

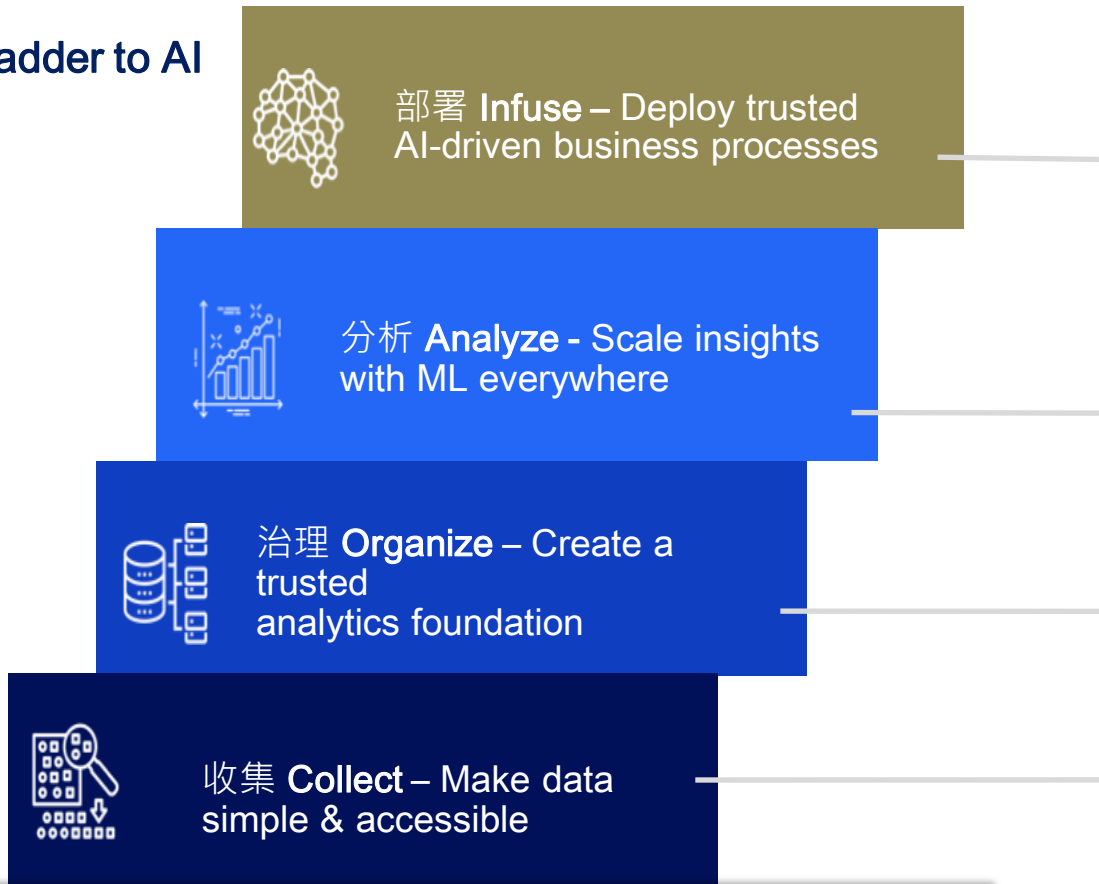
# DATA | 分析大師

全功能且自動化的人工智慧整合平台

## 1. Business Value

## 2. Solution Overview

### Ladder to AI



- Manage models across different stages
- Deploy models and scale for different use cases
- Monitor model and automatically retraining and redeployment
- **Infuse trusted and governed insights**

- Understand and prepare data for analytics
- Build descriptive, predictive & prescriptive models
- Train at scale with support for distributed compute and GPUs
- **Analyze insights on demand**

- Find, Catalog, mask data
- Establish and enforce business policies on data
- Advanced transformation capabilities
- **Organize data so that it can be trusted**

- Collect all Sources of Data
- Connect and discover content from data sources
- Provision new repositories as needed
- **Make it simple & accessible**

