



Sync Lab

Synchronous with technologies



GDP

Gathering Detection Platform

Benvenuti



Cristoforo Decaro

[Email](#)

[Linkedin](#)

- Machine Learning developer at SyncLab
- Phd in Electronic Engineering at Università di Ferrara
- Master degree in Electronic Engineering at Politecnico di Torino



Problema:

In seguito alla pandemia del virus COVID-19, i cittadini sono stati costretti a quarantena e soprattutto ad evitare assembramenti

Soluzione:

Creare una piattaforma in grado di rappresentare, mediante rappresentazione grafica/dashboard delle zone potenzialmente a rischio assembramento.

Come:

Bisogna sfruttare le tecniche di **machine learning** e/o **deep learning** per predire le zone potenzialmente a rischio.



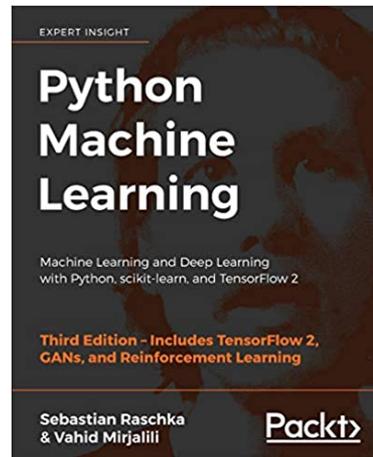
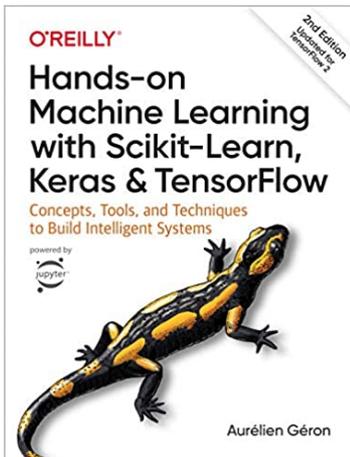
Risorse disponibili:





Risorsse:

