BORDEAUX INP Enseirb-Matmeca

# ENSEIRB–MATMECA

8 Nov 2024

CHAIRE IA
DIGNE DE CONFIANCE

# Agenda

**01**

Introduction

**02**

About Fieldbox

**03**

What is MLOPS?

**04**

Imbalance

# Introduction

01

# Three courses

The three courses will be on three mornings:
→ Nov 8th - 8:30 to 12:30
Nov 22nd - 8:30 to 12:30
Nov 29th - 8:30 to 12:30

If you can't come please let me know by email as soon as you can.

We will have a hands-on lecture, with discussion.

We will alternate between black-boarding and coding in Python.

You can send me your code or commit it to git. Deadline: end of the course.

At the end of each course I will propose ways to go further if you want.

My name is Julien Budynek and I am VP Data Science at Fieldbox.
I work on Data and AI projects delivered for Fieldbox clients, and on internal projects to advance our internal AI R&D roadmap.

You can reach me at jbudynek@fieldbox.ai

The supporting files for this course are available on my github:
https://github.com/jbudynek/teaching-01-imbalance

# About Fieldbox

02

# Fieldbox mission

Our mission is to ensure industrial companies can rely on digital, data & AI to deliver sustainable efficiency bringing resilience & agility.

**DIGITAL DATA & AI**

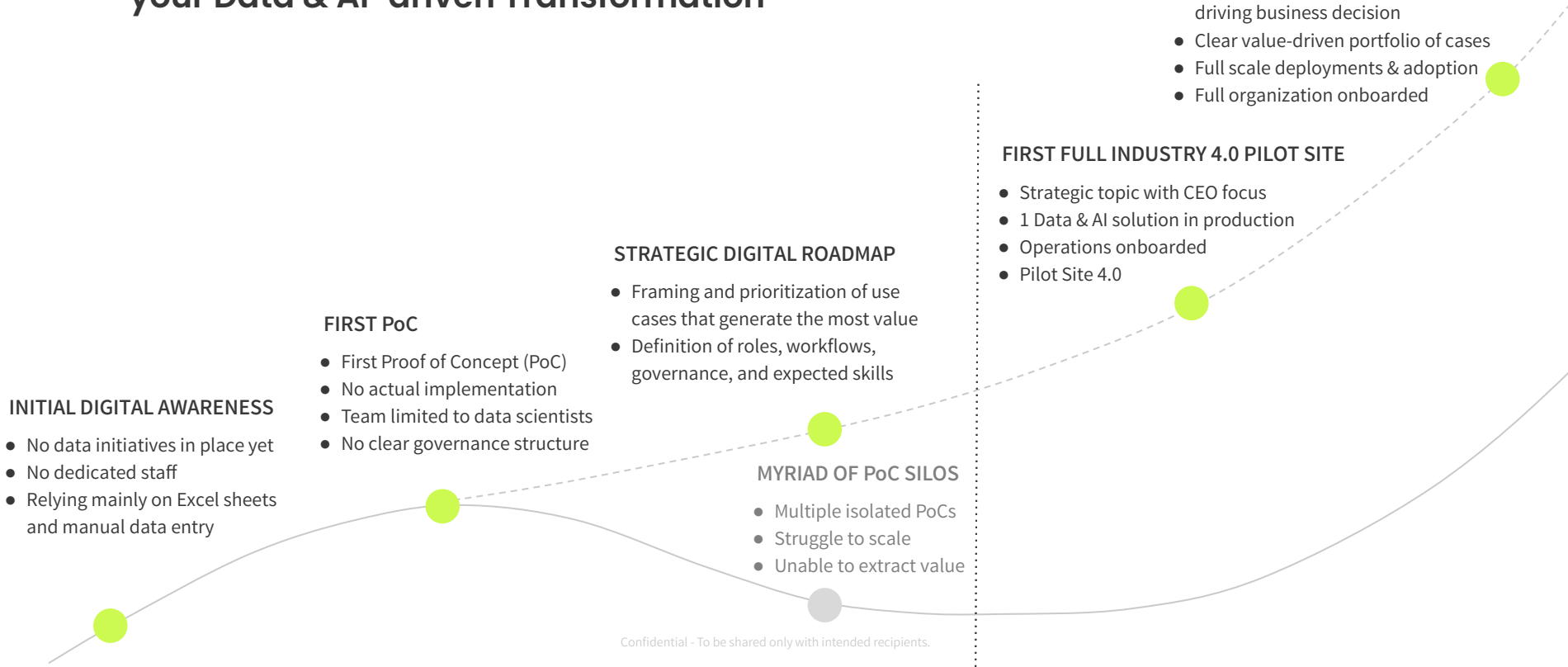are required to manage the increasing complexity of industrial operations.

