there's a section for 'accountinglab' with a 'Add files via upload' button and a timestamp 'Latest commit 9104685 18 minutes ago'. It also includes a 'History' link. Underneath, it says '1 contributor'. The main content area shows the file 'student-0.xlsx' with a size of '206 KB'. To the right of the file name are three icons: 'Download', 'View raw', and a trash bin. A large green arrow points upwards from the bottom right towards the 'Download' button.

Examine your spreadsheet

There are three worksheets in the workbook:

- Purchase Orders
- Goods Received Documents
- Invoices Received

Spreadsheet format

- Rows represent documents
- Columns represent values
- First line identifies the column, for example, "Purchase_Order_Amount_Local"
- Spreadsheet format is a common standard, the US Audit Data Standard (no global standard exists)
<https://www.aicpa.org/interestareas/frc/assuranceadvisoryservices/auditdatastandards.html>

Examine invoices

On the invoice worksheet, identify the first invoice:

- Invoice_ID (70011)
- Invoice_Amount_Local (this is the amount in local currency)
- Purchase_Order_ID (first one is 4500000011)

On the purchase order worksheet, identify the corresponding purchase order

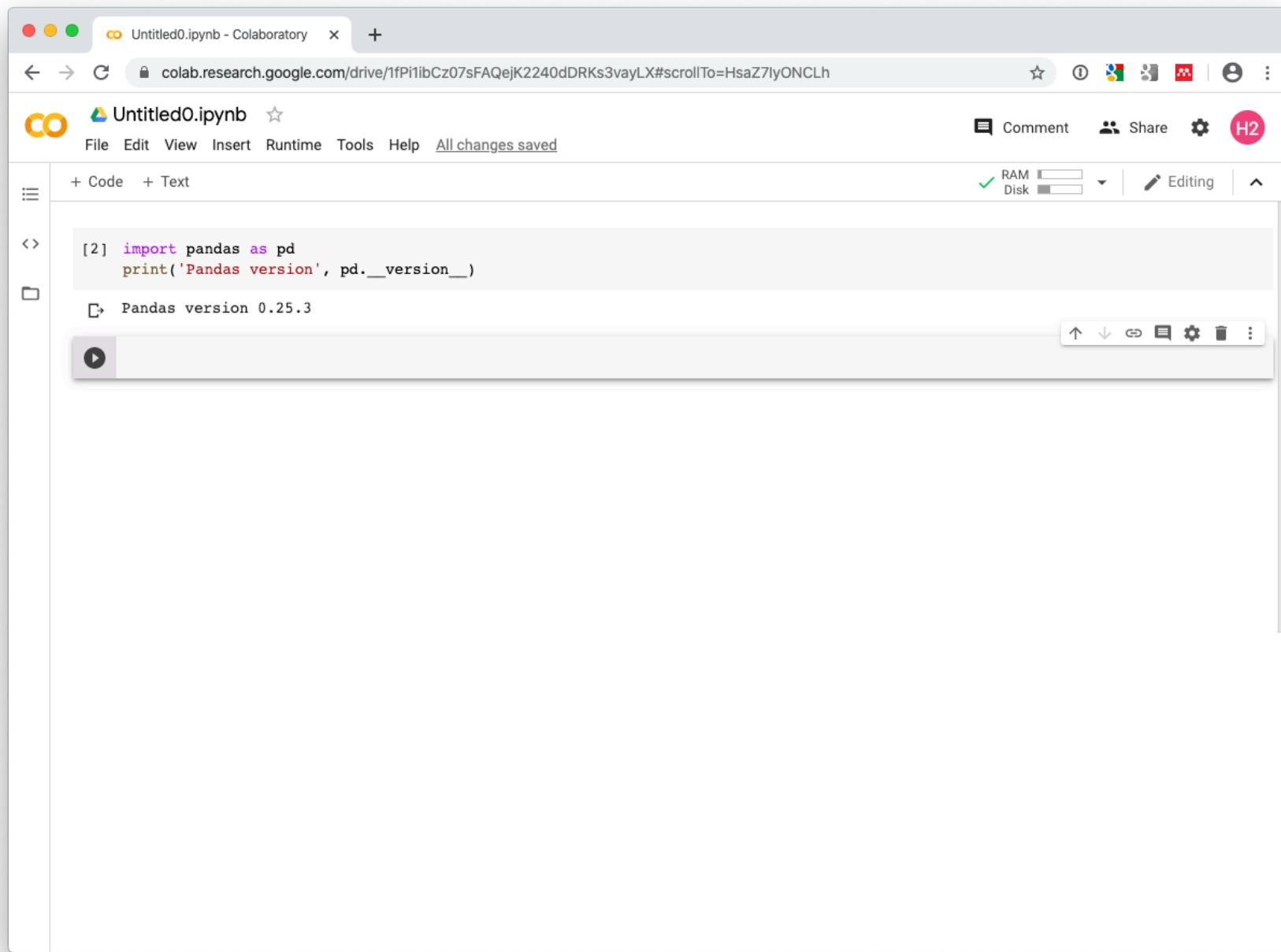
- Purchase_Order_ID (4500000011)
- Purchase_Order_Amount_Local

