

Projet : Analyse de cibles – import, nettoyage, exploration et storytelling

🎯 Objectif global

Vous êtes un attaquant (seul ou en binôme).

Votre mission est d'explorer, de nettoyer et d'analyser les données publics de plusieurs personnes afin de concevoir une attaque de phishing ciblée réaliste. Ces données ont été récupérées lors d'attaque précédente réalisées en amont. Cependant elles ont été réalisée de façon totalement aléatoire et l'objectif ici est de trouver des corrélations entre ces données afin d'entrainer un modèle d'attaque par la suite. Le but n'est donc pas de créer l'attaque en question, mais de **comprendre et justifier** les comportements observés à travers une analyse complète, personnelle et critique des données. En effet, après avoir analysé et trouvé une corrélation entre les différents fichiers, vous devez m'expliquer comment vous attaqueraient une personne en particulier (**StoryTelling**). Afin de me convaincre dans votre choix d'attaque, vous devrez me justifier ce choix par des graphiques et des argumentaires ainsi qu'une vidéo expliquant votre raisonnement.

Pour cela, vous disposez du fichier contenant les résultats des attaques:

results.csv comprend les colonnes suivantes :

- **gaming_interest_score** : correspond au taux d'intérêt d'une personne pour les jeux vidéo.
- **insta_design_interest_score** : correspond au taux d'intérêt d'une personne pour des vidéos de design.
- **football_interest_score** : correspond au taux d'intérêt d'une personne pour le football.
- **age** : correspond à l'âge de la personne.
- **support** : correspond au support utilisé pour le phishing.
- **recommended_product** : correspond au produit recommandé pour le phishing.
- **campaign_success** : indique si le phishing a été un succès ou non (attaque réussie ou non dans le cadre du test).

Ce jeu de données contient de nombreuses **erreurs, incohérences et anomalies** (valeurs aberrantes, manquantes, doublons, formats erronés, etc.). Votre mission est d'en extraire une compréhension pertinente des comportements de chaque utilisateur.

Cette analyse peut vous permettre par la suite d'entraîner un modèle de recommandation / de scoring d'attaque (bonus du tp).

Notes: exercice pédagogique de sensibilisation au phishing et au traitement de données. L'intégralité du jeu de données est synthétique et a été créée spécifiquement pour cet atelier – aucune collecte de données personnelles réelles n'a eu lieu. Utilisation : analyses, simulations consenties.

1. Exploration et compréhension

1. Analyse exploratoire des données (Exploratory Data Analysis - EDA)

Commencez par une **analyse libre et approfondie** de vos données. Cette phase vous permet de **comprendre la structure, la qualité et les comportements du chaque personnes avant toute modélisation**.

Étapes suggérées

1. Importation et inspection

- Charger vos fichiers CSV.
- Vérifier la **qualité des données** : types de colonnes, etc ...
Optimiser la mémoire lors de l'importation (types appropriés, float32...).

Dans le **rappor**t, indiquez et justifiez clairement :

- Comment vous avez optimiser les performances lors de l'import?
- Ce que vous pensez de premier abord sur la qualité de la donnée et de sa pertinence ?

2. Nettoyage et mise en forme

- Vérifier la cohérence des données (aucune informations manquante, etc).
- Gérer les valeurs manquantes : suppression ou imputation justifiée.
- Déetecter et corriger les doublons.
- Transformer les colonnes si nécessaire (Date → datetime, score→ float).

Dans le **rappor**t, indiquez et justifiez clairement :

- Quelles valeurs ont été supprimées ou transformées ?
- Pourquoi ces choix ont été faits ?
- L'impact de votre nettoyage sur la mémoire (avant/après).



2. Détection d'anomalies

Après avoir nettoyé vos données, vous devez effectuer une **analyse des anomalies** afin d'identifier et d'écartez les points atypiques susceptibles de fausser l'attaque.

L'objectif est de **déetecter les comportements anormaux - par exemple un score de gaming trop faible ou trop élevé par rapport aux autres etc ...**

Étapes suggérées:

Proposez et implémentez **votre propre méthode de détection d'anomalies**, en justifiant vos choix. Vous pouvez, par exemple, vous appuyer sur une **Approches statistiques simples** : seuils basés sur l'écart-type, le z-score ...