**80%**

of Data & AI solutions are not put into production and do not produce value.
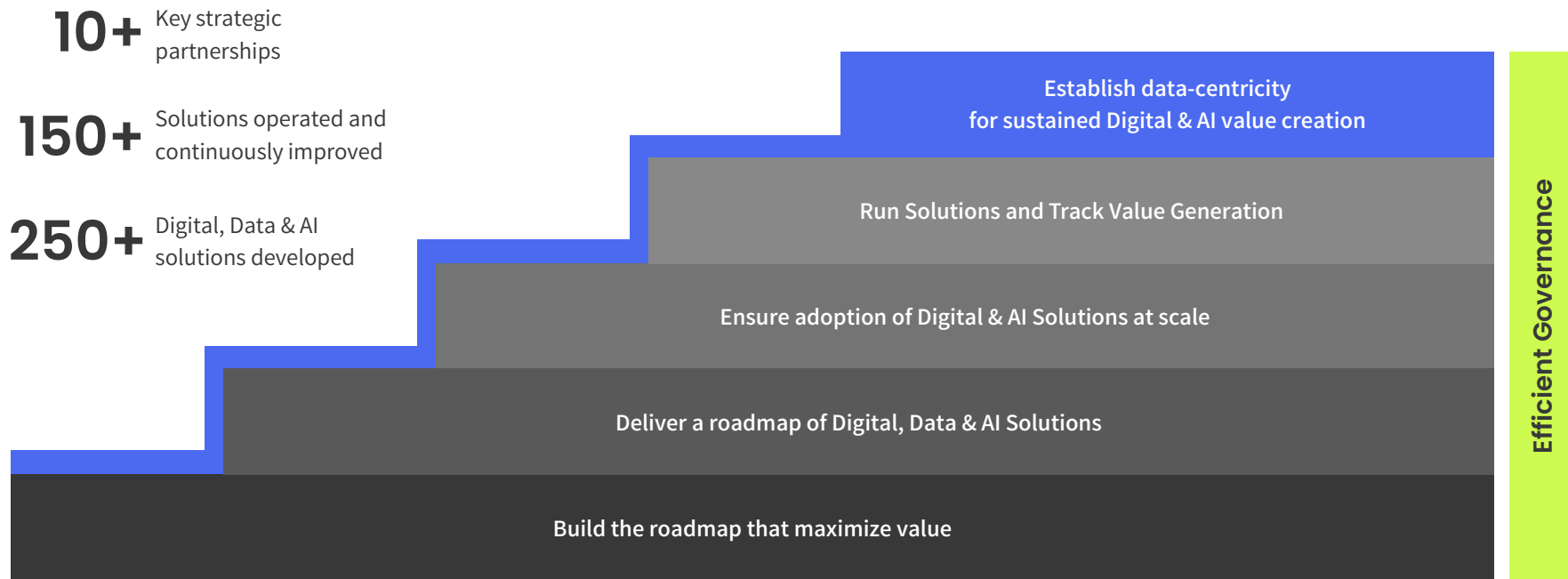
# Go beyond the Wall of Operationalization

## Fieldbox helps you secure & accelerate your Data & AI-driven Transformation

**FULL SCALE INDUSTRY 4.0**

- Data & AI-driven mindset, driving business decision
- Clear value-driven portfolio of cases
- Full scale deployments & adoption
- Full organization onboarded