**Strong Foundation – Built on “OpenShift”**

# DAS資料科學部署流程

COLLECT

ORGANIZE

ANALYZE

INFUSE

Collect, Connect and Access Data

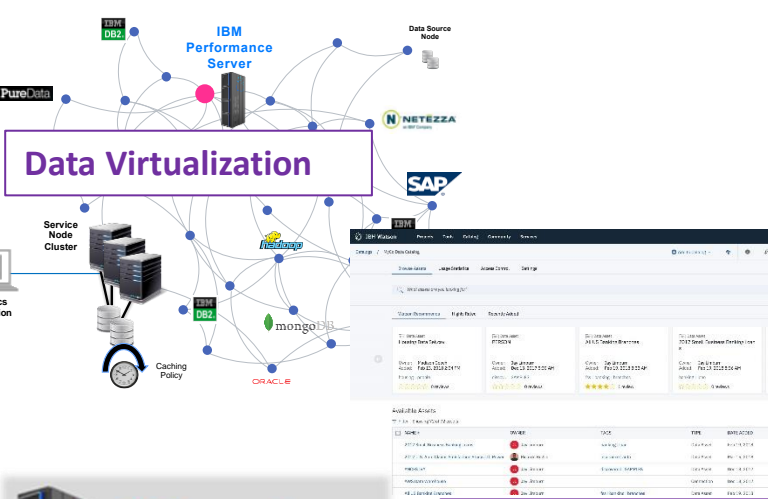
Govern, Search, and Find Data

Understand and Prepare Data for Analysis

Build Descriptive, Predictive, and Prescriptive Models

Model Management and Deployment

Create Analytics Applications



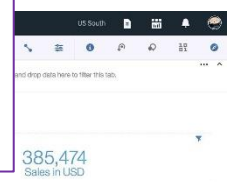
Data Virtualization

- Next Gen Appliances
- Hadoop (Cloudera)
- Db2
- Others

- Watson Knowledge Catalog
- Data Profiling

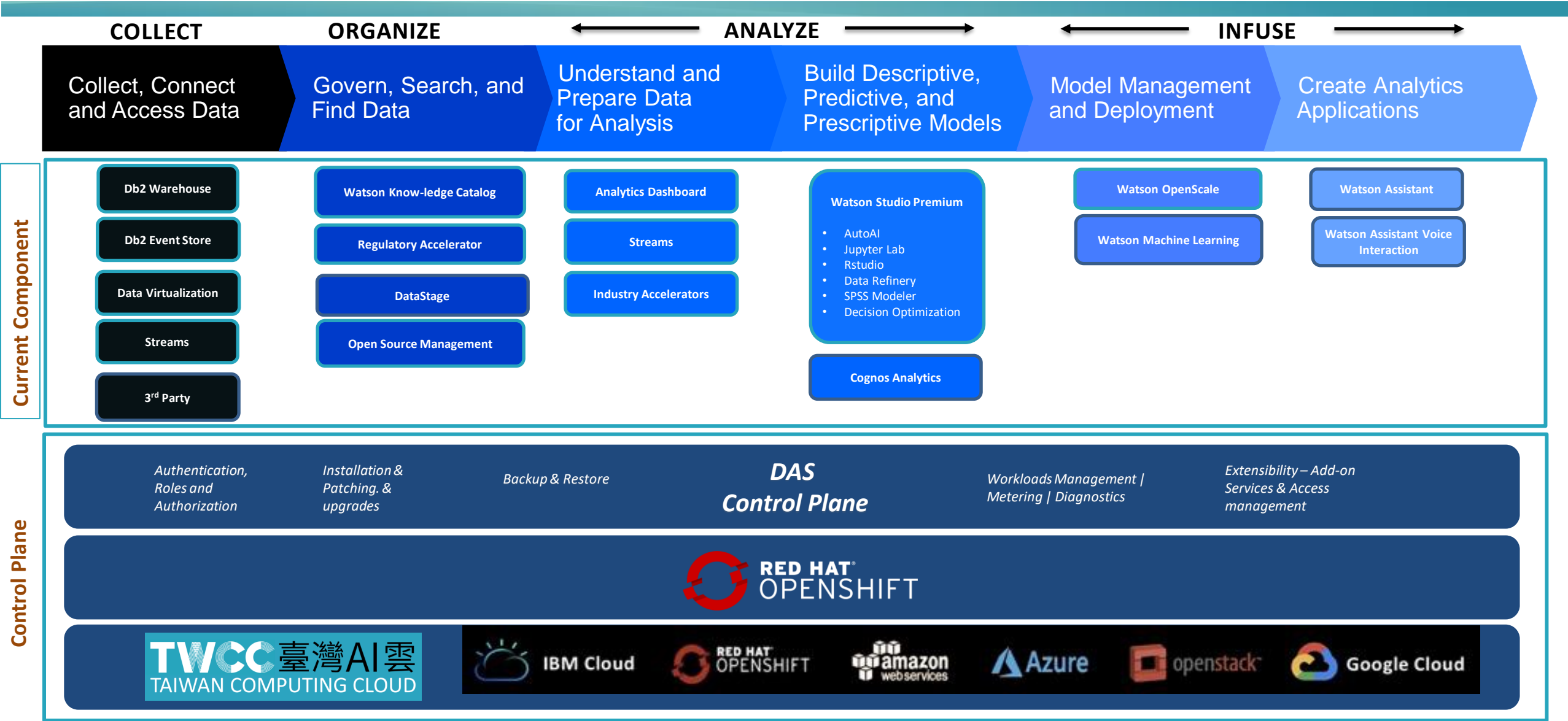


- Watson Studio + AutoAI
- Watson Machine Learning
- Cognos Analytics
- SPSS Modeler Flow



- Watson OpenScale
- Watson Assistant
- Watson Discovery
- Watson APIs
- Others

# DAS底層設計與功能架構



# DAS在資料科學上的循環架構與價值

## 資料推論與監控

在DAS的環境裡，可以透過部署空間的環境來取用API的方式使用模型做資料的推論，另外可以使用分析儀表板來監控資料的狀態和模型推論的成果

## 資料分析

DAS提供了notebook、RStudio、modeler flow、AutoAI等工具可以依需求進行資料分析建立模型，並可設定排程自動化建模

## 資料清理

通常在資料分析的專案中，資料清理大概佔了80%的工作時間，因此DAS提供Data Refinery這個工具，減少程式的撰寫並且提供pipeline的模式進行清理流程的編輯



## 資料收集

提供豐富的資料串接工具與介面，包含MySQL、Oracle、DB2、S3、Dropbox等等資料庫與雲端儲存都可以彙整於DAS的環境中



## 資料治理

DAS除了提供良好的資料上傳介面外，同時也提供了高效率的資料管理介面，包含資料的使用權限、資料的健康度、資料的類別屬性都可以進行統籌管理。

# DAS在國網中心不同層面上的架構與角色



## 產業應用

可應用於金融業異常交易偵測、工業IOT製造異常分析、智慧交通流量預測等等AI與大資料應用場域，提升產業競爭力



## 分析大師

整合IBM與其他軟體工具，在雲端發揮資料收集、治理、分析、部署的價值，提供全功能且自動化的人工智慧整合工具



## 資料科學工具組

從資料收集、資料治理、資料清理、資料分析、模型分析部署都有相對應的軟體設置，可整合內外部的資料與相關人員，提高資料整合的價值

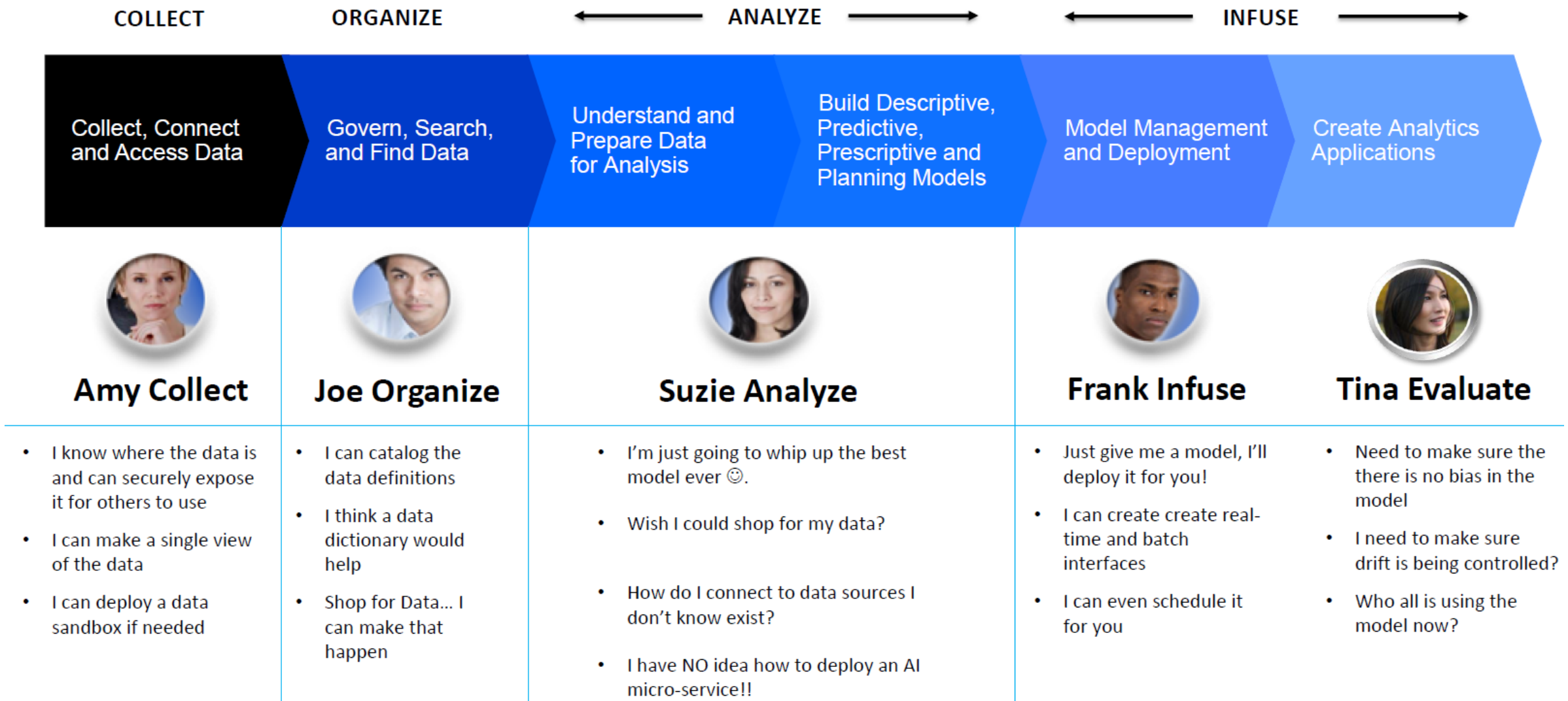
## 硬體層

完整串接 IaaS、PaaS，整合基礎計算設施、使用平台，自助設定資源；支援 Docker、Singularity 等容器技術，並透過 Kubernetes、Slurm 進行資源管理