1) Libri

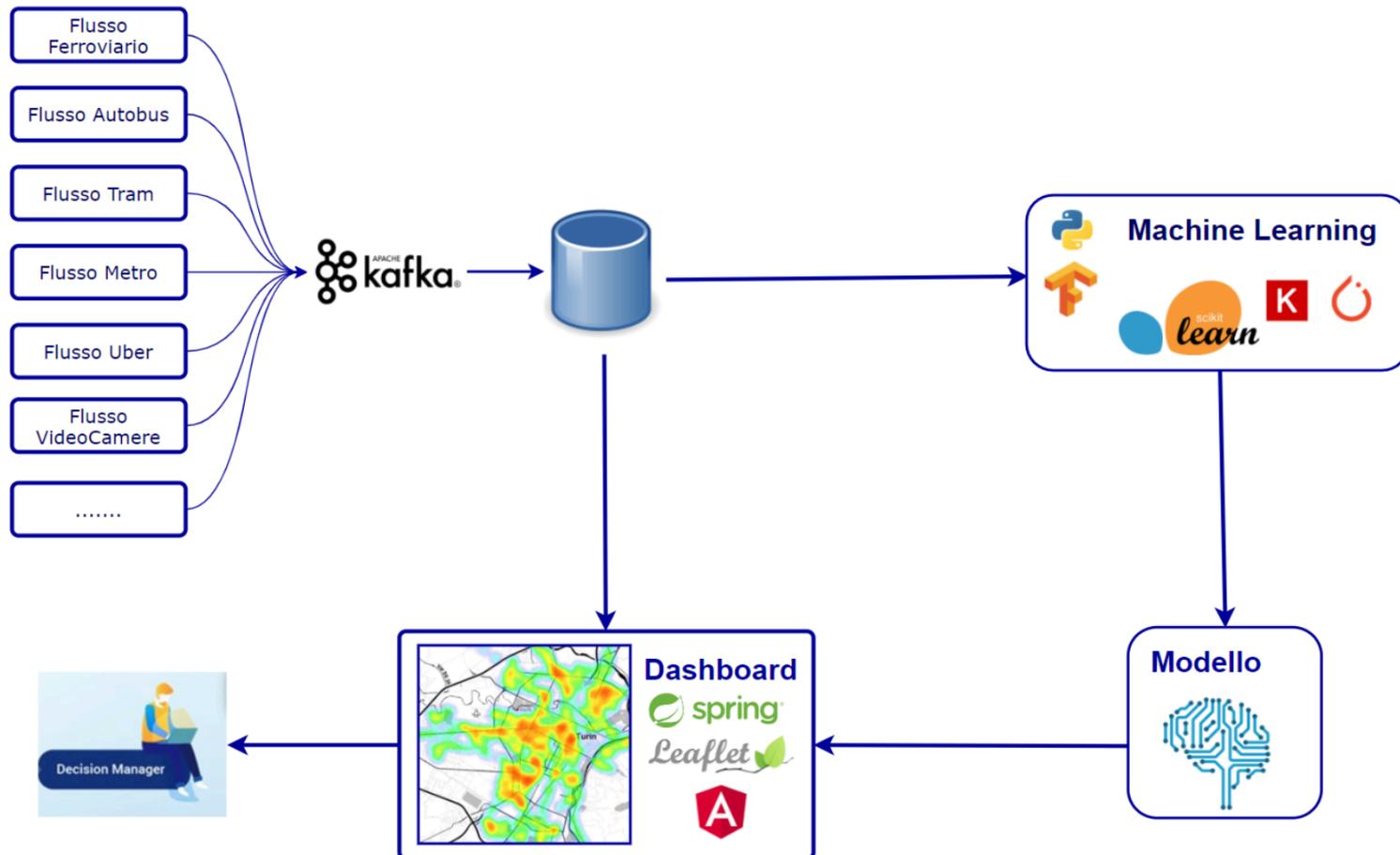


2) Risorsse online

- [Kaggle](#)
- [Aurelien Geron's GitHub](#)
- [Keras documentation](#)
- [Tensorflow guide](#)
- [Pytorch guide](#)
- [Scikit-Learn guide](#)

