

Python, Cloud and Automation

3. Automation, Cloud and Examples

Jack Minchin

Tourism Economics

2022

Table of Contents

- 1 Cloud Technologies
- 2 Example 1: Whitbread
- 3 Example 2: Country Reports
- 4 Example 3: SweRe
- 5 Example 4: APF

What are cloud solutions?

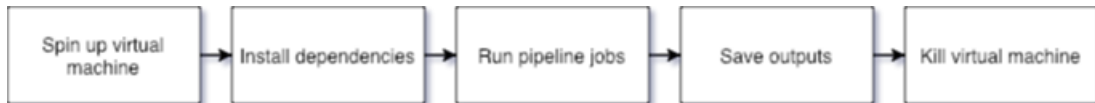
- At the most basic level, cloud computing is the ability to access computers that are not on local premises.
- In practise, cloud technology is an umbrella term for a variety of different computing services (file storage, compute, Platform as a Service (PaaS), cloud functions and pipelines)
- There are a number of cloud providers, at OE we have access to Azure and AWS. Azure is my preferred choice because of its integration with Active Directory (Microsoft) and the more verbose naming conventions (Azure App Service vs. AWS Elastic Beanstalk)

Some examples

- **Azure Portal** The Azure portal is a website where you can manage cloud products.
- **Azure Storage Accounts / Blob Storage** Storage accounts allow the upload and storage of files in the cloud. Blobs (or files) can be downloaded from the Azure Portal or programmatically.
- **Virtual Machines** VMs are computers that run in the cloud they can be used to run application or speed up run times if you are dealing with large datasets, you can interact them with the command line (SSH) or remote desktop (if the VM is running windows)

Pipelines

- Pipelines are automated tasks that run in the cloud they can be triggered by pushing to a git repository, or through an API call.
- In simple terms, they take an ordered set of commands and execute them on a remote virtual machine, the machine is set up specifically for the task and then killed immediately after the pipeline is complete.
- There are number of providers that offer pipelines but I have been using CircleCI because of their free tier.



Setting up a pipeline

- When using circleci, you define pipelines in a text format called YAML.
- This file will define the type of machine that should be created and the steps that the pipeline will include.
- Sometimes, it will also include how the pipeline will be triggered.

```
1 # .circleci/config.yml
2
3 version: 2.1
4
5 # Define the jobs we want to run for this project
6 jobs:
7   exampleJob:
8     executor: # This step defines what type of environment to use
9       name: 'win/default'
10      size: 'medium'
11     steps: # Below we set out the steps to use in the pipeline
12       - checkout
13       - run:
14         name: 'Print some text' # The name of the step
15         command: echo "this is the build job" # This will just return the text
16
17 # Orchestrate our job run sequence
18 workflows:
19   build_and_test:
20     jobs:
21       - exampleJob
```

A completed pipeline

✔ Spin up environment	11s	🔗	⬇
✔ Preparing environment variables	0s	🔗	⬇
✔ Checkout code	8s	🔗	⬇
✔ PIP install	38s	🔗	⬇
✔ Download latest GTS database	23s	🔗	⬇
✔ Download latest Macro database	1m 10s	🔗	⬇
✔ Extract OE model	9s	🔗	⬇
✔ Extract TDM model	9s	🔗	⬇
✔ Run GTS Extract	14s	🔗	⬇
✔ Run Macro Extract	14s	🔗	⬇
✔ Check Extract Exists	2s	🔗	⬇
✔ Install Excel silently	1m 52s	🔗	⬇
✔ Update the links in the excel file	1m 15s	🔗	⬇
✔ Create output Directory	2s	🔗	⬇
✔ Prepare the output file	1m 6s	🔗	⬇
✔ Compress extracts, model and outputs	3s	🔗	⬇
✔ Uploading artifacts	0s	🔗	⬇
✔ Send the output file	3s	🔗	⬇