The amounts should be the same. Are they?

Python

- Manually performing the three-way match (checking purchase order, goods receipt, and invoice) is tedious and time-consuming
- With 11 lines of code in Python, we can achieve the same thing, and check 1000 documents in seconds.
- Let's go...

Import Pandas library



The screenshot shows a Google Colab notebook titled "Untitled0.ipynb". The code cell contains the following Python code:

```
[2] import pandas as pd
      print('Pandas version', pd.__version__)
```

The output cell shows the result of the code execution:

```
Pandas version 0.25.3
```

The interface includes standard Colab controls like a play button, a code editor tab, a text editor tab, and various toolbar icons.

Import spreadsheet

The screenshot shows a Google Colab interface with a Jupyter notebook titled "Untitled0.ipynb". The notebook contains the following Python code:

```
[1] import pandas as pd
     print('Pandas version', pd.__version__)

[2] student_id_ends_with = '0' # <-- Enter the last digit of your student ID inbetween the single quotes.
     print('Your student ID ends with:', student_id_ends_with)

[3] url = 'https://github.com/accountinglab/audit-exercise/raw/master/student-' + student_id_ends_with + '.xlsx'
     print('I will download a spreadsheet from:', url)

[5] pos = pd.read_excel(url, sheet_name='Purchase_Orders')
     goods = pd.read_excel(url, sheet_name='Goods_Received')
     invoices = pd.read_excel(url, sheet_name='Invoices_Received')
```

View sample rows in Python

```
[6] pos.sample(n=5)
```

od	Business_Unit_Code	Supplier_Account_ID	Entered_By	Purchase_Order_Amount_Local	Purchase_Order_Local_Curren
M2	Dairy	Violet & Co	Procurement	1835.35	U
M4	Dairy	Lavender Ltd	Procurement	1285.20	U
M1	Dairy	Lavender Ltd	Procurement	16287.82	U
M5	Dairy	Orchid & Sons	Procurement	5909.74	U
M1	Dairy	Iris & Sons	Procurement	6531.72	U

```
[7] goods.sample(n=5)
```

_Account_ID	Business_Unit_Code	Purchase_Order_ID	Purchase_Order_Date	Entered_By	Receipt_Amount_Local	Receipt_Currency
Azalea & Sons	Dairy	4500000067	2020-01-10	Warehouse	1887.05	U
Iris & Sons	Dairy	4500000297	2020-02-13	Warehouse	1499.30	U
Hibiscus Plc	Dairy	4500000782	2020-04-24	Warehouse	5645.70	U
nolia Importers	Dairy	4500000438	2020-03-05	Warehouse	3465.45	U
Iris & Sons	Dairy	4500000778	2020-04-24	Warehouse	11796.36	U

```
[8] invoices.sample(n=5)
```

	Invoice_ID	Invoice_Number	Invoice Date	Invoice_Fiscal_Year	Invoice_Period	Invoice_Due_Date	Purchase_Order_ID
160	70174	70174	2020-01-27	2020	M1	2020-04-26	45000001
351	70368	70368	2020-02-27	2020	M2	2020-05-27	45000003
130	70144	70144	2020-01-22	2020	M1	2020-04-21	45000001
963	70988	70988	2020-05-26	2020	M5	2020-08-24	45000009
119	70133	70133	2020-01-21	2020	M1	2020-04-20	45000001