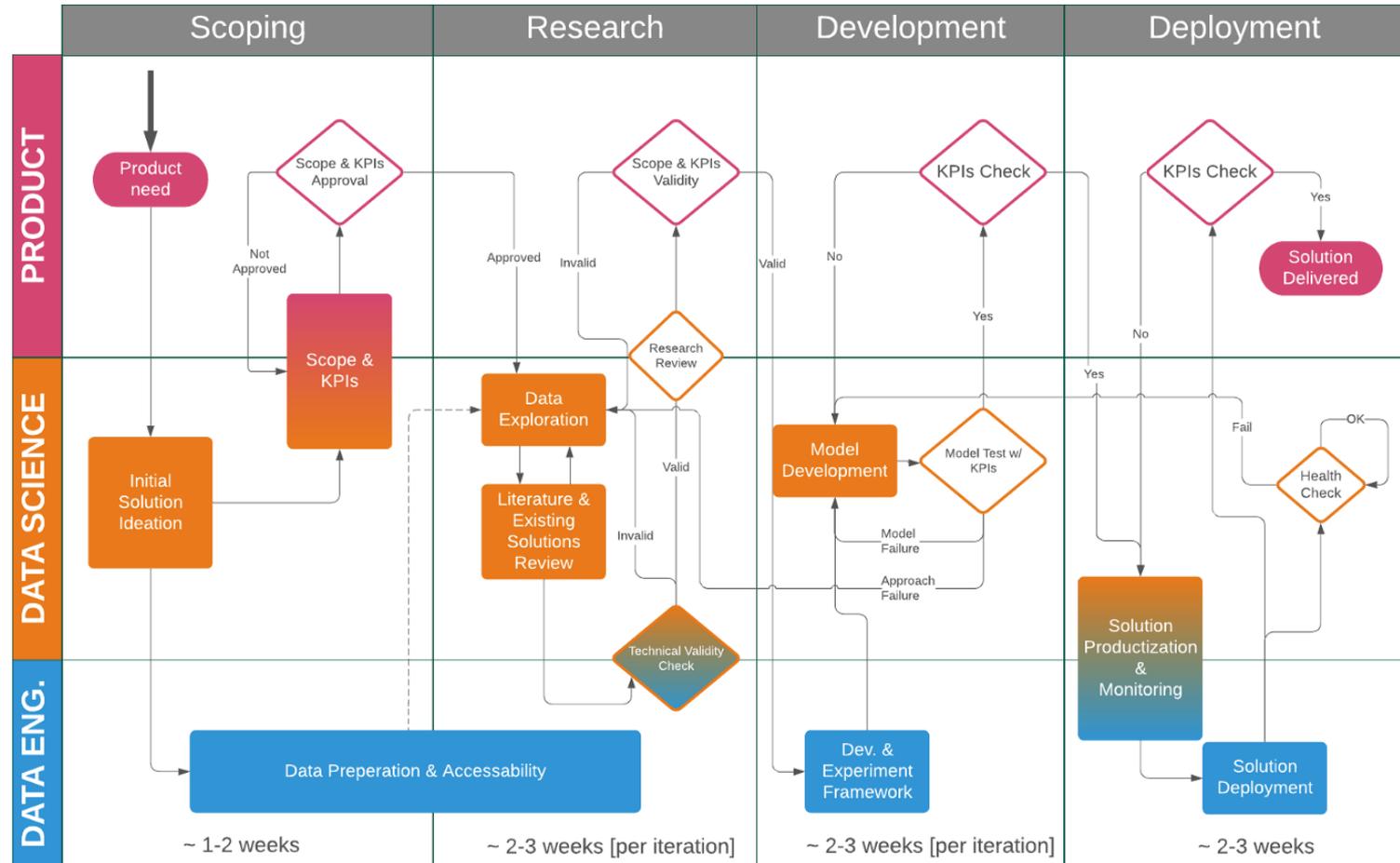
3) Google scholar

Esempio di architettura:



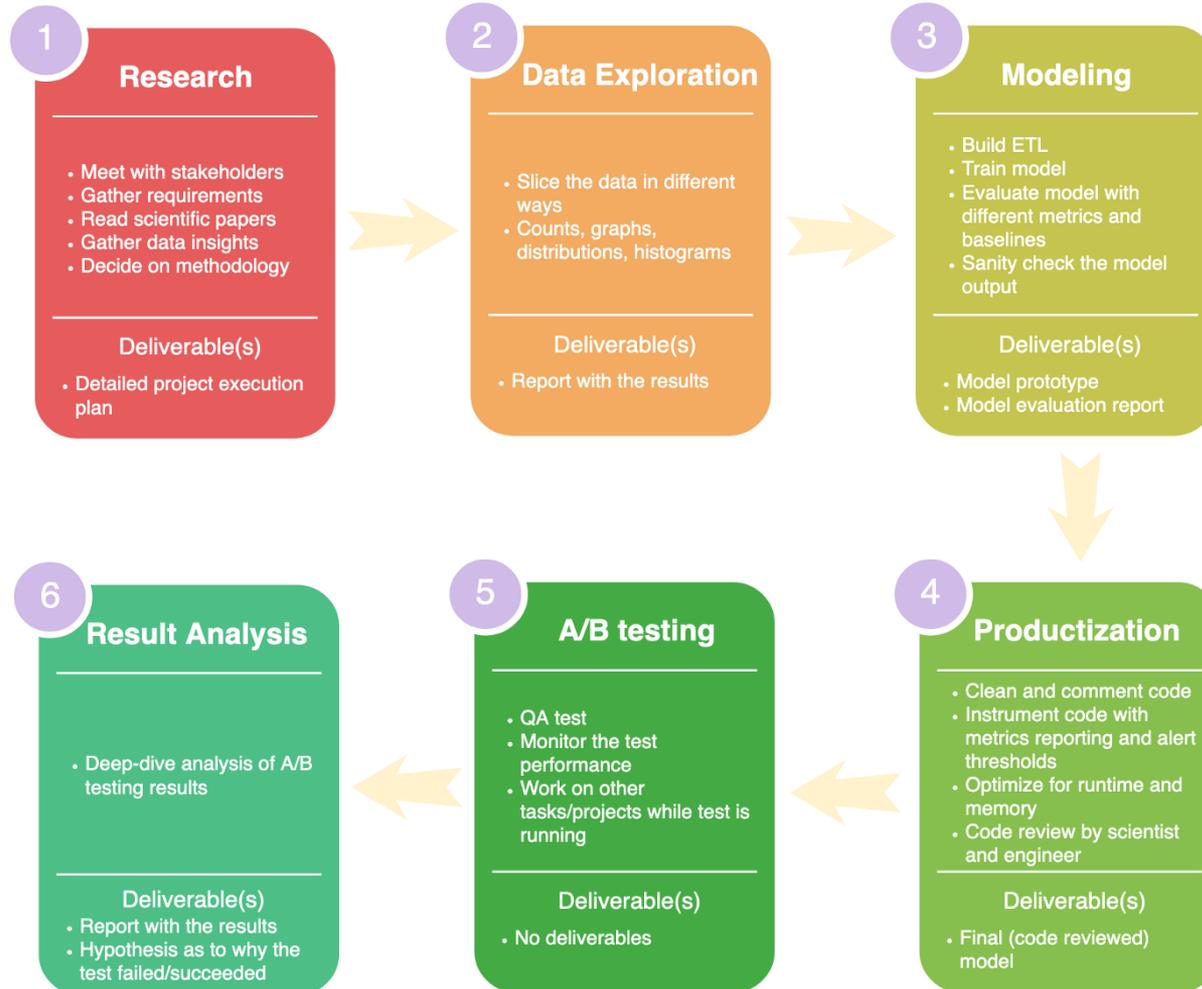


Esempio di architettura complessa:



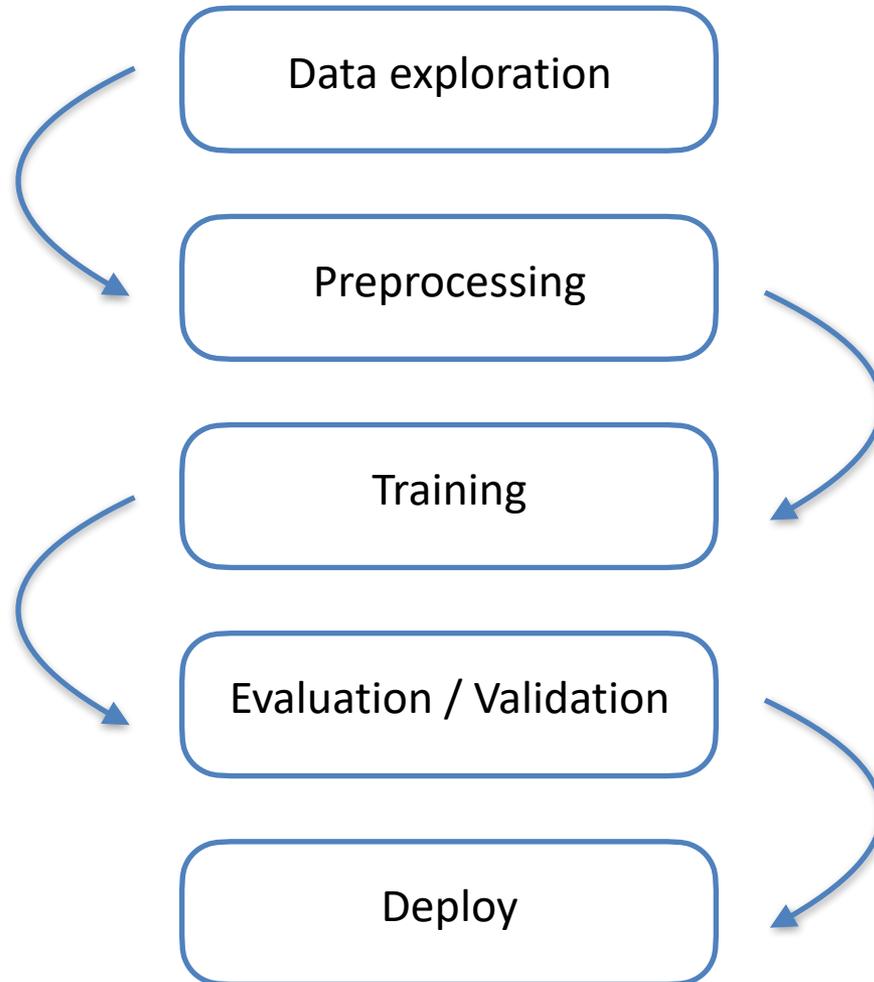


Data science:





Fasi del progetto:

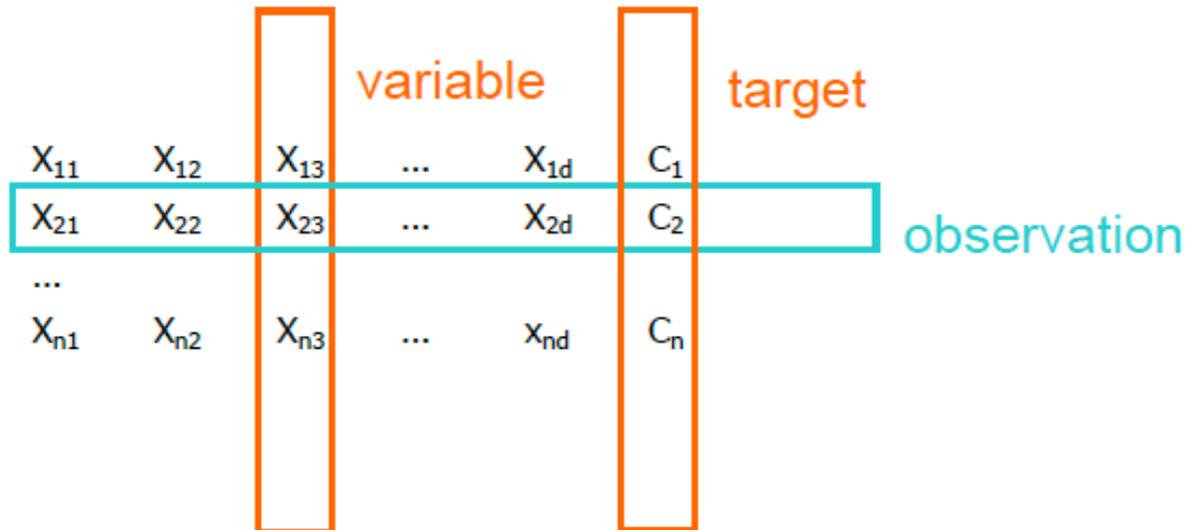




Data exploration

Data exploration:

- Presenza di valori mancanti, outliers, dataset non bilanciati
- Livello di rumore dei dati
- Correlazione dei dati





Preprocessing

Preprocessing:

- Cleaning
- Trasformazione dei dati (normalizzazione, discretizzazione, aggregazione, calcolo di nuove variabili...)
- Feature extraction
- Filtraggio dei dati
- Train / Test set splitting

Come dividere i dati?

- Se i dati hanno una componente **temporale** il training set conterrà i dati **meno recenti**, il validation test e il test set i **più recenti**.
- Se i dati non dipendono dal tempo, la divisione viene effettuata in maniera **casuale**.



