

以Machine Learning技術提升現有 資訊系統智能服務

Wayne Tien

2018.10.12

GSS 勸揚資訊
Galaxy Software Services

國家產業創新獎
卓越中堅企業



2018

李開復：5 秒準則

如果一項工作需要的思考決策，能在**5秒鐘**內做出決定，就有很大的可能被人工智慧技術全部或部分取代

AI是協助我們更高效率地完成手上的某些任務

省下的精力時間去作高價值的事情

大綱

Table Of Content

- 一、機器學習
- 二、案例分享
- 三、續章



國家產業創新獎
卓越中堅企業

Chapter

1

機器學習 (Machine Learning)

機器學習二三事

- 機器學習程序
- 監督式機器學習
(Supervised Learning)
- 非監督式機器學習
(Unsupervised Learning)

機器學習 (Machine Learning)

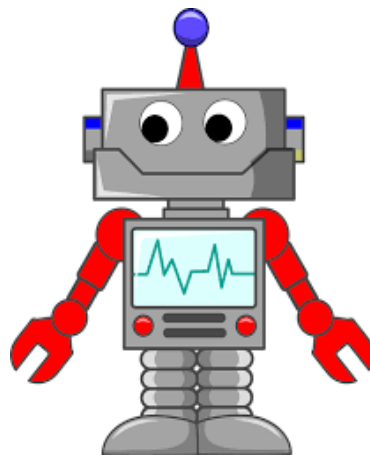
透過從過往的資料
和經驗中學習
並找到其運行規則
最後達到人工智慧
的方法

機器學習包含透過樣本資料來訓練機器辨識出**運作模式(演算法)**，而不是採用特定的**規則(Rule)**。

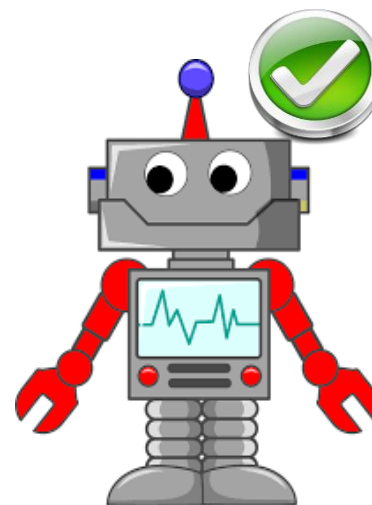
換句話說，機器學習是一種弱人工智慧(narrow AI)，它從資料中得到複雜的函數(或樣本)來學習以創造演算法(或一組規則)，並利用它來做預測。



