

Introduction to Data Science

with Python

Alexis Bogroff

May 25, 2022



Alexis Bogroff

Lecturer and Mentor in Data Science
at Paris 1 Panthéon-Sorbonne, ESILV,
Openclassrooms, EM-Lyon



Alexis Bogroff

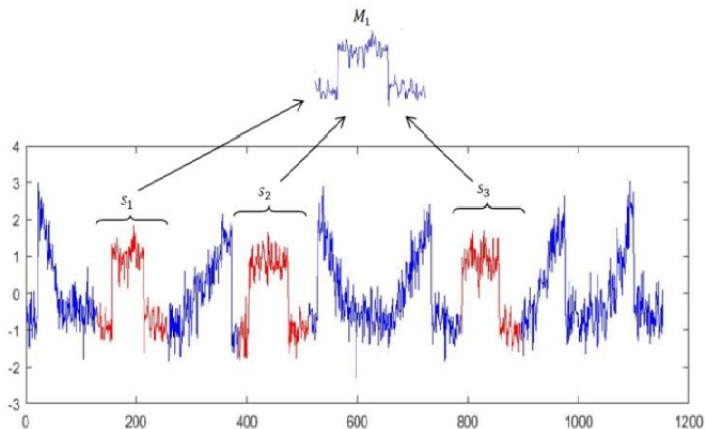
Lecturer and Mentor in Data Science
at Paris 1 Panthéon-Sorbonne, ESILV,
Openclassrooms, EM-Lyon

- 4 years Teaching Assistant and lecturer in VBA, Python for finance, SQL, Data Analysis and Data Science
- 9 months Researcher Assistant at Paris 1 Panthéon-Sorbonne within H2020 European Project
- 1 year Data Scientist at Pléiade Asset Management

Predictions: What does that mean?

What is modeled?

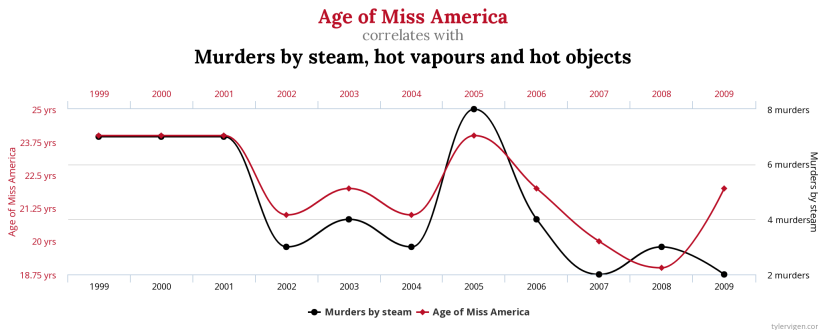
- Continuity (stationarity)
- Correlation (pattern)



Predictions: What does that mean?

What is modeled?

- Correlation vs Causality¹



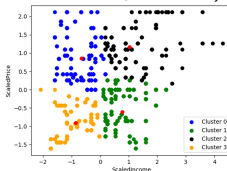
¹<https://www.tylervigen.com/spurious-correlations>

Predictions: Examples

- Present:

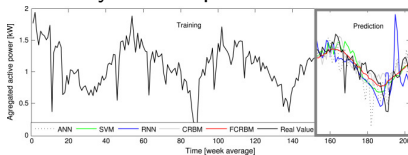
- Electricity consumption based on other cities (e.g. Seattle)
- Missing values (interpolation, extrapolation)

- Client category



- Future:

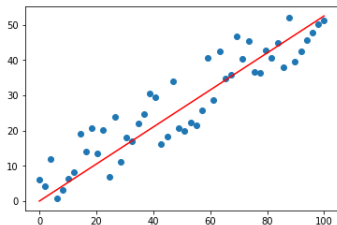
- Electricity consumption next month (time series)



- Client clicking add (recommender sys.)
- Pedestrian and cars trajectories (RL)

Regression

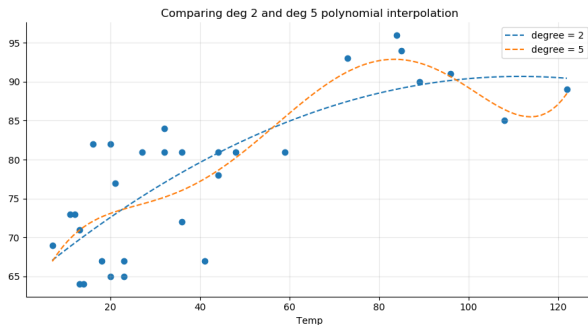
- Linear regression
- Simple: $Y = aX + b$
- Multiple: $Y = a_1X_1 + a_2X_2 + \cdots + a_nX_n + b$