# DAS在推廣上的目標與預期成果

目標	關鍵成果
贏得客戶滿意	<u>透過建置之DAS服務，客戶可以透過這個平台建立自己的資料科學團隊，有效的利用與分享資料和模型，同時完成模型建置與部署，加速公司內部的AI自動化能量</u>
提昇產業效益	AI產業化與產業AI化： <u>建立產業內部資料、模型的彙整環境，加強產業應用資料開發模型與部署的能量及效率</u>
鼓勵創新研發	藉由這個平台，中大型的公司可以彙整其內部的資料與AI量能，創造其內部的資料活化與增值應用
創造社會效益	中心藉由提供這個環境的過程，可以參與並且協助規劃相關AI與大數據產出，讓企業更關注資料分析與團隊合作
增加團隊合作	這項服務需要大量內部跨團隊的合作，同時提供服務後，企業本身也可以整合其內部資料科學的能量，企業與企業間也能夠有一個共有與共同合作的平台

# 組成內外部的DAS團隊





# 案例與功能介紹

## PM2.5分析預測

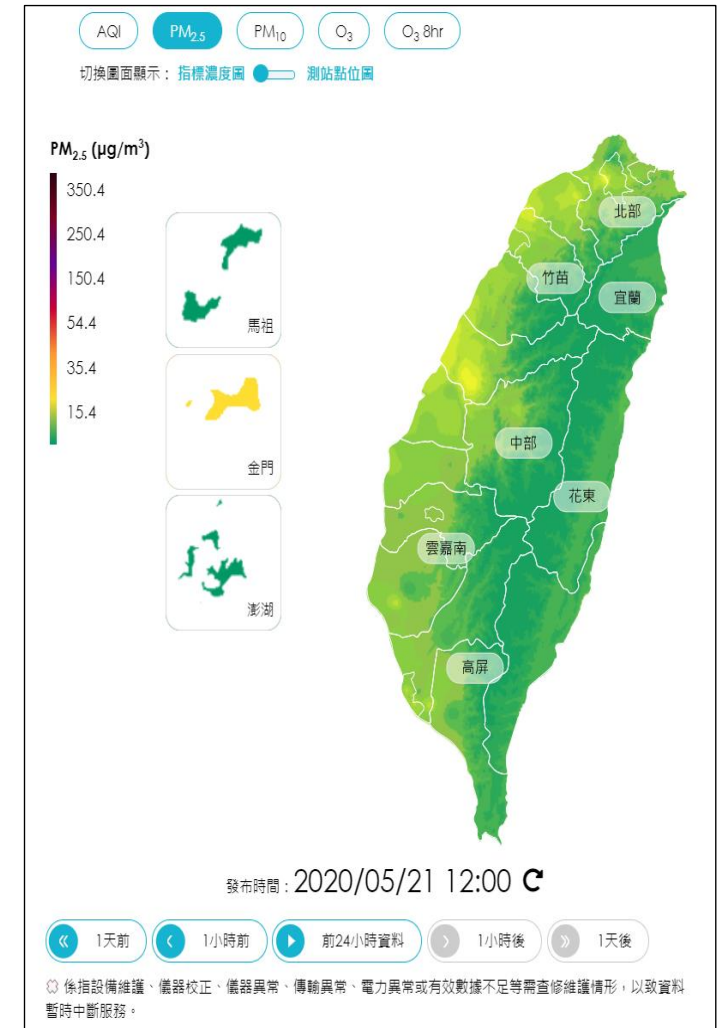
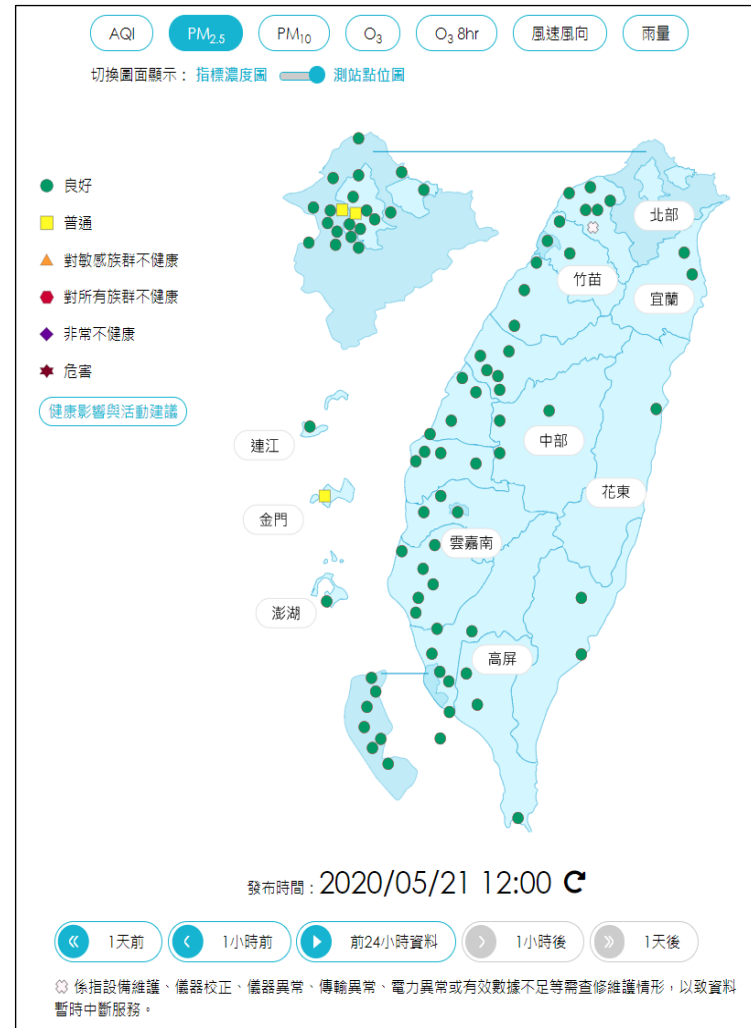
# 問題背景與分析需求