華倫·巴菲特
Warren Buffett



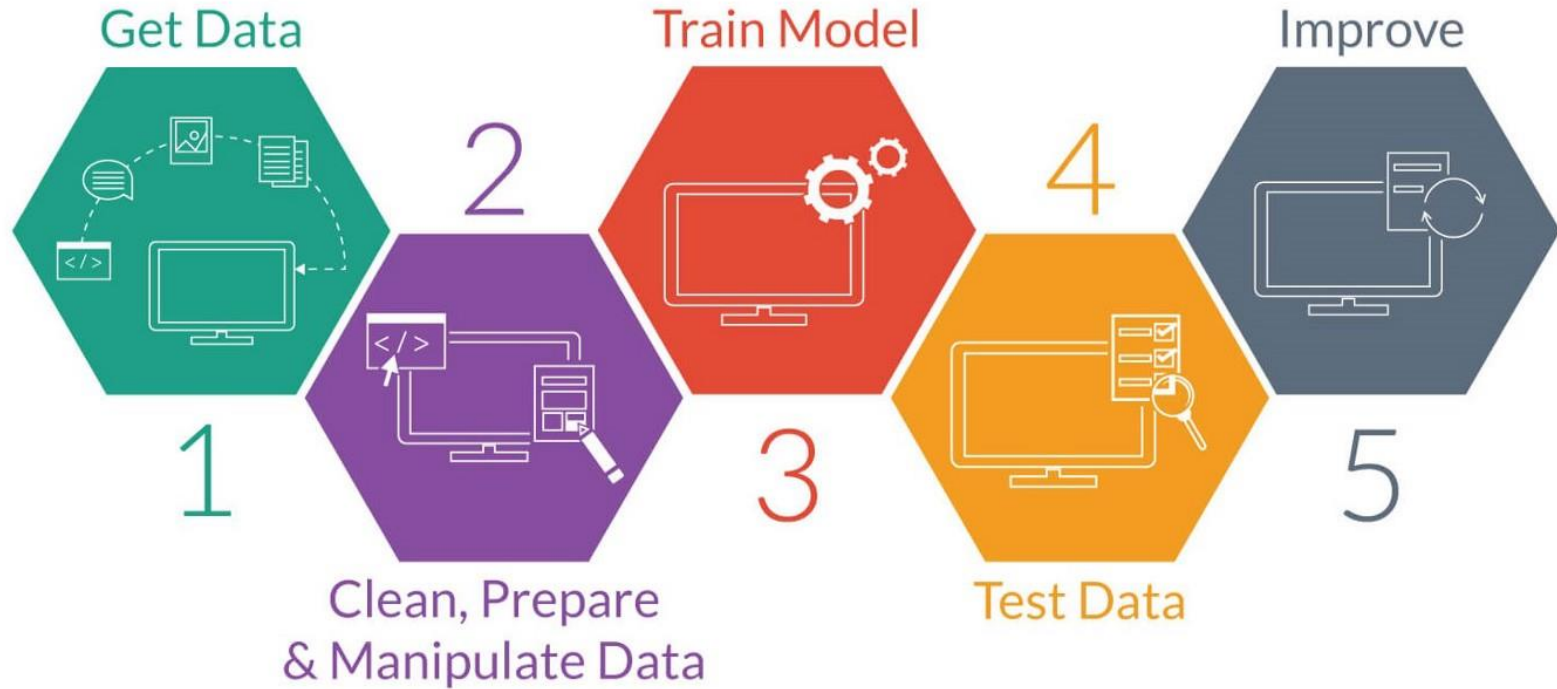
Buffett
選股機器人



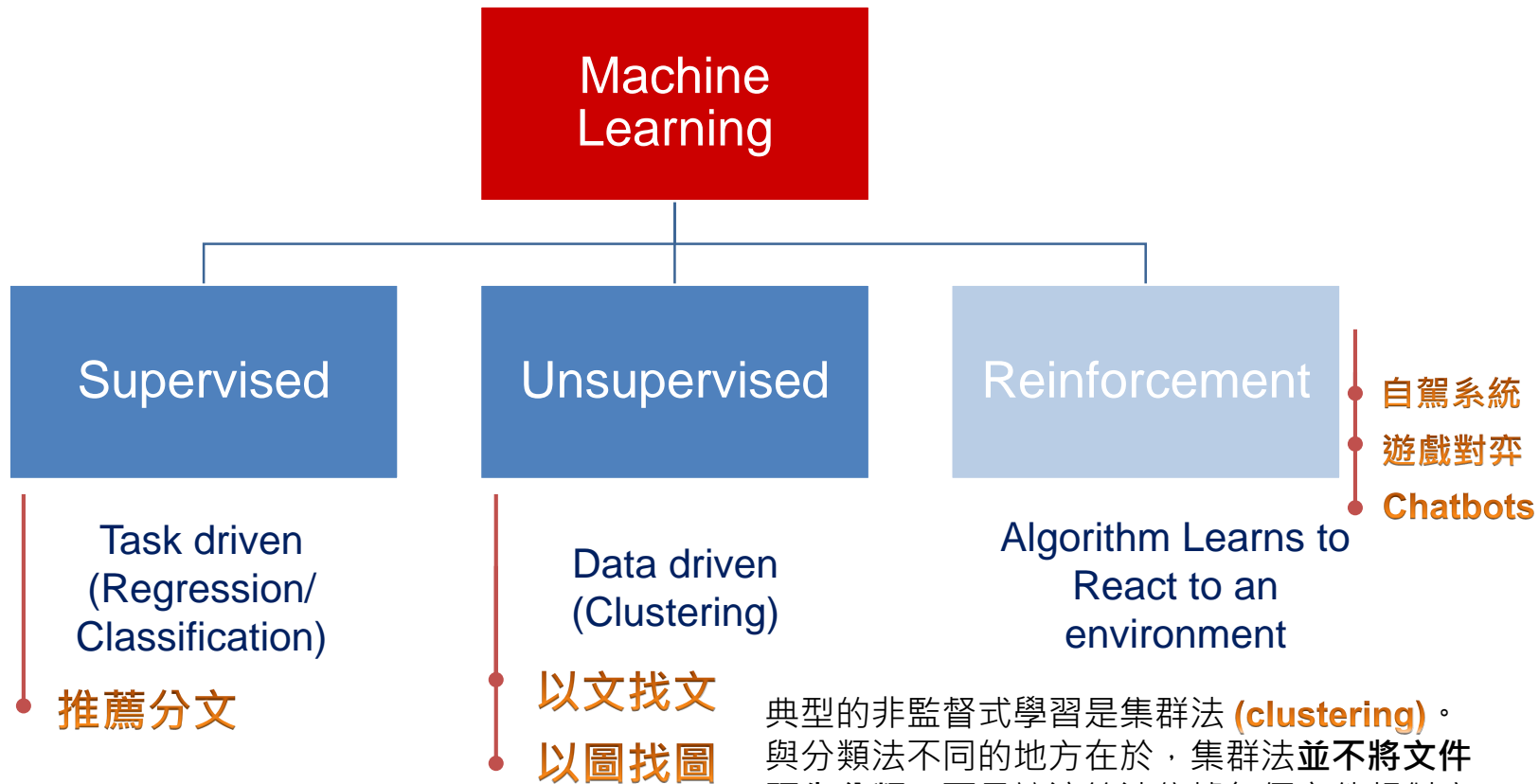
優質
選股機器人

機器學習程序

Typical Machine Learning Workflow



Types of Machine Learning



監督式學習使用經標記之訓練集 (**labeled training dataset**) 來建構系統模型，具有預先定義好之輸入與已知之輸出。

在監督式學習中，典型的任務是**文件分類**和**迴歸分析**。

典型的非監督式學習是**集群法 (clustering)**。與分類法不同的地方在於，**集群法並不將文件預先分類**，而是讓演算法依據每個文件相似度與相異程度，將文件區分成不同之群集。

關鍵在於讓同一群集內文件間的差異盡量小，彼此相似程度高；至於不同群集間的文件則有大差異，從而可將每一文件歸屬至一個特定群集，對於文件檢索之效率與準確度之提升有極大助益。

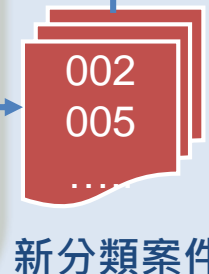
監督式機器學習 (Supervised Learning)