**FIRST FULL INDUSTRY 4.0 PILOT SITE**

- Strategic topic with CEO focus
- 1 Data & AI solution in production
- Operations onboarded
- Pilot Site 4.0

**STRATEGIC DIGITAL ROADMAP**

- Framing and prioritization of use cases that generate the most value
- Definition of roles, workflows, governance, and expected skills

**FIRST PoC**

- First Proof of Concept (PoC)
- No actual implementation
- Team limited to data scientists
- No clear governance structure

**INITIAL DIGITAL AWARENESS**

- No data initiatives in place yet
- No dedicated staff
- Relying mainly on Excel sheets and manual data entry

**MYRIAD OF PoC SILOS**

- Multiple isolated PoCs
- Struggle to scale
- Unable to extract value

# Fieldbox is your Digital & AI Scale-Up Partner

**10+** Key strategic partnerships

**150+** Solutions operated and continuously improved

**250+** Digital, Data & AI solutions developed

Establish data-centricity
for sustained Digital & AI value creation

Run Solutions and Track Value Generation

Ensure adoption of Digital & AI Solutions at scale

Deliver a roadmap of Digital, Data & AI Solutions

Build the roadmap that maximize value

**Efficient Governance**

**Fieldbox support your organization in every steps toward sustainable value generation**

# Using Digital, Data & AI to Reach Operational Excellence in Industrial Operations

| Safety & Risk detection | Energy Management | Predictive Maintenance | Inventory Management |



| Quality Control | Process Control | Production Optimization | Supply Chain Optimization |

# Fieldbox: a Decade of Cutting-Edge Data & AI Solutions for Industry

Fieldbox offices

Customer locations operated with Fieldbox solutions

- Operating on 5 continents
- Independent & organic growth

Bordeaux
Paris
Singapore

**2011**
Creation of the company (as IDMOG)

**2014**
Extension to Africa & South America

**2016**
New AI capabilities
Expansion to new verticals

**2018**
New offices in Paris & Singapore

**2020**
+50% revenue growth

## Selected references

VEOLIA · suez · VALOREM · alperia · HUTCHINSON · SAINT-GOBAIN · LafargeHolcim · EUROVIA

TotalEnergies · BR PETROBRAS · SBM OFFSHORE · TRIDENT ENERGY · KRISENERGY · SNCF · TELESPAZIO

## Partners

bpifrance · Hub · Valley · POLEAVENIA

# Data project lifecycle

| Pre-sales sales | Architect | Data engineers | AI research | DS Team | AI engineers | Software engineers | IT, Ops | Services | Run Team |
|---|---|---|---|---|---|---|---|---|---|

Build project team: Fieldbox team **& client**

**Pre-sales phase**

Defining scientific and technical objectives
Global project planning
Data batch gathering

Train-test splitting

Data prep, Explo, process

Validation

Models training

Feature selection

Hyperparam tuning

**Deliver / Deploy**

Live test
Monitor
Maintain
Re-train
Coach

**Running phase**

**~ months**

# What is MLOPS?

03

# What is MLOPS?

There are several ways to look at MLOPS.

Some see it as the complete process governing the lifecycle of a machine learning model: Data collection - Analytic formulation - Modelization - Test - Deployment - Monitoring - and loopback.

Some say that MLOPS is the practice of making all of those steps easier with tooling.

We will focus on some of the upstream activities, and also on deployment and monitoring.

"MLOps is the standardization and streamlining of machine learning lifecycle management." Wikipedia - march 2022

"MLOps is a set of practices that aims to deploy and maintain machine learning models in production reliably and efficiently"
Wikipedia - nov 2022

https://en.wikipedia.org/wiki/MLOps

Interest over time    Google Trends

● mlops

100

75

50

25

Jun 1, 2017          Apr 1, 2021

Worldwide. 5/11/17 - 11/5/22. Web Search.