## • 展示背景：

1. 利用公開資料，對市民與政府機關提供PM2.5數值未來的預測值，提供查詢與決策。
2. 並對於潛在可能造成汙染的地點與時段(如：工廠、特定路段等)，進行監控並做即時告警。
3. PM2.5數值實時數值及未來五分鐘的預測值。

## • 需求：

依公開資料(氣象資料、雨量資料、PM2.5資料)一預測PM2.5數值在未來5分鐘的變化。

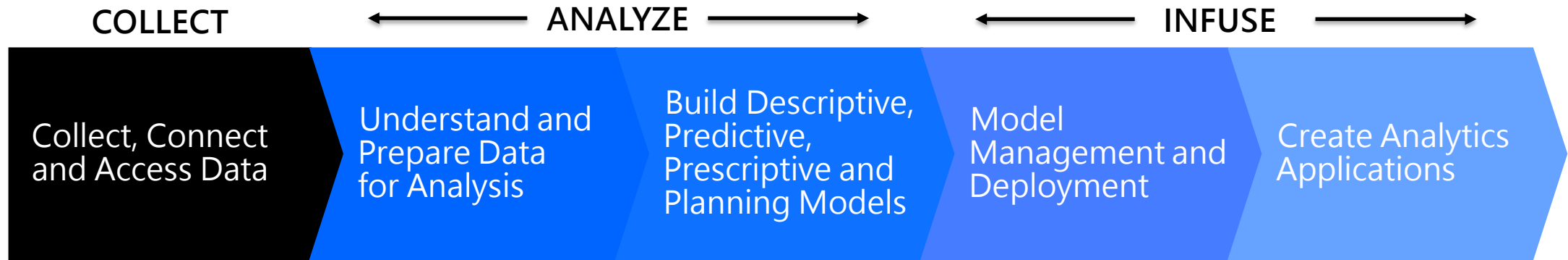


# 空汙與氣象資料來源

- 氣象資料：
  - 局屬氣象站 (共 45 站點)
  - 自動氣象站 (共 432 站點)
- 雨量資料：
  - 自動雨量站 (共 998 站點)
- PM2.5 資料：
  - IASS 社群 (共約 3000 點)



# 分析大師(DAS)Pipeline



- Data Stream:  
氣象資料  
雨量資料  
PM2.5資料

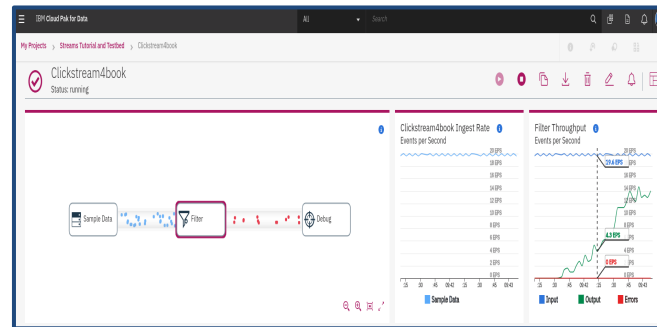
氣象資料：  
局屬自動氣象站  
(1h)  
氣象站(1h)

雨量資料：  
自動雨量站  
(10min)

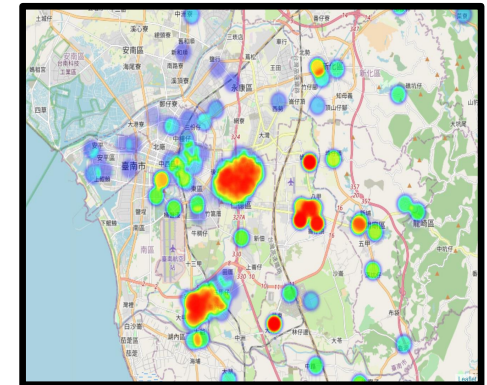
PM2.5 資料：  
lass 社群(5min)

- Data Stream:  
資料清理與融合

- Notebook:  
開發資料處理  
與週期性任務  
的建立

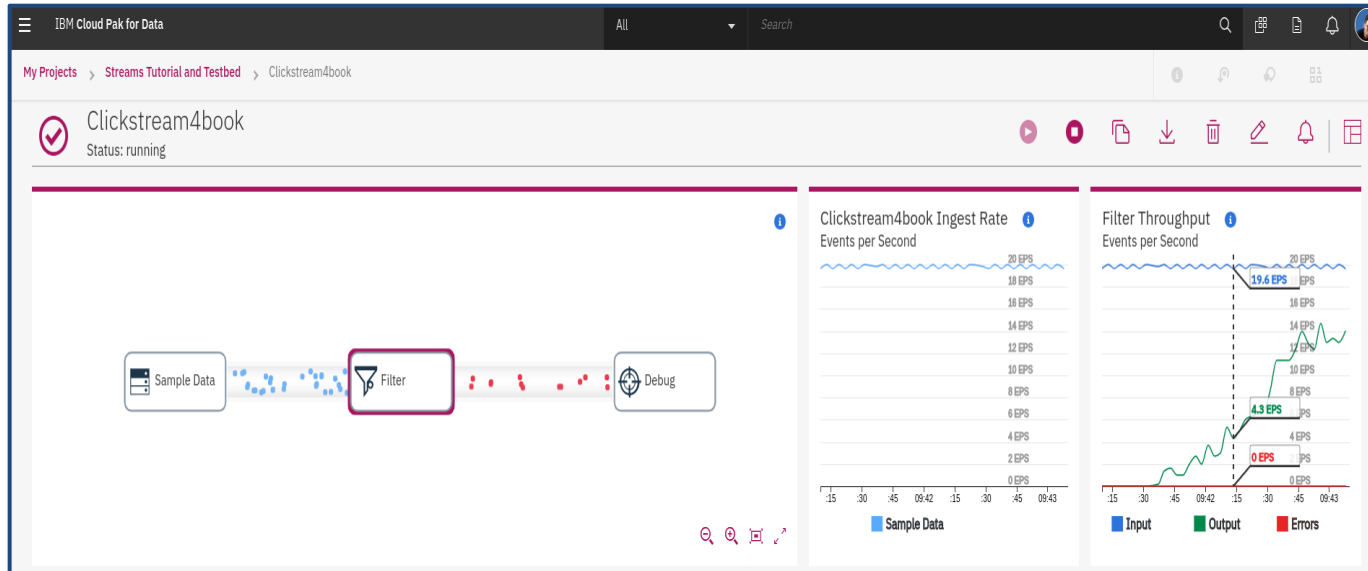


- Deployment:  
部署模型進行  
預測未來變化



- Notebook:  
透過程式繪製合適的空間資料視覺化呈現

# 以Streams功能進行即時資料分析



332: StreamsTutorialandTestbed::pm25\_332: insert\_view

Time	COMMON_ID	COMMON_OBSTIME	DEVICE_ID	HANDLETIME	HUMIDITY	LAT	Lon	PM2_5	RAIN_ID	RAIN_OBSTIME	STATION_ELEV	STATION_H_24R	STATION_HOUR_12	S'
5/20/2020 4:14:14 PM	COC700	2020-05-15 06:00:00.0	08BEACO0882	2020-05-2016:14:13.19	99	24.9505	121.219124	11	A2C560	2020-05-2016:00:00.0	151	0	18.5	11
5/20/2020 4:14:13 PM	467770	2020-05-15 11:00:00.0	74DA38EBF944	2020-05-2016:14:11.783	100	24.269	120.538	7	467770	2020-05-2016:00:00.0	7.2	0	0	2.
5/20/2020 4:14:11 PM	466880	2020-05-15 11:00:00.0	08BEACO0884	2020-05-2016:14:10.377	80	25.019415	121.40534	10	01C570	2020-05-2016:00:00.0	11	0	2	3.
5/20/2020 4:14:10 PM	C0F9T0	2020-05-15 10:00:00.0	74DA3895C2BC	2020-05-2016:14:08.933	100	24.186	120.666	10	C0F9T0	2020-05-2016:00:00.0	111	0	0	7.

