



Apprentissage Automatique Project

Bilal SAID

Ph.D. AI, Research Project Manager



Outline

- Aim
- Scope
- Teams
- Deliverable
- Deadlines
- Evaluation

Aim

- Apply what we learned on new unseen datasets
 - Data discovery, cleaning, exploration, analysis, split, and transformation
 - Model training, evaluation, comparison, and hyperparameter optimization
 - Analysis of the results, critical thinking and outlook for future enhancement
- Get inspired by others (including ChatGPT, GitHub, Kaggle, friends, our lab codes...)
- Master the content
- Show a coherent methodology
- Justify with clear reasoning

Scope

- Target Datasets



- Small/Medium with 10s of features (10-20) & few 10-100Ks instances
- Dense, mostly clean (well structured, not too many missing values...)
- Previously used and verified (known source, not artificial...)
- Numeric AND Categorical features

362 datasets found

dense



>10 features



>10000 instances



verified



- **Optional & Complementary Only:** pure text (e.g., reviews), images (product photo)...



- **Better Excluded:** simulated/generated datasets, time series, dynamic stream...

- Target ML Tasks

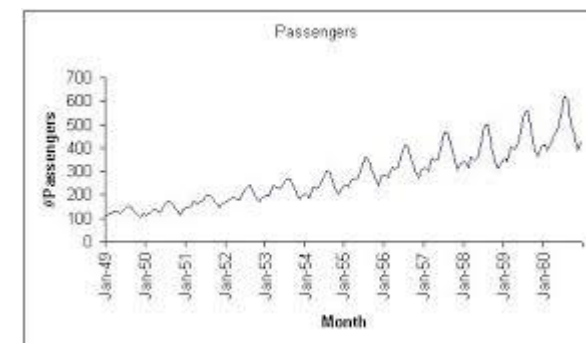
- Regression
- Classification (binary or multi-class)
- Hyper-param tuning



- **Optional & Complementary Only:** deep learning (text, images), LLM...



Dataset	DenseNet	PCN	SLCNet	TopNet	MSTN	WGAN-GP	ED	DTW	DTW	DTW	DTW	DTW	DTW
	-CT_CAM	-CT_CAM	-CT_CAM		-FCN	+HISE	-INN	-INN	-INN	-INN	-INN	-INN	-INN
Anticlock	0.9867	0.9867	0.9867	0.987	0.975	0.99	0.97	0.96	0.987	0.97	0.96	0.987	0.987
WorkRecognition	0.4667	0.4667	0.4667	0.333	0.267	0.333	0.267	0.267	0.2	0.267	0.267	0.267	0.22
Animal	1	1	1	1	0.85	1	0.85	1	0.85	0.85	1	0.85	0.85
Brain-Motion	0.997	0.991	0.8186	0.997	0.945	0.96	0.964	0.968	0.99	0.964	0.968	0.968	0.988
Character	0.5061	0.356	0.5061	0.356	0.515	0.515	0.515	0.515	0.515	0.515	0.515	0.515	0.515
Hand-Motion	0.473	0.4665	0.4459	0.378	0.365	0.365	0.279	0.306	0.236	0.278	0.306	0.236	0.236
Hand-Motion	0.8164	0.8068	0.7196	0.751	0.743	0.717	0.743	0.717	0.717	0.717	0.717	0.717	0.717
Hand-Motion	0.46	0.37	0.64	0.58	0.51	0.5	0.51	0.39	0.5	0.51	0.39	0.5	0.5
NAO-CPS	0.9833	0.9889	0.8833	0.939	0.889	0.87	0.86	0.85	0.883	0.85	0.85	0.883	0.883
PEMS-SF	0.7977	0.7146	0.7146	0.751	0.699	N/A	0.705	0.734	0.711	0.705	0.734	0.711	0.711
RealSigns	0.9817	0.9886	0.986	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98	0.98
Peasants	0.3088	0.3536	0.3536	0.175	0.11	0.39	0.164	0.151	0.151	0.164	0.151	0.151	0.151
SelfRegulation	0.5944	0.5944	0.5889	0.55	0.472	0.46	0.483	0.533	0.539	0.483	0.533	0.539	0.539
SelfRegulation	0.9843	0.9923	0.9841	0.983	0.99	0.982	0.967	0.96	0.963	0.967	0.969	0.963	0.963
StandWalkJump	0.6667	0.5	0.5	0.4	0.067	0.333	0.2	0.333	0.2	0.2	0.333	0.2	0.2
Avg. Value	0.746	0.733	0.711	0.701	0.66	0.666	0.633	0.633	0.666	0.633	0.633	0.633	0.633
Std. Value	0	0	0	0	0	0	0	0	0	0	0	0	0
Avg. Rank	2.067	2.7	4.361	3.631	7.535	5.885	9.5	8.5	1.333	9.7	8.5	1.333	1.333



Teams

- Two students per team (i.e., 12 teams)
- Fill-in this form <https://forms.gle/JEvfWHzqGcu6kvgK9>

Deliverable

- Single shared Google Colab Notebook
- Regularly pushed on a publicly shared GitHub repository
- Well-documented
- Organized into separate sections
- Minimal extra work to load dataset and install extra python libs

Notebook Structure (1/3)

- Text cell for Project & Dataset Description
 - Project aim(s), existing solutions
 - Dataset original source URL(s)
- Code cell containing **ALL** libs imports
 - Comments for libs versions or installation instructions, if any
- Data Access Instructions
 - Dataset GitHub-hosted project-version URL(s)
 - Data Load Code to access it through GitHub (if slow or technically not feasible, then from G. Drive using a Google Collab)

Notebook Structure (2/3)

Dataset Exploratory Analysis

- Metadata: nb. instances, nb. features, types + your comments
- Nb. nulls or na : code for identification, then your interpretation, & possible and chosen strategies to cope with them with arguments
- Features values distributions, scaling & outliers: possible and chosen strategies to cope with them with arguments
- Target feature study: distribution / class balance, your interpretation & decision to cope with it
- Features correlation, selection, dimension reduction & anomalies detections: your explanation, decision on whether to proceed with different data input in # pipelines
- Unsupervised clustering: correlation with target feature,
- Interpretations, hypotheses and conclusions

Notebook Structure (3/3)

ML Baseline & Ensemble Models

- Training / Validation / Test Splits
- Pipelines & Models
- Training / Validation Code
- Test Code
- Interpretation & Discussion of Results

Enhanced Models, Hyperparameter Tuning & Analysis

- Justification of choices
- Interpretation and analysis of the obtained results

Deadlines

Date	Submission
Nov 14	Teams & Topics Data Exploratory Analysis & Unsupervised Exploration
Nov 28	ML Baseline & Ensemble Models
Jan 23	Enhanced Models, Hyperparameter Tuning & Analysis Final Submission
Feb 03	Defense

Evaluation

[Tentative]

- 5% Following Instructions
- 5% Notebook organization
- 30% Data Exploratory Analysis
Unsupervised Exploration
- 30% ML Baseline & Ensemble Models
- 30% Enhanced Models, Hyperparameter Tuning & Analysis