MLOPS

Whole ML lifecycle

Build

Ops

Data exploration

Business question

Model deployment

Model serving

Model creation

Model testing

Feedback

Model monitoring

Improve

Everything

# What is MLOPS?

Kreuzberger, Dominik, Niklas Kühl, and Sebastian Hirschl. "Machine Learning Operations (MLOps): Overview, Definition, and Architecture."
https://arxiv.org/abs/2205.02302

**Roles**
R1 Business Stakeholder
R2 Solution Architect
R3 Data Scientist
R4 Data Engineer
R5 Software Engineer
R6 DevOps Engineer
R7 ML Engineer/MLOps Engineer

**Principles**
P1 CI/CD automation
P2 Workflow orchestration
P3 Reproducibility
P4 Versioning
P5 Collaboration
P6 Continuous ML training & evaluation
P7 ML metadata tracking/logging
P8 Continuous monitoring
P9 Feedback loops



**Tech Components**
C1 CI/CD Component (P1, P6, P9).
C2 Source Code Repository (P4, P5).
C3 Workflow Orchestration Component (P2, P3, P6).
C4 Feature Store System (P3, P4).
C5 Model Training Infrastructure (P6).
C6 Model Registry (P3, P4).
C7 ML Metadata Stores (P4, P7).
C8 Model Serving Component (P1).
C9 Monitoring Component (P8, P9).

Imbalance

04

# About imbalance

We will consider the question of 2-class classification.
In general we consider or we suppose that both classes are equally represented in the dataset.
What is imbalance? It is when one of the classes in less frequent than the other.
In 2-class classification, if one class is <10% we have imbalance.



Imbalanced Dataset



Freq of emojis in a whatsapp conversation

# More imbalance

Imbalance shows up all the time in real life, and in particular in industrial use cases.
**Predictive maintenance**: some equipments have dysfunctions, but sometimes only once every two years!
**Quality control**: by design, the proportion of samples that do not meet the quality is very low.
The fields of anomaly detection and rare events prediction are linked to imbalanced data.

**Finance**: fraud detection datasets commonly have a fraud rate of ~1–2%
**Ad Serving**: click prediction datasets have click-through-rate <1%.
**Medical**: classification of patients with rare condition
**Content moderation**: detection of unsafe content

# Automated pipe condition analysis and reporting



- Training of a computer vision model on a set of categories and sub-categories
- Application for the automatic detection of anomalies
- Co-development of an off-the-shelf product with the client's SMEs

TIME REDUCTION FOR PIPE INSPECTION BY 50%

REPORTING TIME REDUCTION BY 80%

FieldBox.ai

# Hands on coding – expected output

In the following sections we will see three ways to deal with imbalanced data:
1. Dedicated metrics
2. Dedicated models
3. Dedicated data manipulation

Below is the expected output of your notebook:

| #ID | Classifier name | Modifiers | Sampling | dataset size | % imbalance | Training | | | | | | Testing | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | precision | accuracy | recall | f1 | balanced accuracy | training size | precision | accuracy | recall | f1 | balanced accuracy | testing size |
| 1 | Random Forest | N/A | N/A | 1000 | 10 | | | | | | | | | | | | |
| 2 | Random Forest | class_weight | N/A | 1000 | 10 | | | | | | | | | | | | |
| 2 | Random Forest | class_weight | oversampling | 1000 | ... | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |
| | | | | | | | | | | | | | | | | | |

A few sentences of interpretation are also expected.

# Datasets – synthetic data

Let's create and plot several toy datasets that we will use for this course.

Let's try 2D/10D/20D data, two classes, 1-2-5-10-20% imbalance, size 10000.

See notebook **00_imbalanced_synthetic.ipynb**



```
Install imbalanced-learn:
https://imbalanced-learn.org/stable/install.h
tml

Create a Jupyter notebook
Use imbalance learn or sklearn to build and
plot some imbalanced datasets

from sklearn.datasets import
make_classification

from imblearn.datasets import make_imbalance
```

# Splitting

Let's split our datasets and keep an validation set on the side. It should be stratified correctly.