Step1 : Train



再訓練

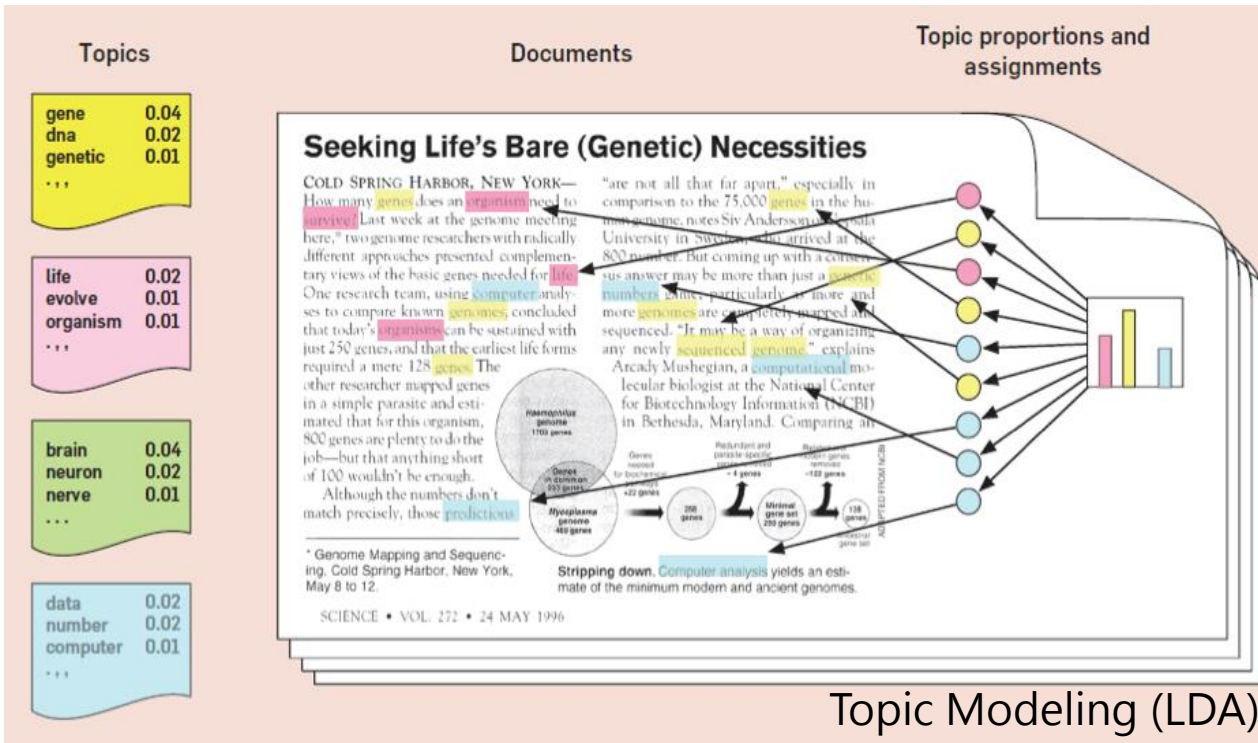
回饋 校正



Step2 : Predict

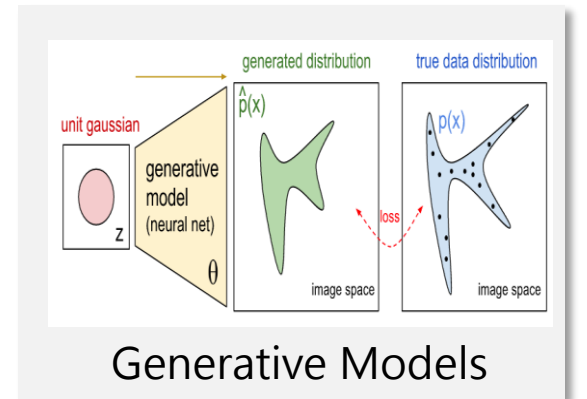
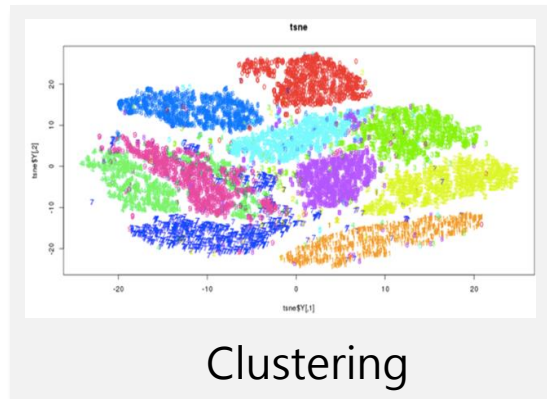
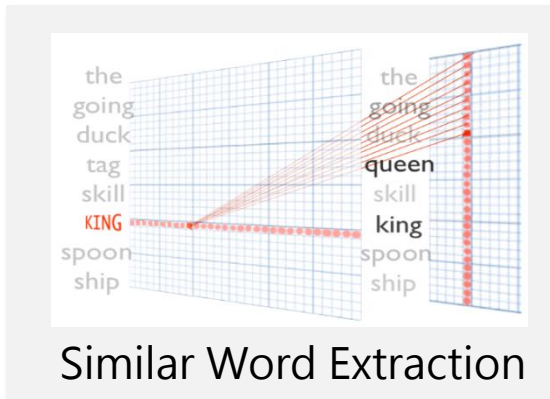
Analytic Engine

非監督式機器學習(Unsupervised Learning)



“Arts”	“Budgets”	“Children”	“Education”
NEW	MILLION	CHILDREN	SCHOOL
FILM	TAX	WOMEN	STUDENTS
SHOW	PROGRAM	PEOPLE	SCHOOLS
MUSIC	BUDGET	CHILD	EDUCATION
MOVIE	BILLION	YEARS	TEACHERS
PLAY	FEDERAL	FAMILIES	HIGH
MUSICAL	YEAR	WORK	PUBLIC
BEST	SPENDING	PARENTS	TEACHER
ACTOR	NEW	SAYS	BENNETT
FIRST	STATE	FAMILY	MANIGAT
YORK	PLAN	WELFARE	NAMPHY
OPERA	MONEY	MEN	STATE
THEATER	PROGRAMS	PERCENT	PRESIDENT
ACTRESS	GOVERNMENT	CARE	ELEMENTARY
LOVE	CONGRESS	LIFE	HAITI

Topic Modeling (LDA)



Chapter

2

案例分享

智慧局產業專利知識平台

- 推薦分類
- 以文找文
- 專同辭庫

UX
Cloud services
AI Bot DevOps
Big data SaaS
IOT Blockchain

智慧局產業專利知識平台實際應用



智慧局產業專利知識平台

專利自動推薦
對於有興趣的專利
可進行收藏

個人檢索/瀏覽歷程
以文找文功能

關鍵字排行

- 驅動系統
- Development
- service
- design
- 風扇馬達
- 有機化合物
- 滾輪乾燥機

標籤排行

- #玻璃砂
- #模造系統
- #破原子有機化學
- #高

瀏覽及收藏排行

追蹤訂閱快速
瀏覽區

智慧局產業專利知識平台

關鍵字檢索

產業別分類主題 檢索及瀏覽

產業專利知識平台

机械 * 冲床

工具機 | 化學 | 電腦 | 光學 | 其他消費品

- 剪床 | 搪床 | 綜合加工機
- 彎曲機 | 攻牙機 | 車床
- 折床 | 放電加工機 | 銑床
- 拉床 | 冲床 | 鑽床
- 插床 | 磨床 | 齒輪加工機

生物科技 | 電機 | 機械

製藥 | 半導體 | 運輸

電子 | 儀器 | 土木工程

通訊 | 醫學技術 | 家具/遊戲

搜尋結果共有 94,585 筆, 花費 0.36 秒

依相關度排序 | 資料範圍

國際分類號
A01G (4669)
B21D (3969)
B23Q (3301)
B29C (2949)
C08L (2906)
MORE (537) v

申請人
CORNING INCORPORATED (67)
3M創新有限公司 (57)
APPLIED MATERIALS, INC. (49)
DOW GLOBAL TECHNOLOGIES LLC (49)
HON HAI PRECISION INDUSTRY CO., LTD. (47)
MORE (995) v

發明人

伺服冲床 監控方法及其裝置

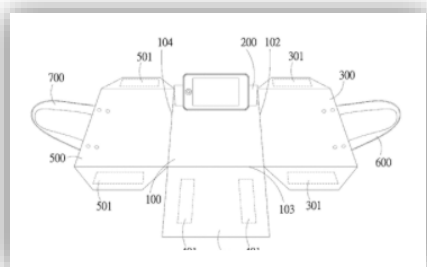
本發明伺服冲床監控方法，包括：提供伺服冲床所要執行的內建於該伺服冲床之理想成形曲線的訊號，驅動伺服馬達運轉以產生扭矩來驅動樞接於該伺服冲床的動應於理想成形曲線的作動；藉由曲軸角度檢出單元接收動力機構的曲軸角度，取得動力機構的實際位置參數，以扭矩回饋單元來接收成形模具中所產生的成形力；藉由出力判斷單元將實際位置參數與成形力計算出實際成形位置，並依據實際成形位置與伺服馬達的運轉速度計算出實際成形曲線，比較理想成形曲線與實際成形曲線，判斷動力機構的運作是否異常。

陳志遠·傅東洪·謝志鴻·林瑞富·林俞廷 - 羅昇企業股份有限公司
公開日：2016-05-01
專利局：TW

冲床之同步伺服送料系統及其運作方法

專利推薦分類 - Why ?