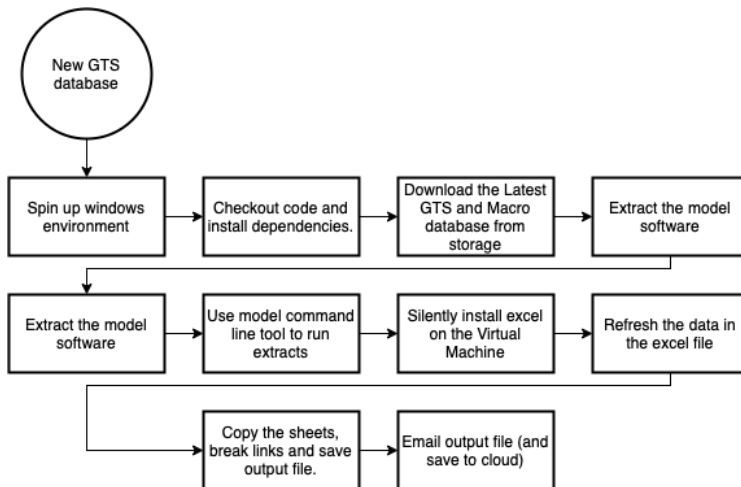
Some considerations

- Pipelines do not only run Python, and many cases Python isn't the best tool to use for orchestrating pipelines. The base level for interacting with pipelines is the command line meaning 99% of tasks can be achieved.
- Any coding language can be used, and there will be a solution allow most programs to be interacted with:
 - Excel → Powershell
 - File structure changes → Powershell
 - Stata → Do files
 - EViews → Scripts
 - OE Models → MDL Command line tool

Example 1: Whitbread Outputs (Basic)

- Whitbread take a very basic quarterly breakdown of some GTS and Macro model indicators.
- The process is:
 - ① Run the macro and GTS extract
 - ② Update the links in the model file
 - ③ Copy the output sheets into a new workbook.
 - ④ Break links and save.
- Manual process time 40 minutes - 1 hour.

Creating a pipeline



Languages & Technologies Used

- Python - for running some of the orchestration
- PowerShell - Interacting with Windows & Excel
- Azure Storage - storing the database files
- Git - Managing the project and codebase
- CircleCI YAML - Writing the pipelines steps

Time saving

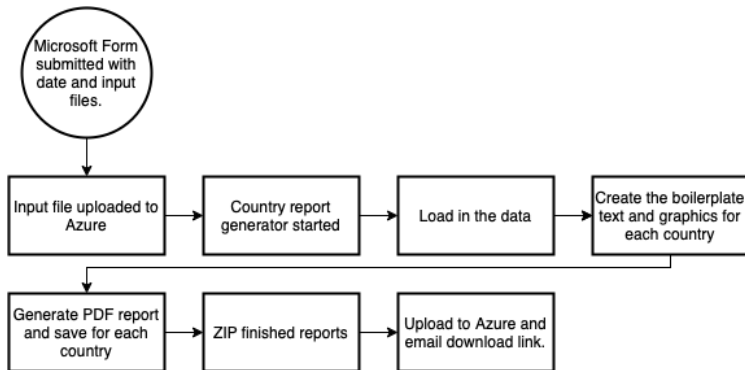
Manual method: 1 hour.

Pipeline method: 8 minutes

Example 2: Country Report Generation (Advanced)

- The country report generator create ~ 40 PDF reports with text and plots.
- The manual process:
 - ① Creation of input file.
 - ② Load excel file for each report.
 - ③ Update the links in each excel file.
 - ④ Rewrite boiler plate text to reflect new data.
 - ⑤ Fix plots where necessary.
 - ⑥ Write custom commentary.
 - ⑦ Save to PDF.

Creating a pipeline



Languages & Technologies Used

- Python - for orchestration, creating text and plots.
- Azure Storage - storing the input and output files
- Azure Functions - triggering the process when a new input file is added.
- Node.JS - application architecture and PDF generation.
- Git - Managing the project and codebase
- CircleCI YAML - Writing the pipelines steps.

Time saving

Manual method: 3 people 3 days.

Pipeline method: 10 minutes + time to write custom text (1 day?).

Sweden Regional TSA PowerPoints

- Creation of 8 regional TSA reports in PowerPoint. The text and graphics are entirely boilerplate.
- Using Python and a package called python-pptx, the code takes the regular output from the Excel model file and outputs the reports in .pptx and PDF.

Time saving

Manual method: Multiple days

Pipeline method: 40 seconds

APF IATA Data Inputs

- We receive a monthly 2020 - To Date data file for PAX and RPK figures from IATA.
- It is in a long form structure with full country names.
- Process:
 - ① Receive input files
 - ② Convert to wide form, change country names to model codes.
 - ③ Extrapolate to the end of the current quarter.
 - ④ Export readin file for model.
- Uses Python notebooks.

APF Automated Checks

- The number of variables in the APF model mean that it is impossible to check each variable series manually.
- We can load an extract from APF and automatically run tests, e.g:
 - ① Ensure no negative values.
 - ② Ensure that current release is with $x\%$ of previous release.
 - ③ Ensure that PAXALL (sum of bilateral indicators) is not greater than PAX (total pax from IATA)