Full dataset          Hold-out for validation
                              (in orange)



In your notebook, create a validation set that we will call hold-out set for each of the synthetic datasets

Use Stratify in train_test_split.
(We will also use Stratification in K fold validation when training)

# Baseline model

Then we can work on a baseline model, a simple classification algorithm such as Random Forests.

sklearn doc: https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html



```
Write a Random Forest classifier for all
datasets, it's ok to keep the default
hyperparameters for now.
Use stratified k-fold cross-validation.
If you like you can also write a Support
Vector Machine (SVC in sklearn).
```



**Random Forest Simplified**





Full dataset    Fold 1    Fold 2    Fold 3    Fold 4    Fold 5

# Training the baseline models

Let's train it, and measure basic performance on our validation set with a confusion matrix
We have TP, TN, FP, FN.



```
my_notebook.ipynb

Code          Python 3 (ipykernel)

[ ]:
Train your baseline models and compute basic
confusion matrix metrics on train set and on
hold-out set.
```

| | | Predicted condition | |
|---|---|---|---|
| Total population = P + N | | Positive (PP) | Negative (PN) |
| **Actual condition** | Positive (P) | True positive (TP), hit | False negative (FN), type II error, miss, underestimation |
| | Negative (N) | False positive (FP), type I error, false alarm, overestimation | True negative (TN), correct rejection |

# Classic Metrics

We can compute the usual metrics stemming from a confusion matrix: precision, accuracy, recall.

Let's see how it goes for all our datasets. Typically precision and accuracy are very high but do not reflect at all what we want.

You get a high accuracy just by predicting the majority class, but you fail to capture the minority class, which is most often the point of the question.

sklearn doc:
https://scikit-learn.org/stable/modules/model_evaluation.html

Wikipedia:
https://en.wikipedia.org/wiki/Confusion_matrix

relevant elements

false negatives | true negatives

true positives | false positives

retrieved elements

**precision** or **positive predictive value** (PPV)

$$PPV = \frac{TP}{TP + FP} = 1 - FDR$$

**accuracy** (ACC)

$$ACC = \frac{TP + TN}{P + N} = \frac{TP + TN}{TP + TN + FP + FN}$$

**sensitivity**, **recall**, **hit rate**, or **true positive rate** (TPR)

$$TPR = \frac{TP}{P} = \frac{TP}{TP + FN} = 1 - FNR$$

|  | Real class 1 | Real class 0 |
|---|---|---|
| Predicted class 1 | 2 | 0 |
| Predicted class 0 | 8 | 100 |

True Positives = 2
False Positives = 0
True Negatives = 100
False Negatives = 8

How many retrieved items are relevant?

$$Precision = \frac{}{}$$

How many relevant items are retrieved?

$$Recall = \frac{}{}$$

$$Accuracy = \frac{True\ Positives + True\ Negatives}{Total\ \#\ samples}$$

$$Accuracy = \frac{2 + 100}{110}$$

Accuracy = 0.92

$$Recall = \frac{True\ Positives}{True\ Positives + False\ Negatives}$$

$$Recall = \frac{2}{2 + 8}$$

Recall = 0.2

$$Precision = \frac{True\ Positives}{True\ Positives + False\ Positives}$$

$$Precision = \frac{2}{2 + 0}$$

Precision = 1

$$F1\text{-}Score = \frac{2 \times (Recall \times Precision)}{Recall + Precision}$$

$$F1\text{-}Score = \frac{2 \times (0.2 \times 1)}{0.2 + 1} = \frac{0.4}{1.2}$$

F1-Score = 0.33

# Alternative metrics more suited to imbalance

Instead of the regular classification metrics, you can use more suited metrics.

F1 score is the harmonic mean of precision and recall.

F-beta score strikes a balance between precision and recall.
If we want to prioritize precision we can use beta=0.5
If we want to prioritize recall we can use beta=2

Balanced accuracy is also a better choice as it takes into account the relative size of the classes.