專利發明



多功能摺疊餐袋

專利細分69,000種分類



專利分類

階層	類號	編排	總數	範例	說明	
一	Section	A~H	部	8個	A	人類生活需要
二	Class	二位數	主類	128個	45	手攜物品或旅行品
三	Subclass	一個大寫的英文字母	次類	628個	C	小包；行李箱；手提袋
四	Main group	1至3位數加/00	主目	69,000個	011	午餐或野餐盒或類似物
五	Sub group	將"/"後的00改為2~4個數字	次目		/20	

專利推薦分類文件特性

中華民國第 I453740 發明專利

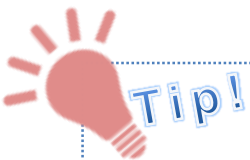
專利名稱：

「光學儲存媒體和從光學儲存媒體閱讀 資料之方法，以及包括閱讀光學儲存媒體上 資料所用的光學拾波器之裝置」。

摘要：

「本案光學儲存媒體包括基片層(2)，和 設置在基片層(2)上之資料層(3)，資料層(3)包 括資料，配置在磁軌內，成為標記和空格。保護碼(PC)寫碼於一磁軌或磁軌之一部份，該碼包括第一尺寸(w1)之標記(P)和較小的第二尺寸(w2)之標記(P)。第二尺寸(w2)之標記尤 其是寬度，較第一尺寸(w1)的標記之寬度 小。欲得保護碼(PC)，使用此方法包括步驟 為，以第一雷射功率閱讀磁軌或磁軌(T)之一 部份，得第一資料訊號(D1)，在另一步驟以 不同於第一功率之第二功率，閱讀同樣磁軌 或磁軌(T)之同樣部份，得第二資料訊號 (D2)，並考量第一和第二資料訊號，計算保 護碼。保護碼之計算，可例如利用各資料閱 讀裝置之微處理器進行。」。

本件專利係關於DVD之專利，然無論「發明名稱」或「摘要」皆未出現「DVD」、「CD」、或「CD-R」等用語。



Tip!

不同知識領域資料內容特性多有不同，有其特殊詞彙或是專有名詞

針對不同知識領域會採取不同的對策，並經多過次模擬去決定策略及演算手法

沒有一種模型 (model) 可直接套用於各個知識領域

專利推薦分類流程



推薦分文機器學習(Supervised)流程一

資料搜集 整備解析

- 讀取資料
- 取得資料分類資訊

分詞/斷詞

- 使用 **Jieba** 搭配**對照詞庫**，對文件進行斷詞
- Jieba 原理
 - 產生句子所有斷詞可能的組合
 - 計算出最有可能的組合
 - 未知詞採用 HMM 模型

Stop Words 移除

- 把**常出現**，且**沒有意義**的詞去掉，例如「你」、「我」、「的」等等

特徵化 向量轉換

- 要將文件轉換成**數字向量**，所採用的方法稱為 Bag of Words
- 將文件表達為一個向量，每個文件的向量長度是一樣的，向量維度代表出現在文件的每一個詞，數字大小即這個詞在文件中出現的次數

推薦分文機器學習(Supervised)流程二

特徵值演算 向量權重

- 套用 **TF-IDF 演算** 對詞的向量權重進行調整

向量維度 調整

- 套用 **Word2Vec 演算**，將特徵化的**向量維度降低**，來加快訓練的速度

各分類 獨立訓練

- 訓練前會將分類的資料標示出來，讓**分類器**去學習如何區分把標的分類及其他分類
- 為避免多次訓練後會有**記憶體不足**的問題，每個分類訓練都會開啟新的 process，以確保訓練結束後記憶體被完全釋放。訓練完成後，即將該分類的模型儲存起來

預測

- 要對新文件進行分類預測時，由各分類 model 來計算出屬於該分類的機率，取機率最高的前三名作為預測答案

以文找文



使用者的想法

「尋找資料花太多時間，想快點找到相關的資訊」

「不太清楚要下什麼關鍵字，才能找到我要的東西」

「也許給我一些相關(可能沒看過或一時沒想到)的內容也不錯」

以文找文實作 (輸入)

輸入一段文字
或一篇文件

文字
前處理

演算分析

主題
關聯延伸

相似度
權重排列

檢索
結果呈現

根據本發明一實施例，提供一種半導體晶片封裝構件，包含有一基板，具有一晶片安裝面；複數個焊接墊，設於該晶片安裝面上；一第一虛設接墊以及一與該第一虛設接墊間隔開的第二虛設接墊，設於該晶片安裝面上；一防焊遮罩，設於該晶片安裝面上，並部分覆蓋該焊接墊、該第一虛設接墊與該第二虛設接墊；一晶片封裝，安裝在該晶片安裝面上，並透過設於該複數個焊接墊上的複數個錫球電連接該基板；一分立元件，設於該晶片封裝與該基板之間，該分立元件具有一第一連接端與一第二連接端；一第一焊錫，將該第一連接端、該第一虛設接墊與該晶片封裝連接起來；以及一第二焊錫，將該第二連接端、該第二虛設接墊與該晶片封裝連接起來。其中該分立元件包含一基板側電容、一去耦電容、一電阻或一電感。

檢索

取消

半導體 晶片 封裝 構件 基板 晶片 安裝 面 複數 焊接 墊 設於 晶片 安裝 面上 虛設 接墊 需 設 接 墊 間 隔 開 虛 設 接 墊 設 於 該 晶 片 安 裝 一 面 上 防 焊 遮 罩 設 於 晶 片 安 裝 面 上 部 分 覆 蓋 焊 接 墊 虛 設 接 墊 虛 設 接 墊 晶 片 封 裝 安 裝 在 晶 片 安 裝 面 上 透 過 設 於 複 數 焊 接 墊 上 複 數 錫 球 電 連 接 基 板 分 立 元 件 設 於 該 晶 片 封 裝 基 板 之 間 分 立 元 件 連 接 連 接 端 焊 錫 將 連 接 端 虛 設 接 墊 晶 片 封 裝 連 接 焊 錫 連 接 端 虛 設 接 墊 晶 片 封 裝 連 接 分 立 元 件 基 板 側 電 容 去 耦 電 容 電 阻 電 感

以文找文實作 (分析學習)