Regression

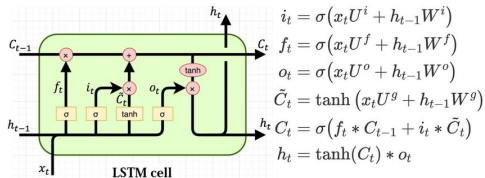
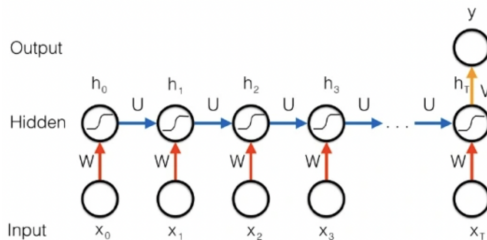
- Polynomial regression
- Simple: $Y = a_1X + a_2X^2 + \dots + a_nX^n + b$
- Multiple:

$$Y = a_{11}X_1 + a_{21}X_2 + a_{n1}X_n + a_{12}X_1^2 + a_{22}X_2^2 + a_{n2}X_n^2 + \dots + a_{1m}X_1^m + a_{2m}X_2^m + a_{nm}X_n^m + b$$



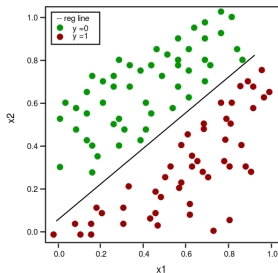
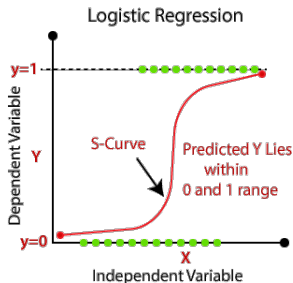
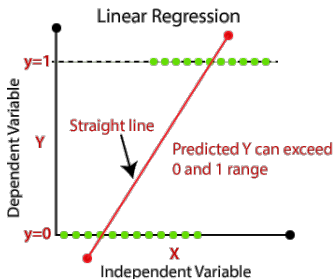
- SARIMA: seasonal, auto-regressive, integrated, moving average
- Parametric model that capture patterns like:
 - Trend
 - Cycle
 - Season

- Deep Learning (RNN, LSTM)



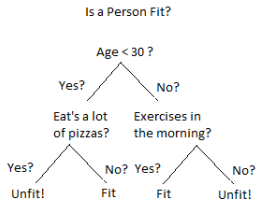
Classification

- Logistic regression $\frac{1}{1+e^{-z}}$ with z linear (polynomial) regression

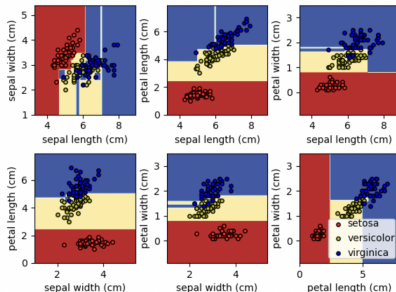


Classification

- Tree (ensemble models, RF, XGBoost)

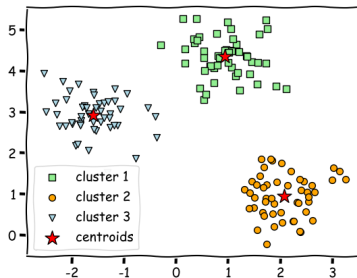


Decision surface of decision trees trained on pairs of features



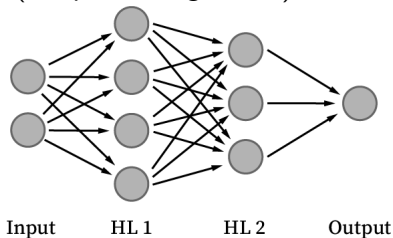
Classification

- K-means
 - Positionnement
 - ① Capture (iter 1)
 - ② Recentrage (iter 2)

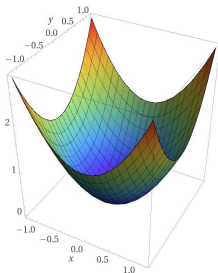


Classification

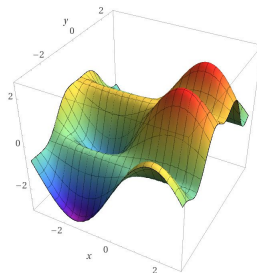
- Neural Network (Deep Learning: CNN)²



- Gradient descent



Computed by WolframAlpha



Computed by WolframAlpha

- Clustering (grouping)
- Dimensionality reduction
- Reducing multicollinearity
- Models:
 - K-Means: entropy minimisation principle (min var intra, max var inter)
 - Hierarchical Clustering
 - PCA

- Parameters
- Hyperparameters
- Train, cross-validate, test
- Feature importance
- Data Leakage

- Generalization
- Complexity
- Over/Under-sampling
- Unbalanced Datasets (weights on cost function, SMOTE, Auto-encoder)

- Regression
 - RMSE
 - R^2
- Classification
 - Accuracy
 - AUC, ROC Curve
 - Other metrics based on confusion matrix

Transfer Learning, Why?

- Training can be complicated, long and expensive
- Specific but complex (and similar) task (NLP)
- Few samples

What has been learnt?

- Weights value (or centroids)
- Hyperparameters

- Optimize target objective on long term, intermediate steps on short term:
 - Increase task difficulty gradually
 - Better generalisation: Multi-task learning (RL, learn recognize unrelated objects)
 - Improve Neural Network architecture (genetic algorithms)

Some code examples

- Sklearn simple 4 lines of code
- More advanced Sklearn
- Deep Learning with Pytorch