Visualisations suggérées :

Pour chaque type de données, on peut tracer les valeurs normales en bleu et les **anomalies** en rouge. La moyenne et les bornes peuvent être indiquées en pointillés. Il est possible de **zoomer** sur certaines zones pour examiner plus finement les anomalies.

Dans le **rappor**t, indiquez et justifiez clairement :

- La ou les méthodes pour détecter les anomalies? Et Pourquoi ces choix ont été faits ?
- Quelles sont les anomalies retenues ?



3. Phase d'analyse statistique

Cette étape vise à approfondir la compréhension du comportement de chaque personne à travers une **analyse** des données. L'objectif est d'identifier les **caractéristiques statistiques** de chaque groupe de personne avant d'aller vers une prédition d'attaque pour un groupe en particulier.

Étapes suggérées:

- **Calculer des indicateurs (KPI)** tels que le taux de réussite global, le taux de réussite par produit, le taux de réussite par segment d'intérêt, le taux de réussite par support, ou le taux de réussite par tranche d'âge. **Tracer ces KPI sous forme de graphiques** à l'aide de matplotlib.
- **Analyser la corrélation entre plusieurs KPI**, afin de comprendre les relations entre eux (par exemple, intérêt pour un segment et taux de réussite). Ces corrélations

peuvent ensuite être visualisées à l'aide de graphiques ou de matrices de corrélation.

Dans le **rappor**t, indiquez et justifiez clairement :

- Décrivez les tendances ou observations principales que vous avez identifiées.
- Soulignez les points forts et les limites de votre première approche.
- Indiquez les hypothèses ou pistes à approfondir pour l'étape suivante.

4. Datatelling et création de l'attaque:

À partir des premières analyses de vos données, rédigez une section de *data telling* montrant comment les tendances observées permettent d'anticiper les types d'attaques auxquels un groupe de population pourrait être exposé.

Votre récit doit :

- Présenter les données clés identifiées.
- Mettre en évidence les comportements ou vulnérabilités propres au groupe étudié.
- Expliquer comment ces éléments permettent de prévoir les méthodes d'attaque susceptibles de les viser.
- Justifier clairement la cohérence entre les données et les scénarios d'attaque anticipés.

Exemple (non détaillé) d'un datattelling : Les données montrent que 59 % des enfants âgés de 0 à 3 ans utilisent régulièrement des applications ou contenus inspirés de jeux vidéo populaires tels que FIFA, avec un score moyen d'intérêt avoisinant les 80 %. Par ailleurs, leur niveau de vulnérabilité apparaît élevé : la campagne de simulation de phishing menée dans ce cadre pédagogique indique que 78 % des enfants de cette tranche d'âge ont été trompés par le message testé. Le moyen de support utilisé pour lancer l'attaque sera Instagram.

Dans le **rappor**t, rédiger clairement :

- Votre Datatteling argumente à l'aide de chiffre clé et de vos courbes.

Une **vidéo explicative (3-4 min)** est également attendue, car un *data telling* est toujours plus parlant lorsqu'il est présenté de manière visuelle et concrète.

5. Bonus: Automatiser des attaques:

L'objectif est ici de concevoir un modèle statistique ou de machine Learning capable d'identifier le type de message auquel une personne serait la plus susceptible de réagir, en fonction de ses caractéristiques (scores et âge). Un tel modèle permet d'anticiper les techniques qu'un attaquant pourrait exploiter, afin de renforcer les mesures de prévention et de protection adaptées à chaque profil.

Exemple: Jean a 12 ans avec un score d'intérêt pour les jeux vidéos important. Le résultat générera sera: Créer une attaque de phishing pour le jeu Fortnite via une campagne sur Tiktok.



Rendu :

L'ensemble des livrables doit être déposés sur un dépôt GitHub PUBLIC. Il doit contenir, l'**ensemble de votre code (.py), le rapport (pdf) et la vidéo (mp4)**.