Caso Predizione

1) Classificazione (Categorical target)

Naive Bayes

Decision Trees

Neural Networks

KNN

Support Vectors Machine

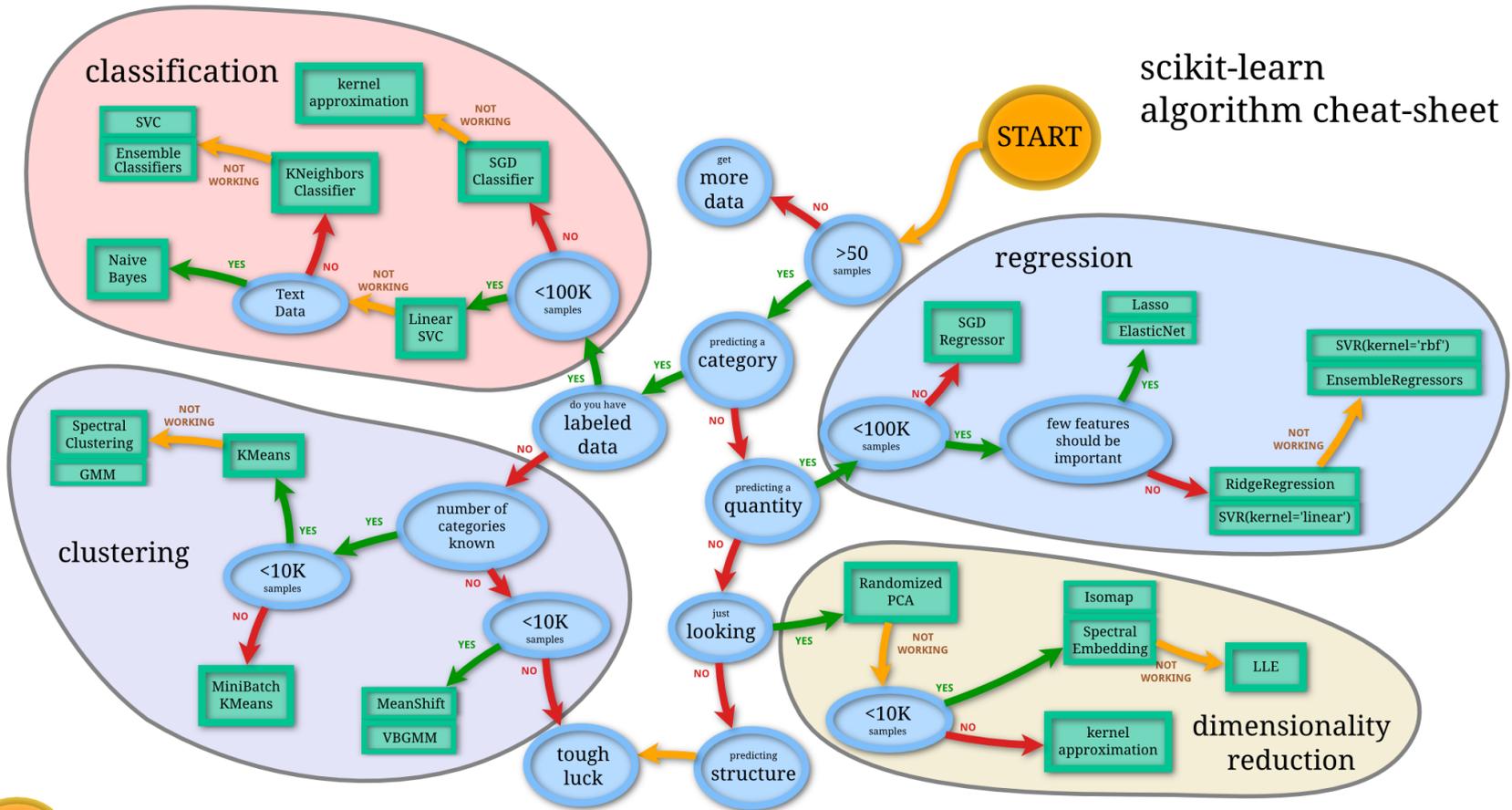
...