輸入一段文字
或一篇文件

文字
前處理

演算分析

主題
關聯延伸

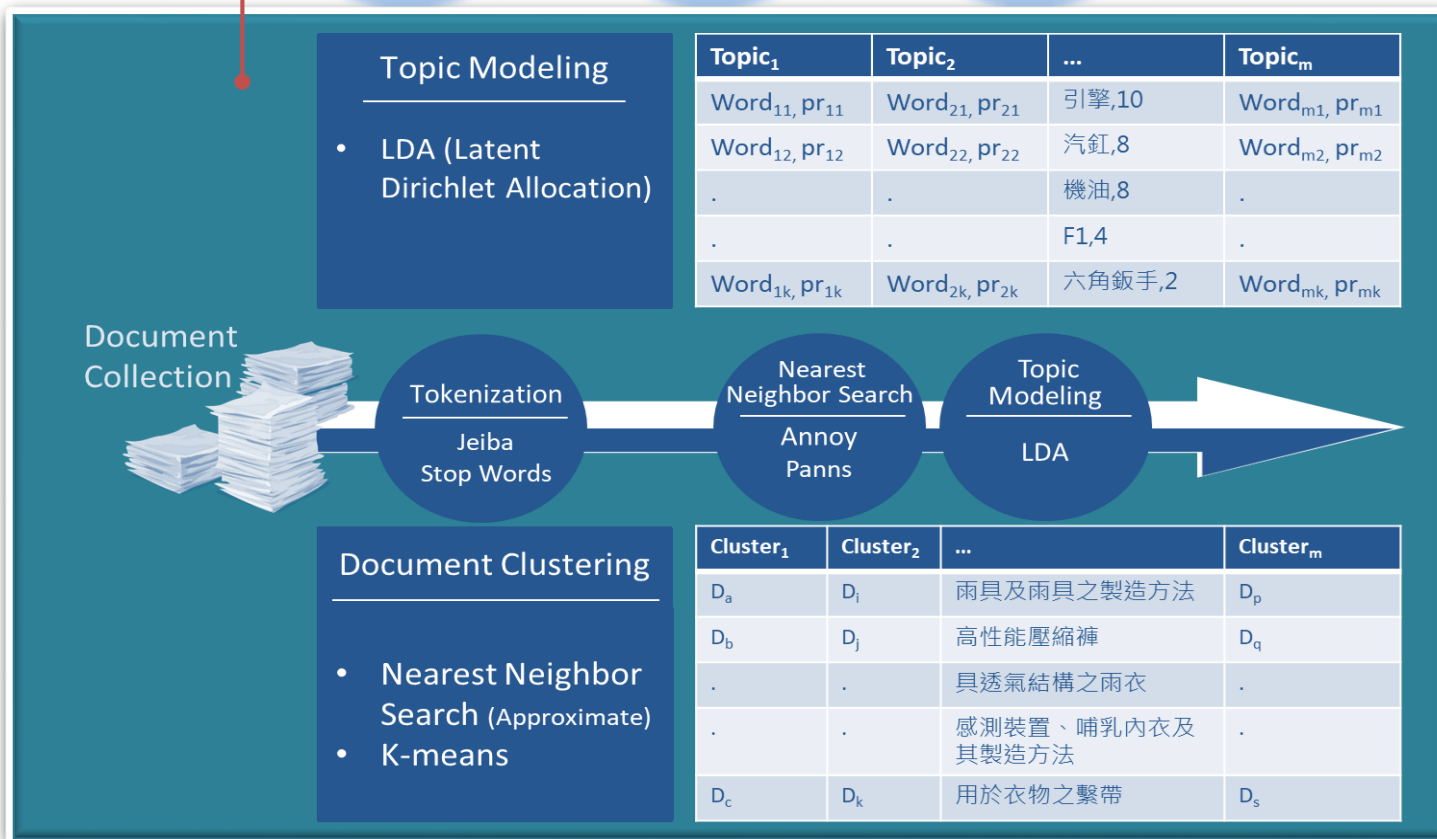
相似度
權重排列

檢索
結果呈現

A 文件特徵向量化

B 文件特徵索引

C 最鄰近搜尋



多國專利文件轉置成果

106年1月 ~ 106年7月
發明專利公開公報轉置資料數



193,640



88,907



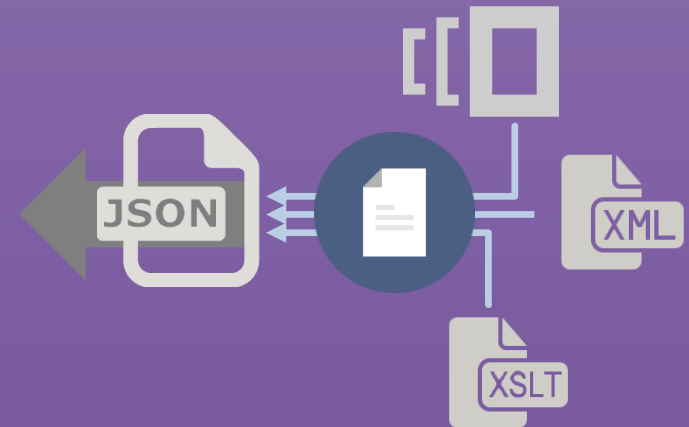
104,195



61,484



575,080



技術辭庫學習流程

辭庫分析

詞庫整理

- 採用專利技術名詞 **中英對照詞庫**
- 刪除不必要之詞彙及片段
- 另加入 **Jieba** 工具簡體中文詞量

分詞與詞性標註的媒合

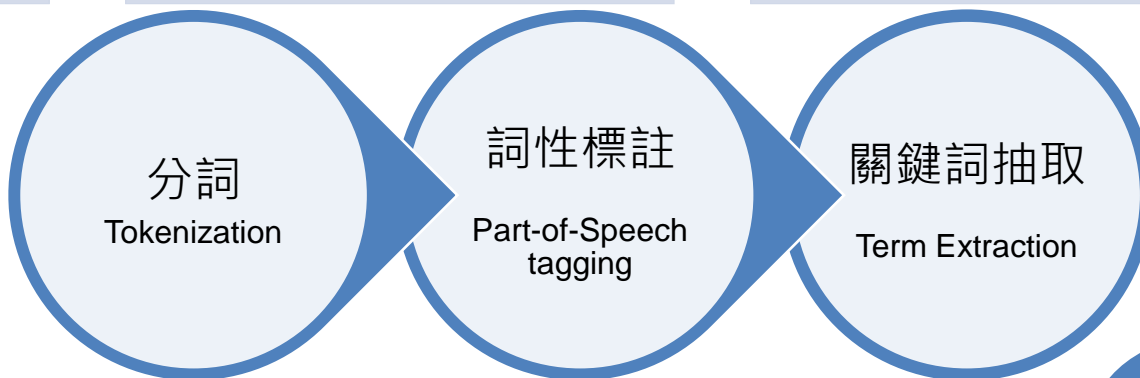
- **詞性標註** 是指對於句子中的每個詞指派適合的詞性；名詞、動詞、形容詞等等
- 詞庫與標註詞性在訓練模組時需要將兩者調整為相同標準

微調基礎詞庫

- 詞庫內容混雜，需要由關鍵詞抽取結果來找尋詞彙是否錯誤



- 分詞 (Tokenization)
 - Jieba 演算法工具
- 標註詞性 (POS tagging)
 - OpenNLP 演算法工具
- 關鍵詞抽取 (Term-Extraction)
 - Atr4s 演算法工具