IBM Cloud Pak for Data

My Projects > Streams Tutorial and Testbed > TWM-rain-Station-Streams

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | Python 3.6.0

### 即時雨量資料擷取與更新

In this notebook, you'll see rain data access by streams:

1. [Setup your Streams instance](#)
2. [Create the Rain Station application](#)
3. [Submit the application](#)
4. [Connect to the running application to view data](#)
5. [Stop the application](#)

#### 1. Setup

##### 1.1 Add credentials for the IBM Streams service

In order to submit a Streams application you need to provide the name of the Streams instance.

1. From the navigation menu, click **My instances**.
2. Click the **Provisioned Instances** tab.
3. Update the value of `streams_instance_name` in the cell below according to your Streams instance name.

```
In [1]: # from icpd_core import icpd_util
streams_instance_name = "cp4d-streams-instance" ## Change this to Streams instance
cfg=icpd_util.get_service_instance_details(name=streams_instance_name)
cfg
```

##### 1.2 Import the streamsx package and verify the package version

```
In [2]: import streamsx.topology.context
from streamsx.topology.schema import StreamSchema
from streamsx.topology.topology import *
from streamsx.topology.context import *
from streamsx.topology.schema import StreamSchema
import streamsx.database as db
print("INFO: streamsx package version: " + streamsx.topology.context.__version__)

#For more details uncomment line below.
#!pip show streamsx

INFO: streamsx package version: 1.13.14
```

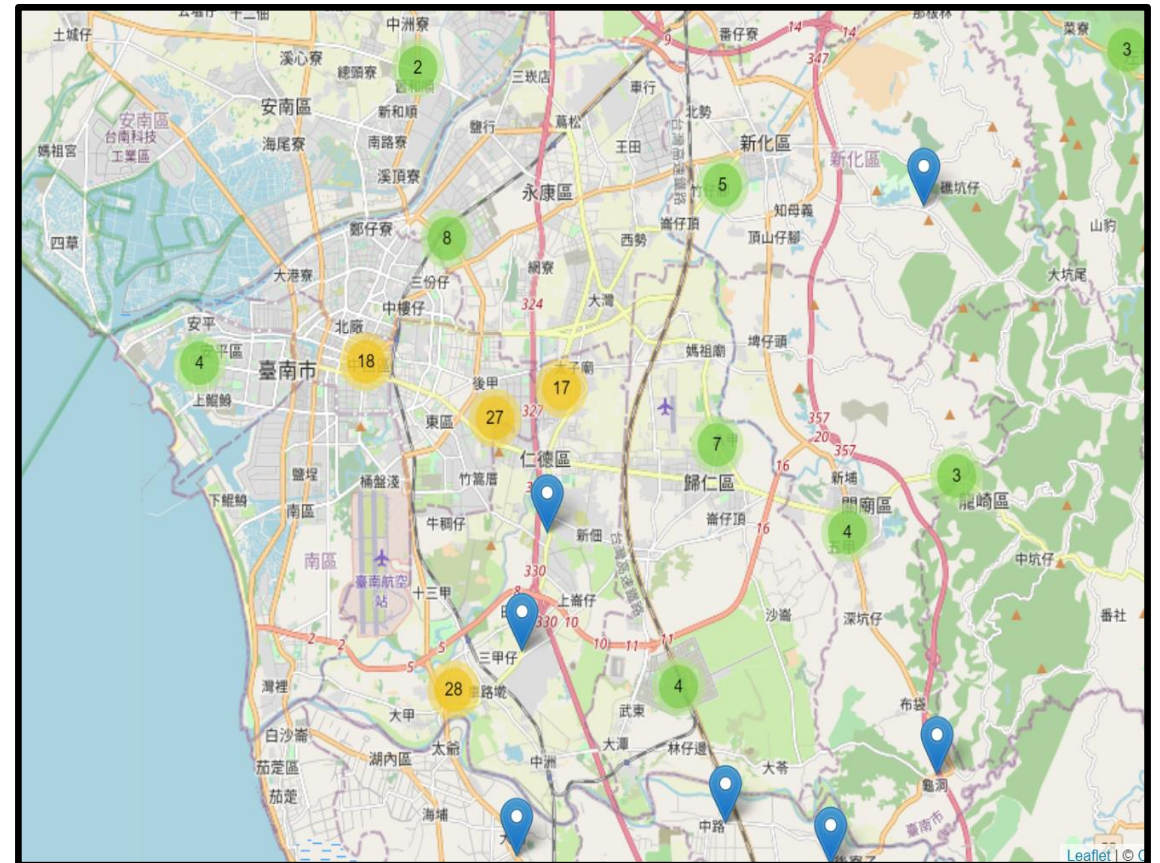
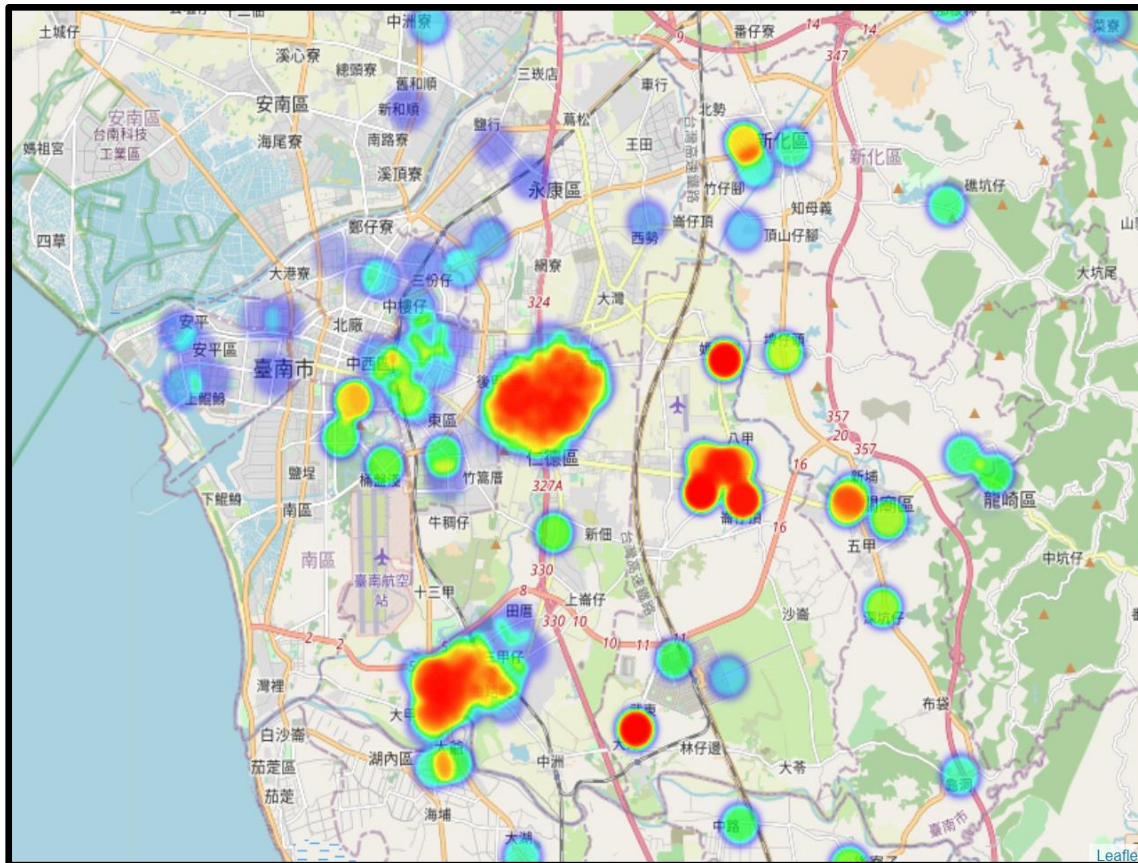
#### 2. Create the application

```
In [3]: #create topology...change the name to something more meaningful
from streamsx.topology.topology import Topology
import streamsx.topology.context

topo = Topology(name="rainStation")
```

