

Hackathon  
QMI/LFIS/SESAMm  
11-12 Mars 2022

**Sujet 1 : Construction de portefeuille long-short robuste et performant**

Le but du projet est de construire le portefeuille Long-Short avec la meilleure performance absolue sous contrainte de drawdown maximum.

**Objectif :** Trouver l'allocation avec le meilleur ratio Performance absolue divisée par Drawdown maximum sur un an roulant. Le drawdown maximum est la performance entre le point haut et le point bas sur l'historique d'une allocation.

**Data :** Les données disponibles sont les suivantes :

- données de marché : Des données observables comme le prix, le beta ou encore la volatilité réalisée
- données de sentiment : Des données exploitant les travaux de Sesamm sur le sentiment des articles traitant d'une société.
- données ESG : Des données exploitant les travaux de Sesamm sur l'aspect environnemental, social et gouvernance d'une société .
- données académiques : Des données exploitant les travaux de LFIS sur les différentes primes académiques

Vous pourrez soumettre vos prédictions durant la première phase de la compétition et recevoir automatiquement la performance absolue et le drawdown maximum de votre allocation ainsi que leur ratio (voir ci-dessous le processus d'évaluation).

**Matériel :** Durant la phase 1 vous disposerez des fichiers suivants :

- x\_train\_init.csv : les données avec le return associé sur la période d'entraînement.
- x\_valid\_init.csv : les données sans le return associé sur la période de validation.
- x\_exemple\_init.csv : Un exemple de poids à envoyer au bot.

Vous pourrez soumettre au bot des allocations sur la période de validation en nombre illimité. Pour la phase 1, chaque poids doit être 1% ou 0% ou -1%. Ne peuvent être traitées pour un jour donné que les sociétés entrant dans la composition de l'indice.

**Évaluation :** Le fichier `x_test_init.csv` contient les données sur les observations de test. Vous ne pourrez pas recevoir votre score automatiquement sur ces échantillons durant la première phase de la compétition. Les organisateurs préviendront samedi matin les participants, sur le chat Telegram du Hackathon, pour annoncer la fin de la phase de validation et le début de la phase de test. Vous pourrez alors soumettre vos prédictions sur les échantillons de test, mais ne recevrez pas votre score. Il sera révélé à l'issue des soumissions de la totalité des équipes.

**Phase finale :** Après le déjeuner, les 4 meilleures équipes seront sélectionnées. L'épreuve concerne un autre univers géographique et les poids pourront varier entre 1% et 5% pour les poids positifs, -1% et -5% pour les poids négatifs. Les équipes disposeront alors des fichiers ci-dessous et ne pourront faire que 5 soumissions de validation. Un peu avant la fin de l'épreuve, elles auront deux tentatives pour contribuer sur le set de test et la dernière sera sélectionnée.

**Matériel :** Durant la phase finale vous disposerez des fichiers suivants :

- `x_train_final.csv` : les données avec le return associé sur la période d'entraînement.
- `x_test_final.csv` : les données sans le return associé sur la période de test.
- `x_valid_final.csv` : les données sans le return associé sur la période de validation.

**Features :**

**Données de marché :**

- **Date :** Un repère temporel correspondant à l'observation.
- **Sousjacent :** Une transco correspondant à un sous jacent donné.
- **secteurs :** Une transco correspondant au secteur.
- **Presence Indice** Un booléen indiquant si le sous jacent appartient à l'indice pour la date donnée.
- **vol :** La volatilité réalisée sur les 252 derniers jours
- **beta :** Le beta sur les 252 derniers jours.

**Données académiques :**

- **indic\_mom\_eu :** Un indicateur entre -1 et 1 définissant l'exposition à la prime momentum.
- **indic\_qual\_eu :** Un indicateur entre -1 et 1 définissant l'exposition à la prime qualité.
- **indic\_val\_eu :** Un indicateur entre -1 et 1 définissant l'exposition à la prime value.

### Données de sentiment :

- **neutral\_weighted\_article** : La part de neutralité dans les articles traitant de la société.
- **neutral\_weighted\_entity** : La part de neutralité dans les phrases traitant de la société.
- **agreement\_entity** : Un indicateur de la dispersion entre phrases positives et phrases négatives pour une société à l'échelon phrase.
- **agreement\_average** : Un indicateur de la dispersion entre phrases positives et phrases négatives pour une société moyenné sur l'article et la phrase.
- **positive\_weighted\_article** : Un indicateur de la positivité moyenne à l'article pour une société donné.
- **positive\_weighted\_entity** : Un indicateur de la positivité moyenne à la phrase pour une société donné.
- **positive\_index\_article** : Un second indicateur de la positivité moyenne à l'article pour une société donné.
- **positive\_index\_entity** : Un second indicateur de la positivité moyenne à la phrase pour une société donné.
- **volume\_article** : Nombre d'articles pour le jour traitant d'une société donné.
- **volume\_article\_weighted** : volume\_article weighté par la la similarité moyenne pour un jour donné pour une société donnée.