技術辭庫筆數



553,605



529,619



332,853



238,004

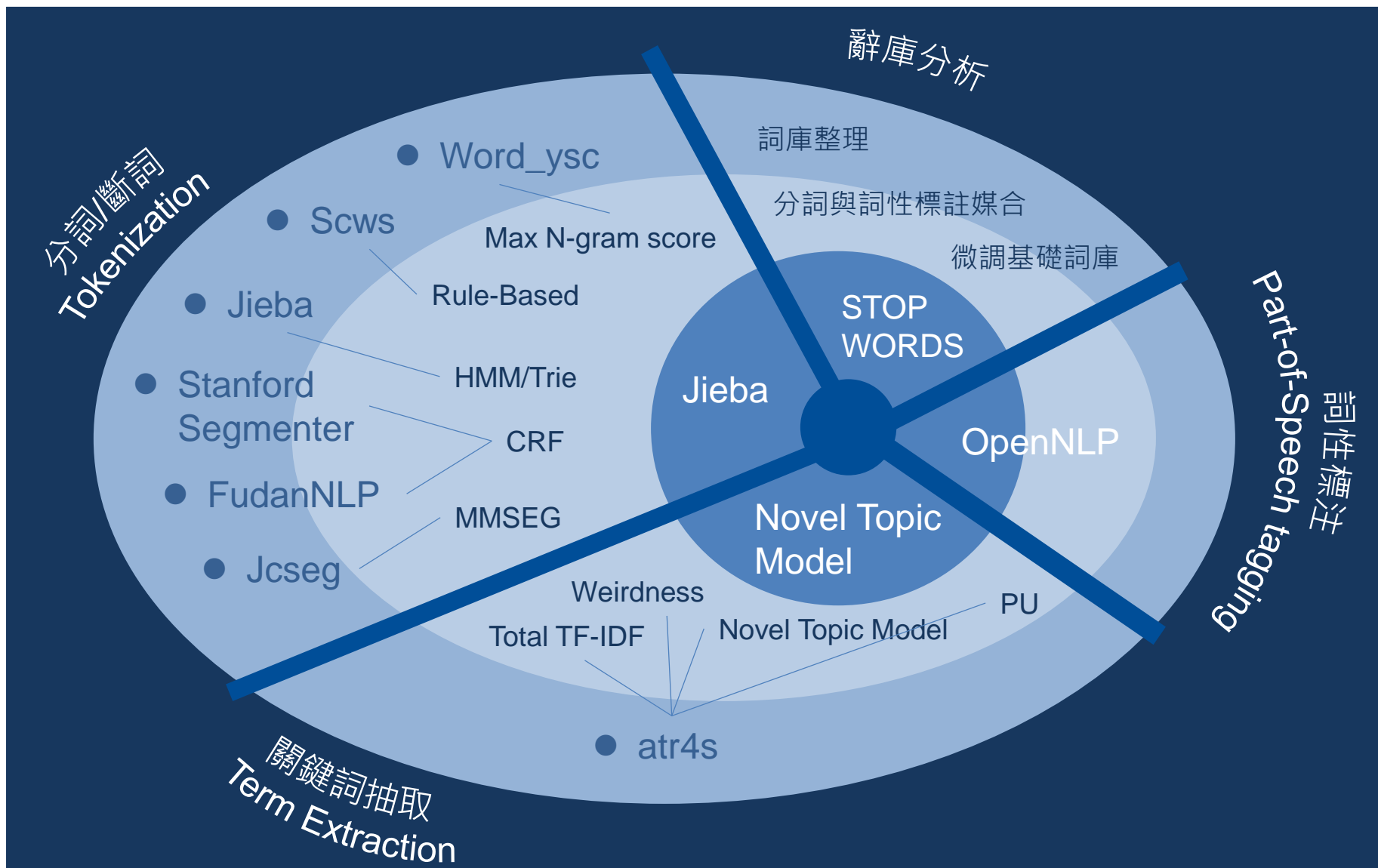


542,888



235,207

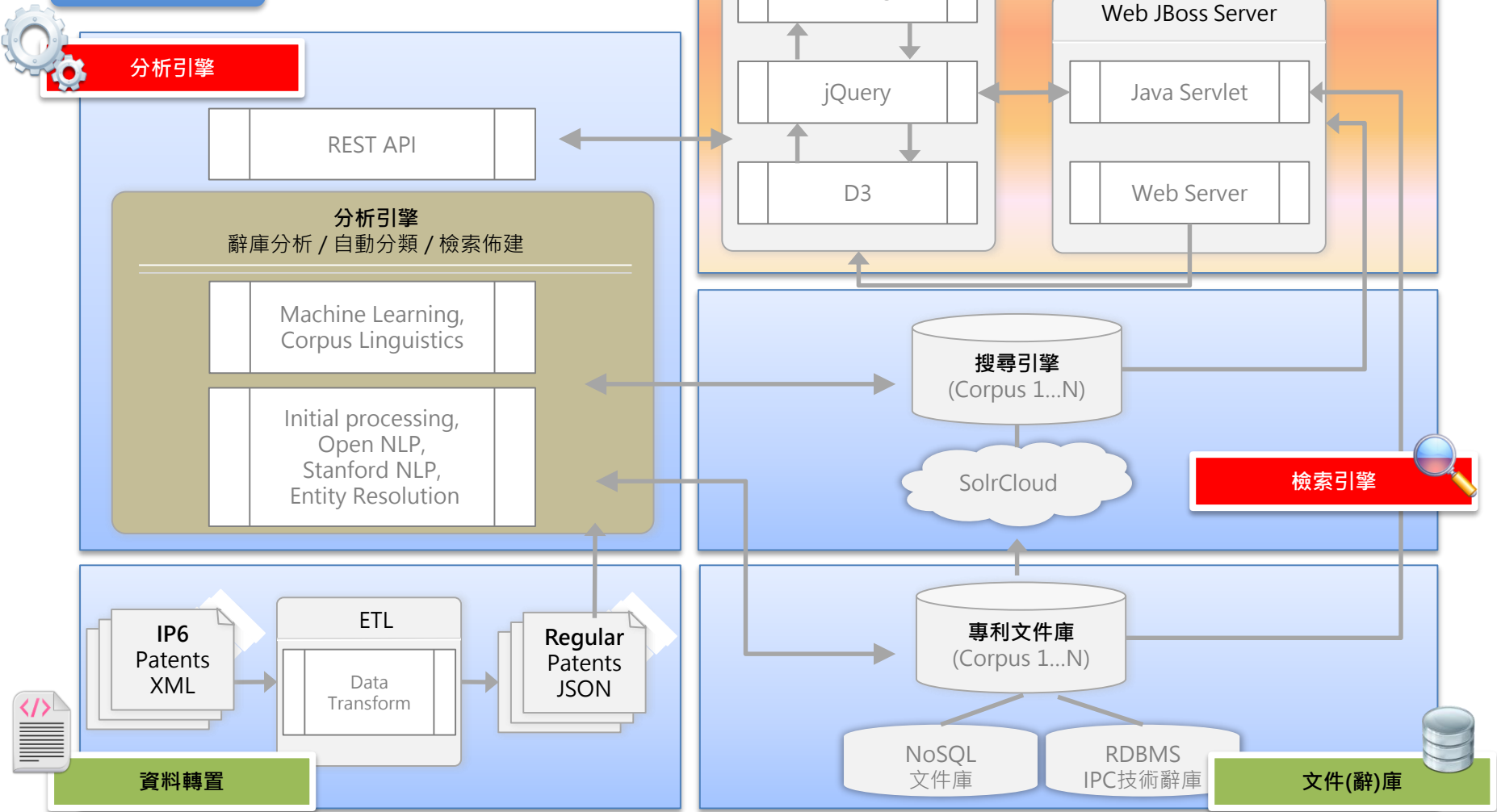
技術辭庫學習



應用系統架構

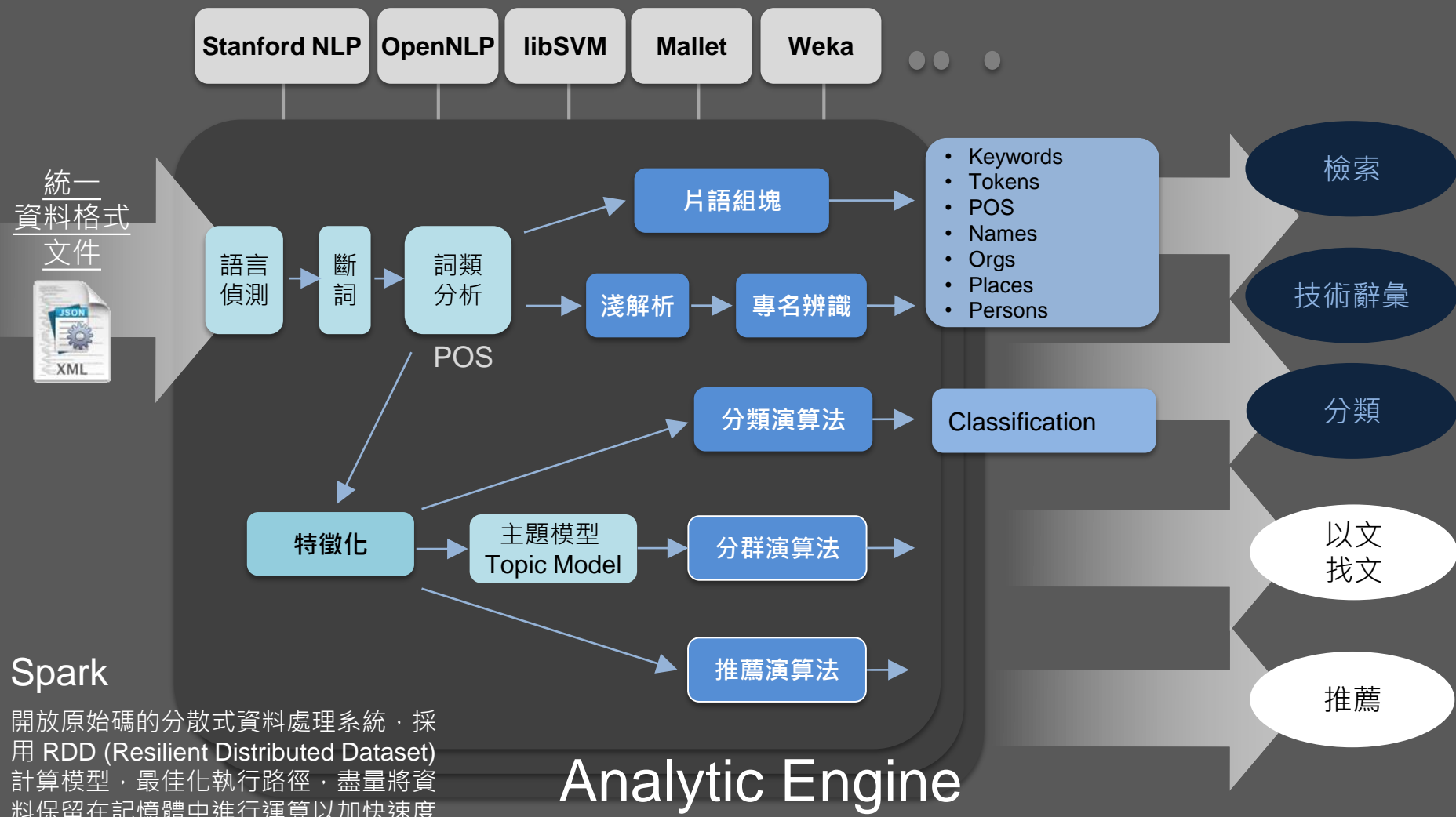
IPKM
產業專利知識平台

應用系統架構示意圖



產業專利知識平台

分析引擎模型



Spark

開放原始碼的分散式資料處理系統，採用 RDD (Resilient Distributed Dataset) 計算模型，最佳化執行路徑，盡量將資料保留在記憶體中進行運算以加快速度

Chapter

3

續章

關於 AI 的迷思
我們的現在與未來

UX
Cloud services
AI Bot DevOps
Big data SaaS
IoT Blockchain

AI 常見的迷思



「透過 google 就可以取得海量資料？」

「得到的這些資料，就可以提供機器學習使用？」

「AI 就是要省人工，如果資料得先整理不就開倒車？」

