

## ■ Analyse rapide du dataset

Dataset: <https://www.kaggle.com/datasets/datasetengineer/teenage-online-behavior-and-cybersecurity-risks/data>

### █████ résumé

- collecté depuis des logs provenants d'outils d'e-safety
- 2017-2024
- Texas, Californie
- urban / suburban

### █████ features intéressantes

- Device type: le type d'appareil utilisé lors de la session
- Malware détection (Yes/No): si un malware a été détecter sur le device
  - peut être corrélér avec le type de device > système d'exploitation
- Phishing attempts: nombres de tentatives de phishing lors de la session
- Social Media Usage (Low, Medium, High): fréquence d'utilisation de média sociaux
  - définir média social (diff de réseau social ?)
- VPN Usage (Yes, No): utilisation d'un vpn lors de la session
  - peut être corrélér avec divers facteurs de risque pour déterminer si l'utilisation d'un vpn réduit les risques cyber.
- Cyberbullying Reports: nombres d'incidents lié au cyberharcèlement sur cette session
- Parental Control Alers: Nombres d'alertes déclenchés par des programmes de contrôle parental installer sur la machine.
- Firewall Logs: nombre de connections bloquée par le firewall
  - le firewall est un programme permettant de bloquer les ports d'un pc ou l'accès internet de certains programmes. Certains firewall se basent sur une liste de logiciels malvenant à bloqué tandis que le plus basique se charge de bloqué certains ports / protocoles. Une bonne pratique sur un serveur / PC est de bloqué tout les ports par défaut et d'ouvrir uniquement ceux dont on a besoin (ex, le port 22 pour ssh.)
- Login attempts: nombres de tentatives de connexions lors de la session
  - usage ?
- Download Risk (Low, Medium, High): Niveau de risque associé aux fichiers téléchargés
  - Basé sur ?
- Password Strength (Weak, Moderate, Strong): Résistance des mots de passe utilisés
  - la qualification de la résistance d'un MDP se base sur sa capacité à résister à un bruteforce sur la durée. Pour cela, on va regarder si ce dernier est bien aléatoire, qu'il ne comporte pas de pattern identifiable, pas de mots du dictionnaire, une variété de caractères, ect...
- Data Breach Notifications: Nombre d'alertes lié aux informations personnelles de cet identifiant, ex: mot de passe, address, téléphone, ...
- Online Purchase Risk: Similaire au **download risk** mais pour les achats en ligne
- Education Content Usage: Fréquence d'engagement de l'utilisateur avec du contenu éducatif
- Age Groupe (under 13, 13-16, 17-19): catégorie d'âge auquel appartient le sujet
- Geolocation: localisation de l'utilisateur
- Public Network Usage (Yes, No): Est ce que l'activité enregistrée était elle sur un réseau public ?
  - Permettrait de corrélérer le niveau de risques des différents types de réseau
- Network Type: Type de connexion utilisée. ex: wifi, 5G, ect...
- Hours Online: Total d'heure passée sur la session
- Website Visits: Nombres de sites visités lors de la session
- Peer Interactions: Nombres d'interaction en Peer-To-Peer durant la session
  - Le Peer-To-Peer est une méthode de communication où deux PC communiquent directement entre eux sans intermédiaires. ex: le protocol bittorrent
- Risky Website (Yes, No): Est ce que des visites sur des sites à risque ont été enregistrées ?
  - Source ?
- Cloud Service Usage (Yes, No): Est ce que des services dans le *cloud* ont été utilisés lors de cette session ?
  - Définir cloud
- Unencrypted Trafic (Yes, No): Est ce que du trafic réseau non chiffré a été utilisé lors de cette session.
  - La plupart du trafic sur internet est chiffré via le protocole SSL afin d'éviter qu'un acteur mal intentionné puisse avoir accès aux données de cette connection. Il est possible de piéger un utilisateur à utiliser une connection non chiffrée afin de lui subtiliser des données.
- Ad Clicks (Yes, No): Est ce que des publicités ont été cliquées lors de cette session ?
- Insecure Login Attempts: Nombre de tentatives de connexions sur des réseaux non sécurisé.