2) Regressione (Numerical target)



Scegliere l'algorithmo giusto:

scikit-learn
algorithm cheat-sheet





Evaluation and validation:

Valutazione

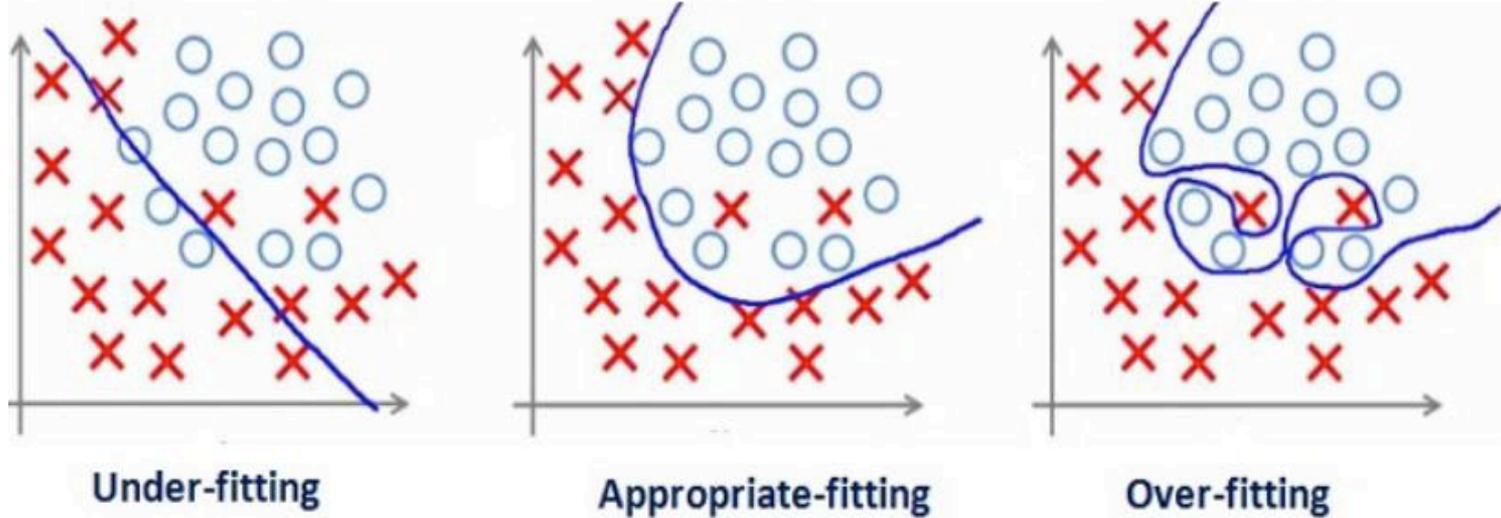
- *Metriche di valutazione: r^2 , MSE, MAE, correlation matrix, ROC curve ...*

Ottimizzazione

- Ottimizzazione degli iperparametri: trials and errors, grid-search, random-search,...



Overfitting:



Per risolvere il problema dell'**overfitting** di solito il dataset viene diviso in tre parti:

- **Training set**
- **Validation set**
- **Test set**

I diversi **modelli** (ottenuti cambiando **algoritmo** o il valore degli **iperparametri**) vengono addestrati sul training set e testati sul validation set per testare la loro capacità di generalizzazione.



Checklist:

- Data Assumptions
- Preprocessing
- Training
- Loss/Evaluation Metric
- Overfitting
- Runtime



Outcomes:

- Report:
 - Project scope and product need
 - Algoritmo e approccio scelto
 - Dati utilizzati, preprocessing
 - Modelli esaminati, training
 - Comparazione dei risultati tramite scelta delle metriche adatte
 - Ottimizzazione degli iperparametri
- Models:
 - Salvataggio del modello
 - Load del modello pre-allenato



Contatti:

Ing. Fabio Pallaro

email: f.pallaro@synclab.it

Cristoforo Decaro

email: cristoforo.decaro@synclab.it