AI 需要一定量、能有效解析的樣本資料作為學習，隨之而來的
前置作業會有一大堆資料「淨化」工作，並不是只要有大量資
料就可以。

在很多領域裏，它的效能並不會隨著訓練樣本數量的增加而逼
近完美。能在訓練資料中抓出所有相關的輸入特徵尤為重要。

機器學習的用意在於模擬人類的分類和預測能力，不同知識領
域資料內容特性多有不同，有其獨特性與技術專門性。

針對不同知識領域（資料內容），需要採取不同的對策，並經
多次模擬後，才能決定策略及演算手法。

對 AI 的正確認知



「AI 改變了什麼？沒改變什麼？」

「分析引擎建立後，只會越來越聰明？」

「AI 這麼棒，為什麼不盡量把問題都丟給它解？」

AI 不會創造新的商業模式，但提高了效率。

機器學習不是魔術，你自己都不知道使用者的問題如何解決的時，資料科學是不會跑出結果的。

沒有「開箱即用」的 AI 應用；資料不同、解析/學習結果不同。沒有一種模型 (model) 可直接套用於所有知識領域。

不是「找到分析模型」就可結案，而不需再投入；相反地，學習成果的維持需要持續使用、評估、調整。

AI 模擬人，人會出錯，機器也會出錯，人可作即時的修正就馬上的補正，但 AI.....