# PM2.5 資料預測結果:以台南市為例

- 將分析結果以RStudio繪圖並部署到資料市集，可針對資料坐即時的監控



# 案例與功能介紹

## Youbike借用應用

# 問題分析與需求

- Youbike公司背景：

1. Youbike公司在借用服務對市民提供即時各站目前剩餘車輛數及可還車空位數供市民查詢。
2. 對於熱門站點於特殊時段(如：公館捷運站出口放學時段需求量大增)需求只能依經驗作調度。
3. 調度量及各站間調度面臨困難。

- 需求：

依各站借用歷史資料，結合氣象資料(氣溫、降雨)發展一預測模式作為車輛調度使用。





# 資料內容與結構

- Youbike:

- sno : 站點代號
- tot : 場站總停車格
- sbi : 場站目前車輛數量
- bemp : 空位數量
- act : 全站禁用狀態
- utime: 系統時間

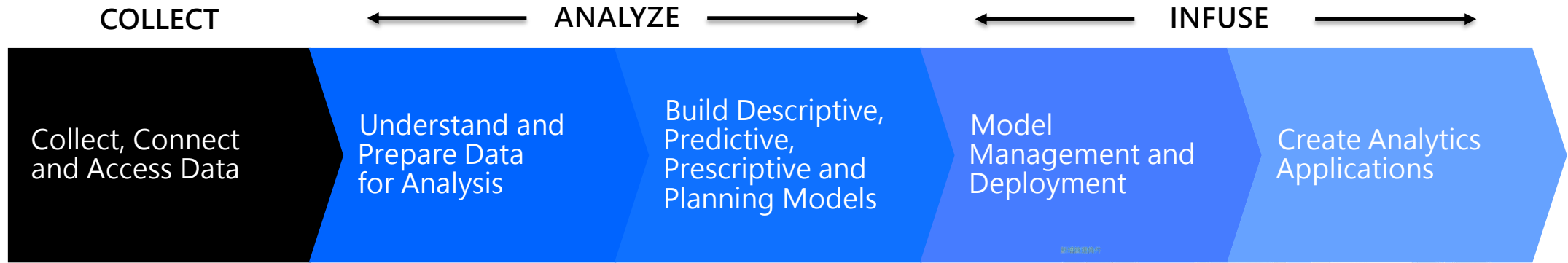
- 天氣:

elementName	中文說明
ELE	高度，單位 公尺
WDIR	風向，單位 度，一般風向 0 表示無風
WDSD	風速，單位 公尺/秒
TEMP	溫度，單位 攝氏
HUMD	相對濕度，單位 百分比率，此處以實數 0-1.0 記錄
PRES	測站氣壓，單位 百帕
H_24R	日累積雨量，單位 毫米
H_FX	小時最大陣風風速，單位 公尺/秒
H_XD	小時最大陣風風向，單位 度
H_FXT	小時最大陣風時間，yyyy-MM-ddThh:mm:ss+08:00
D_TX	本日最高溫，單位 攝氏
D_TXT	本日最高溫發生時間，hhmm (小時分鐘)
D_TN	本日最低溫，單位 攝氏
D_TNT	本日最低溫發生時間，hhmm (小時分鐘)
CITY	縣市
CITY_SN	縣市編號
TOWN	鄉鎮
TOWN_SN	鄉鎮編號
補充說明	-99 皆表示 該時刻因故無資料。

<https://tcgbusfs.blob.core.windows.net/lobyoubike/YouBikeTP.json>

<https://opendata.cwb.gov.tw/api/v1/rest/datastore/O-A0001-001?Authorization=rdec-key-123-45678-011121314>

# 分析大師(DAS)Pipeline



- Data Connection:
  - Youbike資料
  - 氣象觀測資料
  - 氣象預報資料

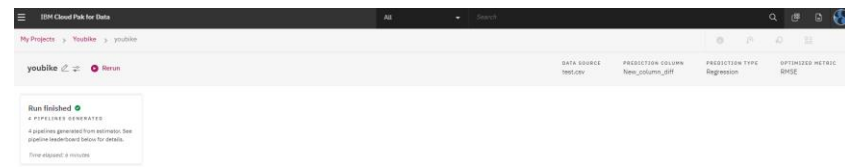
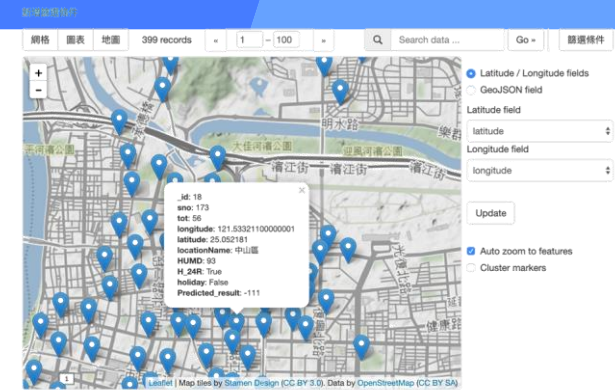
- Data Refinery:
  - 計算每天車輛借用情況
  - 清理氣象資料