Let's calculate it for all our datasets.

**F1 score**

is the harmonic mean of precision and sensitivity:

$$F_1 = 2 \times \frac{PPV \times TPR}{PPV + TPR} = \frac{2TP}{2TP + FP + FN}$$

$$F_\beta = (1 + \beta^2) \cdot \frac{precision \cdot recall}{(\beta^2 \cdot precision) + recall}.$$

**specificity, selectivity** or **true negative rate** (TNR)

$$TNR = \frac{TN}{N} = \frac{TN}{TN + FP} = 1 - FPR$$

**balanced accuracy (BA)**

$$BA = \frac{TPR + TNR}{2}$$

my_notebook.ipynb

Code    Python 3 (ipykernel)

[ ]:

Compute all metrics for your baseline model!

# Tweaks to classification models to handle imbalance

Some classification models can behave a bit better than others with imbalanced data, through the use of "class_weight" parameter (Decision Tree, Random Forest, SVC)

The "balanced" mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as
`n_samples / (n_classes * np.bincount(y))`

sklearn doc:
https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html



```
my_notebook.ipynb        ●    +

[ ]:

Try class_weight parameter in Random Forest
See influence on the metrics
```

# Alternative classification models more suited to imbalance

Besides the metrics change you can also use different machine learning models, that are more suited to imbalance.

For instance **Balanced Random Forests** is a variation on Random Forests that deals with class imbalance natively.

For each iteration of RF, take a bootstrap of minority class, then take the same number of observations in majority class

Reference: Chao **Chen**, Andy Liaw, Leo Breiman, and others. Using random forest to learn imbalanced data. University of California, Berkeley, 110(1-12):24, 2004.

Imbalanced-learn ref:
https://imbalanced-learn.org/stable/references/generated/imblearn.ensemble.BalancedRandomForestClassifier.html



Imbalanced train set

Train weak classifiers — Bootstrap sample 1 — Bootstrap sample 2 — Bootstrap sample N

Aggregate predictions — Ensemble classifier

Bootstrap samples examples:

❌ Imbalanced          ✔️ More balanced

# Alternative classification models more suited to imbalance

```
my_notebook.ipynb          ●    +
💾  +  ✂  📋  📋  ▶  ■  ⟳  ⏩  Code     ∨  🕐  git  🐞  Python 3 (ipykernel)

[ ]:                                           📋  ↑  ↓  ⬆  ⬇  🗑

Try balanced random forests from
imbalanced-learn

See influence on the metrics
```

If you are in a strongly imbalanced case, you can decide to use one-class classification models such as isolation forests or single class svm.
You can also use anomaly detection techniques.

# Change the input data

Then, a powerful method to address class imbalance is to use undersampling and oversampling techniques.
If the minority class has a sufficient number of representants, then you can try undersampling the majority class at random.

Otherwise, you must find clever ways to oversample the minority class. You can simply duplicate data.

You could also try to add noise to members of the minority class.



```
Use imbalanced-learn to randomly undersample
and oversample.

See influence on the metrics
```

# SMOTE

Finally, you could use **SMOTE** -
Synthetic Minority Oversampling TEchnique.

SMOTE allows you to interpolate between members of your minority class in order to create additional data points.

For each sample in the minority class:
Compute its K-Nearest Neighbors (KNN).
Select one of them randomly.
Synthesize a new observation by linear interpolation.

Source of schema -
https://www.researchgate.net/publication/287601878_A_Novel_Boundary_Oversampling_Algorithm_Based_on_Neighborhood_Rough_Set_Model_NRSBoundary-SMOTE
Source of algorithm
https://www.cs.cmu.edu/afs/cs/project/jair/pub/volume16/chawla02a-html/node6.html