Loop through invoices

```
[ ] pos = pd.read_excel(url, sheet_name='Purchase_Orders')
    goods = pd.read_excel(url, sheet_name='Goods_Received')
    invoices = pd.read_excel(url, sheet_name='Invoices_Received')

[ ] for index, invoice in invoices.iterrows():
    invoice_id = invoice['Invoice_ID']
    invoice_amount = invoice['Invoice_Amount_Local']
    po = invoice['Purchase_Order_ID']
    po_amount = pos[pos['Purchase_Order_ID']==po]['Purchase_Order_Amount_Local'].values[0]
    goods_amount = goods[goods['Purchase_Order_ID']==po]['Receipt_Amount_Local'].values[0]
    print(invoice_id, 'Purchased:', po_amount, 'Received:', goods_amount, 'Invoiced:', invoice_amount)
```

- 70011 Purchased: 4368.65 Received: 4368.65 Invoiced: 4368.65
- 70012 Purchased: 6467.04 Received: 6467.04 Invoiced: 6467.04
- 70013 Purchased: 7939.6 Received: 7939.6 Invoiced: 7939.6
- 70014 Purchased: 1318.35 Received: 1318.35 Invoiced: 1318.35
- 70015 Purchased: 1473.45 Received: 1473.45 Invoiced: 1473.45
- 70016 Purchased: 16792.31 Received: 16792.31 Invoiced: 16792.31
- 70017 Purchased: 3947.4 Received: 3947.4 Invoiced: 3947.4
- 70018 Purchased: 1370.05 Received: 1370.05 Invoiced: 1370.05
- 70019 Purchased: 8648.4 Received: 8648.4 Invoiced: 8648.4
- 70020 Purchased: 5393.25 Received: 5393.25 Invoiced: 5393.25

The final routine

The screenshot shows a Google Colab notebook titled "practiceset-2.ipynb". The code cell contains the following Python script:

```
[1] import pandas as pd
     print("Pandas version", pd.__version__)

[2] Pandas version 0.25.3

[3] student_id_ends_with = '0' # <-- Replace with last digit of your student ID
     print('Your student ID ends with:', student_id_ends_with)

[4] Your student ID ends with: 0

[5] url = 'https://github.com/accountinglab/audit-exercise/raw/master/student-' + student_id_ends_with + '.xlsx'
     print('I will download a spreadsheet from:', url)

[6] I will download a spreadsheet from: https://github.com/accountinglab/audit-exercise/raw/master/student-0.xlsx

[7] pos = pd.read_excel(url, sheet_name='Purchase_Orders')
     goods = pd.read_excel(url, sheet_name='Goods_Received')
     invoices = pd.read_excel(url, sheet_name='Invoices_Received')

[8] for index, invoice in invoices.iterrows():
        invoice_amount = invoice['Invoice_Amount_Local']
        po = invoice['Purchase_Order_ID']
        po_amount = pos[pos['Purchase_Order_ID']==po]['Purchase_Order_Amount_Local'].values[0]
        is_po_ok = invoice_amount == po_amount
        if not is_po_ok:
            print('PO anomaly for invoice:', invoice['Invoice_ID'])
            goods_amount = goods[goods['Purchase_Order_ID']==po]['Receipt_Amount_Local'].values[0]
            is_goods_ok = invoice_amount == goods_amount
            if not is_goods_ok:
                print('Goods anomaly for invoice:', invoice['Invoice_ID'])

[9]
```

Found an anomaly, now what?

If all goes well, Python will uncover one or more suspect invoices

The routine will only print the suspect invoice ID

You need to look back into the spreadsheet to see what the problem was with the invoice.

For example

- The suspect invoice ID is 70044
- You look up the invoice in the spreadsheet, and you note down the amount, let's assume it was £3,557.25.
- You also note the Purchase Order ID.
- With the purchase order ID, you look up the purchase order amount in the purchase order worksheet, and the goods receipt amount in the goods received amount. You notice these amounts are £1,557.25.
- You conclude: you were over-invoiced by £2000.

Assessment

Try it out!

- Use Practice Set 1 to help you understand the basics of Python
- Use Practice Set 2 to help you code the three-way match
- Use the Assessment Form to complete your exercise and submit
- All files are on
<https://github.com/accountinglab/audit-exercise>