- AutoAI:
  - 利用現有借用資料及氣象觀測資料
  - 透過AutoAI訓練預測模型

- Deployment:
  - 部署預測模型進行預測

- Dashboard:
  - 互動式監控面板

- Data sharing:
  - 資料市集、open data

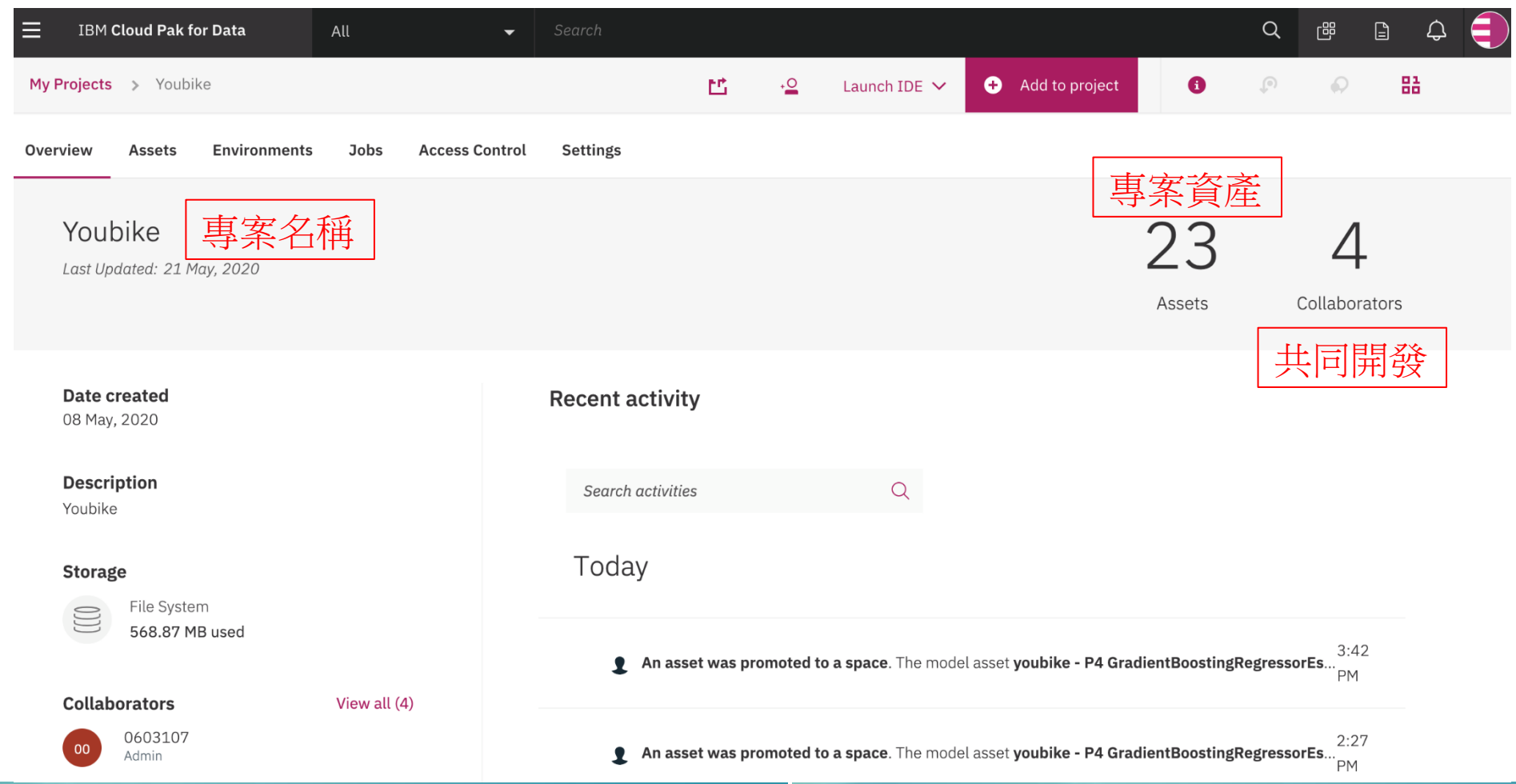


Pipeline leaderboard

Rank	Name	Estimator	RMSE	Enhancements	Build time
1	Pipeline 2	Random forest regressor	46.405	(HPO-1) (FE)	00:02:26
2	Pipeline 4	Random forest regressor	46.405	(HPO-1) (FE) (HPO-2)	00:02:14
3	Pipeline 1	Random forest regressor	49.122	None	00:00:02