適合 AI 的應用與問題？



「AI 會讓我的工作消失？」

「先決定明年是 AI 元年，再來看怎麼 AI？」

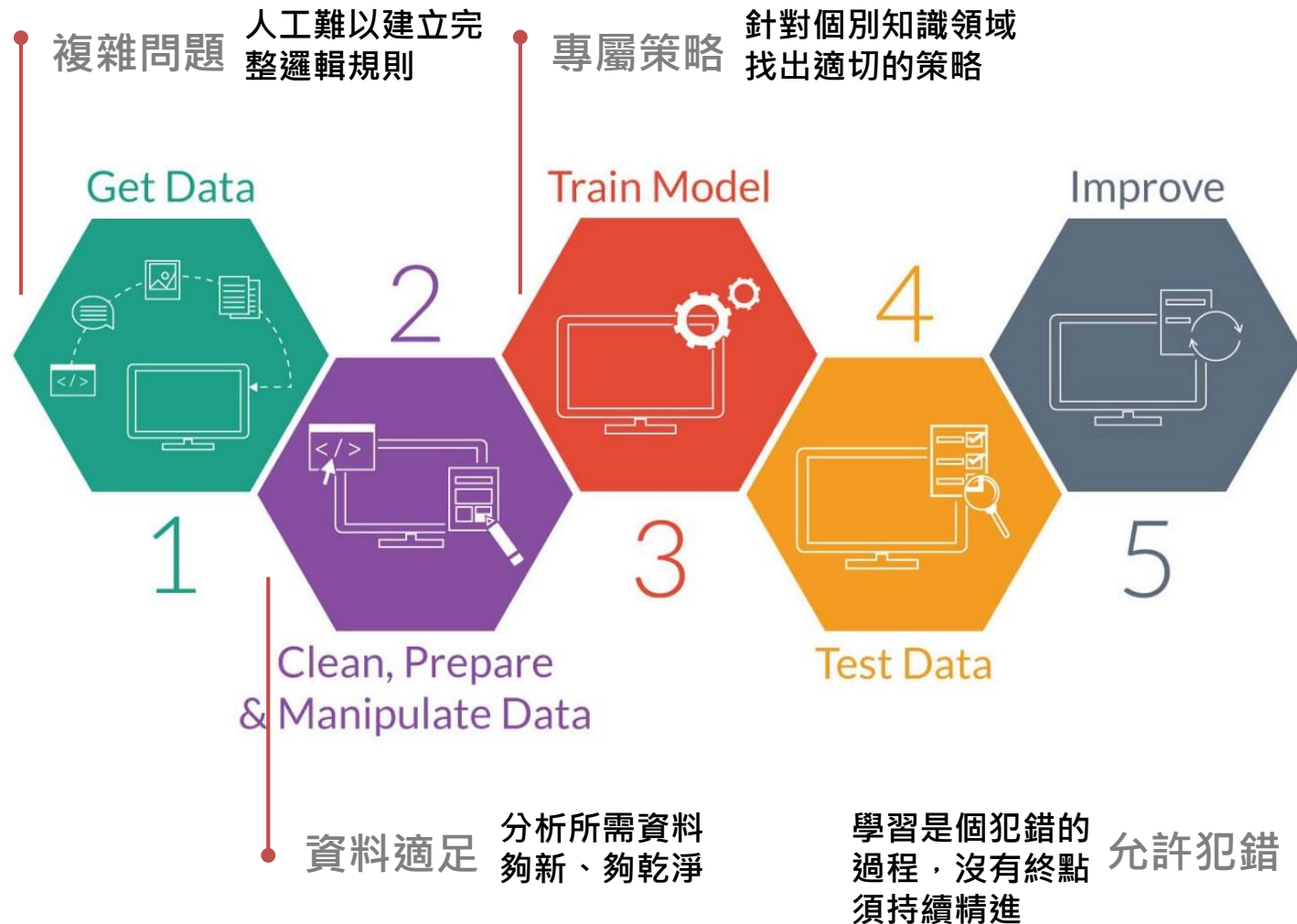
理想應用標的：「如果我們能知道『...』就好了」，鎖定有價值的問題，而非「找一件事讓 AI 解決」，價值愈高的問題，吸引愈多投入企圖與資源。

商業中需要複雜數學能力或 AI 解決的問題還很少。大部分事情透過 AI 解決，反而把事情複雜化。構建這個能力的工作量可能是直接解決它的難度的一百倍、一千倍。

設定易達成的短期目標；
不急於投入大量基礎建設；但要從主要目標分解出易於達成及驗證成果的階段性計畫。

適合用 AI 處理的問題特徵

Typical Machine Learning Workflow



哪些任務最適合 AI ？

標記界定明確的輸入和輸出，能以學習函式將其對應起來的任務

分類、預測；例如「分析一份入學申請來預測未來合格的可能性」

有明確的目標和度量標準、提供清晰反饋的任務

訓練資料明確度有多高或有明確界定的全系統表現尤為重要

例如「優化全內湖範圍內而不是民權東路交叉路口的交通流量」

不需要對於決策過程進行細緻解釋的任務

系統可提供建議，但不代替決策，人類本身就會做出不完美的決策！

例如「醫療診斷」

能夠容忍錯誤、不需要可證實的正確度或最優解決方案的任務

幾乎所有的機器學習演算法都是從統計學和概率上得出解決方案

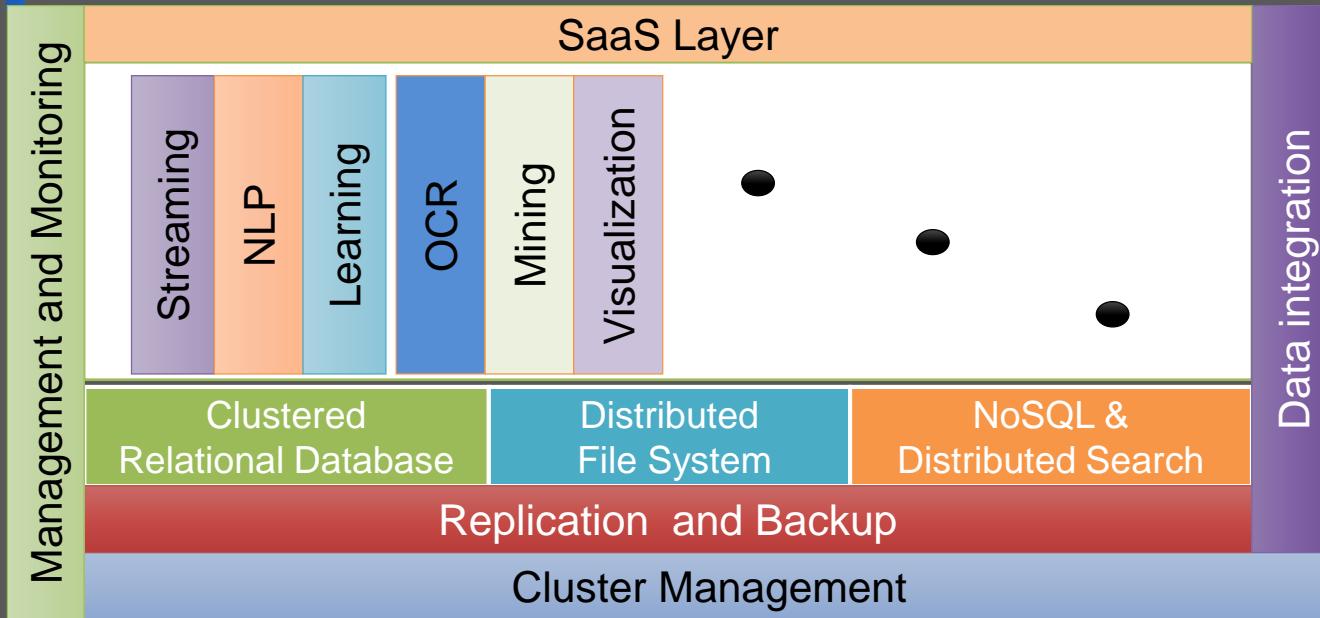
因而，要把它們訓練到百分之百的準確度幾乎不可能



深度研發、創新研究



創新研究所 – NLP Topics



大型分散式系統

NE Recognition (專名辨識)

Parser (解析器)

Tokenization (斷詞)

Recommender Systems (推薦系統)

自然語言機器學習
相關研究主軸

POS Tagging (詞類標示)

Summarization (總結)

Cluster/Classification (叢集/分類)

Keyword Extraction (關鍵字擷取)

Emotion Detection (情緒偵測)

2016

2018

Thank You

感謝您的聆聽，敬請指教

Cloud services
SaaS
AI
DevOps
Big data
Blockchain
IOT Bot

GSS 叢揚資訊
Galaxy Software Services

國家產業創新獎
卓越中堅企業



GSS 叢揚資訊



Vital 雲端服務家族



GSS 技術部落格