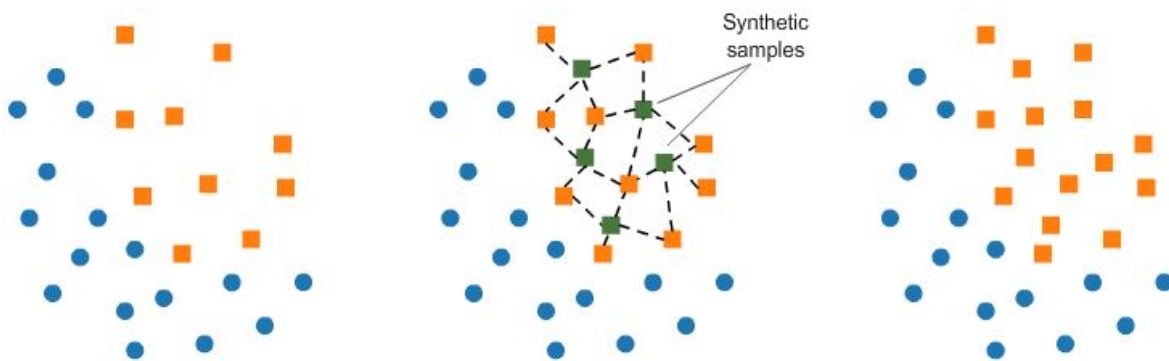Majority class samples
Minority class samples
Synthetic samples

# SMOTE

Let's use SMOTE on our toy datasets and plot them.
Then let's learn on these augmented datasets and see how the test metrics change.

```
Use imbalanced-learn to carry out SMOTE.
Plot datasets.

See influence on the metrics.
```
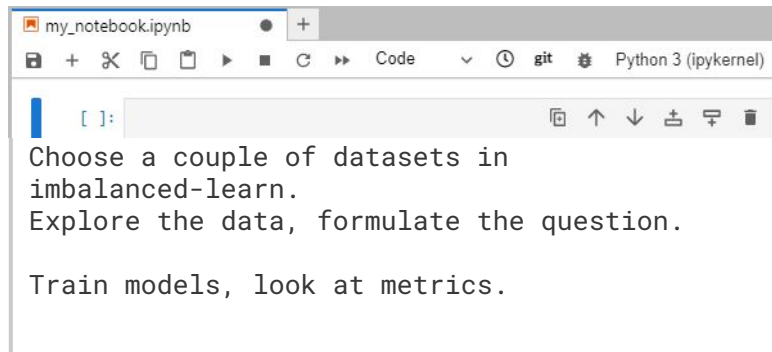


Synthetic samples

# Datasets – real data

There are several public datasets that exhibit imbalance, now that we went all the way with synthetic data, we can use them too.

A few of them are available in **imbalanced-learn** directly.

https://imbalanced-learn.org/stable/references/generated/imblearn.datasets.fetch_datasets.html

See notebook
**01_imbalanced_public_dataset.ipynb**



my_notebook.ipynb

```
Choose a couple of datasets in
imbalanced-learn.
Explore the data, formulate the question.

Train models, look at metrics.
```

# SECOM dataset

https://archive.ics.uci.edu/ml/datasets/SECOM

A complex modern **semi-conductor manufacturing process** is under consistent surveillance via the monitoring of signals/variables collected from sensors and/or process measurement points. The dataset presented in this case represents a selection of such measurements where each example represents a single production entity with associated measured features and the labels represent a simple **pass/fail yield** for in house line testing, where –1 corresponds to a pass and 1 corresponds to a fail and the data time stamp is for that specific test point.

**Quality control question.**

See notebook **02_imbalanced_secom.ipynb**

my_notebook.ipynb

Code   Python 3 (ipykernel)

```
[ ]:
```

```
Download the data.

Explore the data, formulate the question.

Choose one model, metric, and data
augmentation technique.
Train your model with cross-validation, look
at metrics.

import pandas as pd
dataset_url =
"http://archive.ics.uci.edu/ml/machine-learni
ng-databases/secom/secom.data"
label_url =
"http://archive.ics.uci.edu/ml/machine-learni
ng-databases/secom/secom_labels.data"
data_df = pd.read_csv(dataset_url, sep=" ",
header=None, names=[f"F{i+1:03d}" for i in
range(590)])
label_df = pd.read_csv(label_url, sep=" ",
header=None, names=["Label", "Time"],
parse_dates=["Time"], dayfirst=True)
  df = pd.concat([label_df, data_df], axis=1)
```