# 專案管理清單

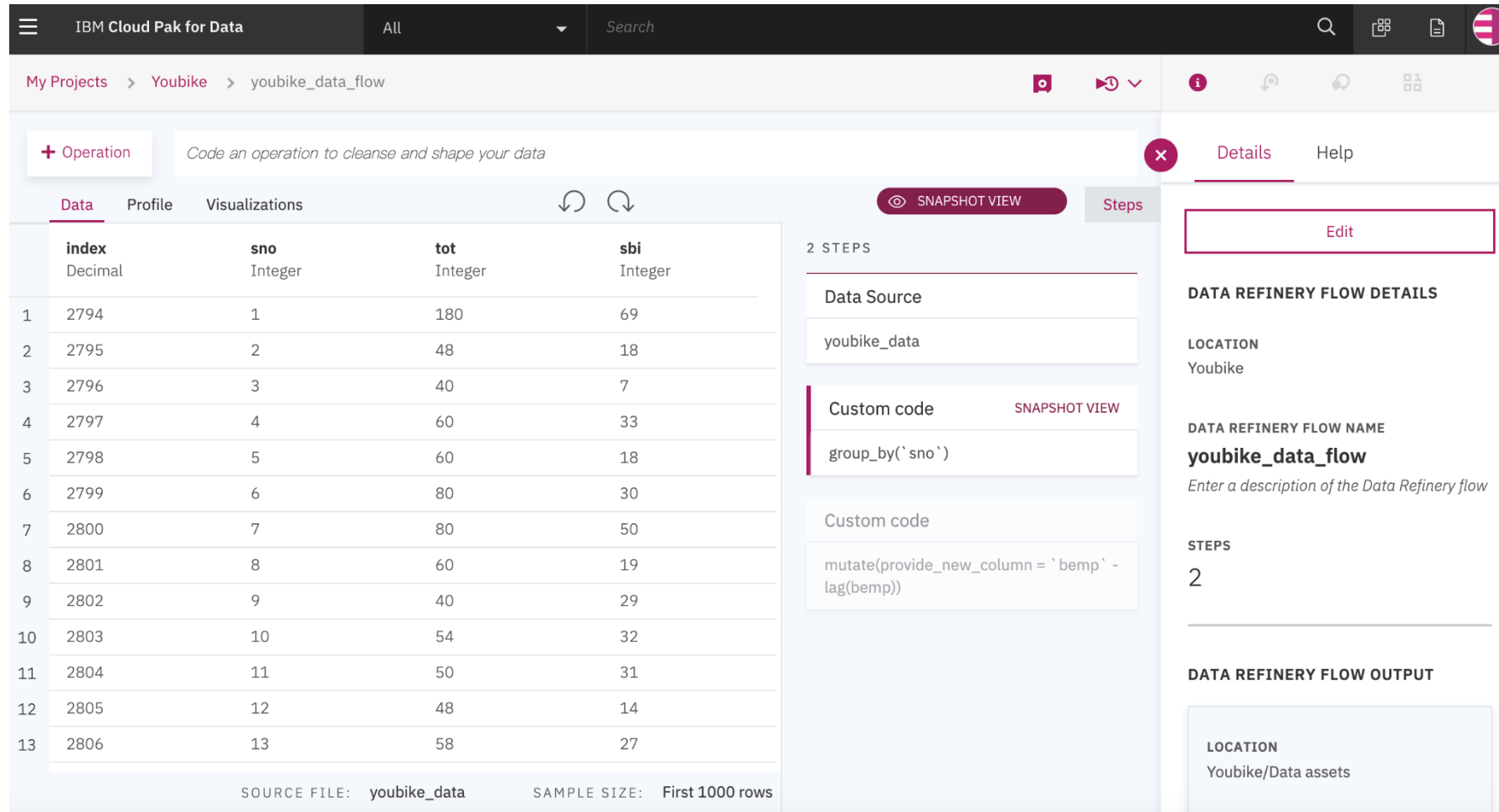
- 建立專案並加入共同開發者



The screenshot shows the IBM Cloud Pak for Data interface for a project named 'Youbike'. The top navigation bar includes 'IBM Cloud Pak for Data', 'All', and a search bar. Below the navigation, there are tabs for 'Overview', 'Assets', 'Environments', 'Jobs', 'Access Control', and 'Settings'. The main content area displays the project name 'Youbike' with a red box around it labeled '專案名稱'. To the right, it shows '23 Assets' and '4 Collaborators', with a red box around the '4 Collaborators' labeled '共同開發' and another red box around the '23 Assets' labeled '專案資產'. Below this, there are sections for 'Date created' (08 May, 2020), 'Description' (Youbike), 'Storage' (File System, 568.87 MB used), and 'Collaborators' (0603107 Admin, with a 'View all (4)' link). The 'Recent activity' section shows two entries: 'An asset was promoted to a space. The model asset youbike - P4 GradientBoostingRegressorEs...' at 3:42 PM and another similar entry at 2:27 PM.

# Data Refinery

- 將資料導入並使用data refinery進行資料清理



IBM Cloud Pak for Data | All | Search

My Projects > Youbike > youbike\_data\_flow

+ Operation *Code an operation to cleanse and shape your data*

Details Help

	index	sno	tot	sbi
	Decimal	Integer	Integer	Integer
1	2794	1	180	69
2	2795	2	48	18
3	2796	3	40	7
4	2797	4	60	33
5	2798	5	60	18
6	2799	6	80	30
7	2800	7	80	50
8	2801	8	60	19
9	2802	9	40	29
10	2803	10	54	32
11	2804	11	50	31
12	2805	12	48	14
13	2806	13	58	27

SOURCE FILE: youbike\_data | SAMPLE SIZE: First 1000 rows

2 STEPS

Data Source: youbike\_data

Custom code: `group_by(`sno`)` (SNAPSHOT VIEW)

Custom code: `mutate(provide_new_column = `bemp` - lag(bemp))`

DATA REFINERY FLOW DETAILS

LOCATION: Youbike

DATA REFINERY FLOW NAME: **youbike\_data\_flow**  
*Enter a description of the Data Refinery flow*

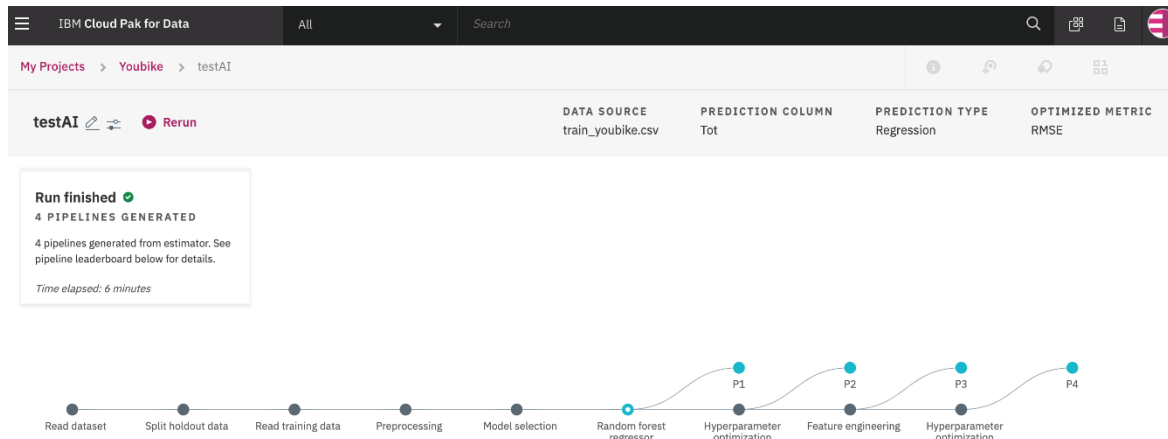
STEPS: 2

DATA REFINERY FLOW OUTPUT

LOCATION: Youbike/Data assets

# AutoAI

- 將清理好的資料導入 AutoAI
- 自動選取回歸模式
- 由系統自動選取最佳的模型，並可觀察模型分析細節



Pipeline leaderboard Compare pipelines Ranking based on: Root Mean Squared Error (RMSE) (Optimized) ▾

Rank	Name	Estimator	RMSE	Enhancements	Build time
> ★ 1	Pipeline 3	Random forest regressor	0.540	HPO-1 FE	00:02:17
> 2	Pipeline 4	Random forest regressor	0.540	HPO-1 FE HPO-2	00:03:10
> 3	Pipeline 1	Random forest regressor	1.267	None	00:00:02
> 4	Pipeline 2	Random forest regressor	1.267	HPO-1	00:00:30

IBM Cloud Pak for Data | All | Search

My Projects > Youbike > testAI

← Back to testAI

RANK 1 Pipeline 3 ▾ RMSE (Optimized) 0.54 ESTIMATOR Random forest regressor ENHANCEMENTS HPO-1 FE BUILD TIME 00:02:17 Save as model

### RandomForestRegressor Model Evaluation Measures ⓘ

TARGET : TOT

EVALUATION

- Model Evaluation Measures
- MODEL VIEWER
- Model Information
- Feature Transformations
- Feature Importance

	Holdout Score	Cross Validation Score
Root Mean Squared Error (RMSE)	-0.271	-0.540
R <sup>2</sup>	1.000	0.998
Explained Variance	1.000	0.998
Mean Squared Error (MSE)	-0.073	-0.294
Mean Squared Log Error (MSLE)	-0.000	-0.000



Thank you for your attention.

分析大師 DAS 試用申請: <https://event.nchc.org.tw/2020/dasapply/>  
聯絡我們: [das@narlabs.org.tw](mailto:das@narlabs.org.tw)