# NASA Turbofan dataset

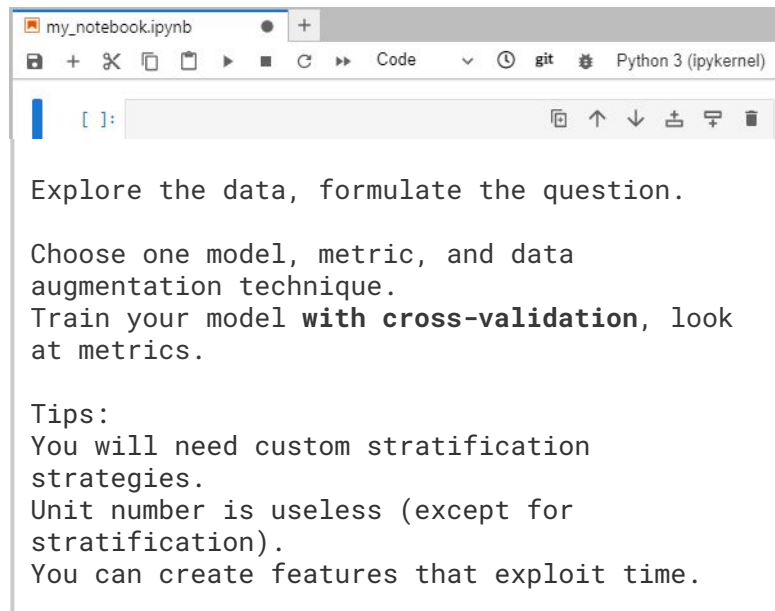https://www.kaggle.com/datasets/behrad3d/nasa-cmaps

Use only FD001

Dataset presents Run-to-Failure simulated **data from turbo fan jet engines**. It consists of multiple multivariate time series. Each time series is from a different engine. Each engine starts with different degrees of initial wear and manufacturing variation which is unknown to the user.
The engine is operating normally at the start of each time series, and develops a fault at some point during the series. In the dataset, the fault grows in magnitude until system failure.

**We want to predict whether the unit will fail within the next 5 cycles.**

See notebook **03_imbalanced_turbofan.ipynb**



Explore the data, formulate the question.

Choose one model, metric, and data augmentation technique.
Train your model **with cross-validation**, look at metrics.

Tips:
You will need custom stratification strategies.
Unit number is useless (except for stratification).
You can create features that exploit time.

# Conclusion and perspectives

Data imbalance is frequent and you can deal with it in several ways.
You can :
- change your data,
- use appropriate models,
- use appropriate metrics.
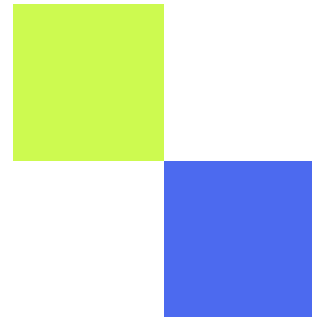
More about SMOTE:
- There are other approaches that can be used for regression questions, in particular one is SMOTER.
- At Fieldbox we also studied an algorithm for sequence-to-sequence questions, which is called SMOTEST.
- SMOTE has been discussed in a recent paper on medical trials and seem to be counterproductive sometimes.

Ideas of complementary work:
- Find other datasets and work with them
- Study SMOTER and implement it in Python
- Study this paper that tempers SMOTE:
van den Goorbergh, Ruben, et al. "The harm of class imbalance corrections for risk prediction models: illustration and simulation using logistic regression." Journal of the American Medical Informatics Association (2022).
- Watch this video of Guillaume Lemaitre at euroscipy2023 about some advanced scikit-learn features, some of which deal with imbalance.

# FIELDBOX

# Thank you

# ENSEIRB–MATMECA

8 Nov 2024