### Données ESG :

- **volume\_wordcount\_environment\_aggregated** : Nombre d'articles contenant au moins un keyword du topic environment agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_governance\_aggregated** : Nombre d'articles contenant au moins un keyword du topic governance agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_risk\_other\_aggregated** : Nombre d'articles contenant au moins un keyword du topic risk\_other agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_social\_aggregated** : Nombre d'articles contenant au moins un keyword du topic social agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_all\_risks** : Nombre d'articles contenant au moins un keyword du topic all\_risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.

- **volume\_wordcount\_bankruptcy** : Nombre d'articles contenant au moins un keyword du topic bankruptcy moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_child\_labor** : Nombre d'articles contenant au moins un keyword du topic child\_labor moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_corruption** : Nombre d'articles contenant au moins un keyword du topic corruption moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_cybersecurity** : Nombre d'articles contenant au moins un keyword du topic cybersecurity moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_data\_privacy** : Nombre d'articles contenant au moins un keyword du topic data\_privacy moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_discrimination\_racism\_sexism** : Nombre d'articles contenant au moins un keyword du topic discrimination\_racism\_sexism moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_forced\_labor** : Nombre d'articles contenant au moins un keyword du topic forced\_labor moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_fraud\_embezzlement\_crime** : Nombre d'articles contenant au moins un keyword du topic fraud\_embezzlement\_crime moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_genocide** : Nombre d'articles contenant au moins un keyword du topic genocide moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_governance** : Nombre d'articles contenant au moins un keyword du topic governance moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_health\_hazards** : Nombre d'articles contenant au moins un keyword du topic health\_hazards moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_human\_rights** : Nombre d'articles contenant au moins un keyword du topic human\_rights moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_human\_trafficking** : Nombre d'articles contenant au moins un keyword du topic human\_trafficking moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_insider\_trading** : Nombre d'articles contenant au moins un keyword du topic insider\_trading moyenné par la similarité moyenne pour un jour donné pour une société donnée.

- **volume\_wordcount\_kyc** : Nombre d'articles contenant au moins un keyword du topic kyc moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_lawsuit** : Nombre d'articles contenant au moins un keyword du topic lawsuit moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_money\_laundering** : Nombre d'articles contenant au moins un keyword du topic money\_laundering moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_political\_risk** : Nombre d'articles contenant au moins un keyword du topic political risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_regulatory\_risk** : Nombre d'articles contenant au moins un keyword du topic regulatory risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_scandal** : Nombre d'articles contenant au moins un keyword du topic scandal moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_sexism** : Nombre d'articles contenant au moins un keyword du topic sexism moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_social** : Nombre d'articles contenant au moins un keyword du topic social moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_tax\_avoidance** : Nombre d'articles contenant au moins un keyword du topic tax\_avoidance moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **volume\_wordcount\_whistleblower** : Nombre d'articles contenant au moins un keyword du topic whistleblower moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_environment\_aggregated** : Part négative des articles contenant au moins un keyword du topic environment agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_governance\_aggregated** : Part négative des articles contenant au moins un keyword du topic governance agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_risk\_other\_aggregated** : Part négative des articles contenant au moins un keyword du topic risk\_other agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_social\_aggregated** : Part négative des articles contenant au moins un keyword du topic social agrégé moyenné par la similarité moyenne pour un jour donné pour une société donnée.

- **negative\_wordcount\_all\_risks** : Part négative des articles contenant au moins un keyword du topic all\_risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_bankruptcy** : Part négative des articles contenant au moins un keyword du topic bankruptcy moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_child\_labor** : Part négative des articles contenant au moins un keyword du topic child\_labor moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_corruption** : Part négative des articles contenant au moins un keyword du topic corruption moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_cybersecurity** : Part négative des articles contenant au moins un keyword du topic cybersecurity moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_data\_privacy** : Part négative des articles contenant au moins un keyword du topic data\_privacy moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_discrimination\_racism\_sexism** : Part négative des articles contenant au moins un keyword du topic discrimination\_racism\_sexism moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_forced\_labor** : Part négative des articles contenant au moins un keyword du topic forced\_labor moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_fraud\_embezzlement\_crime** : Part négative des articles contenant au moins un keyword du topic fraud\_embezzlement\_crime moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_genocide** : Part négative des articles contenant au moins un keyword du topic genocide moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_governance** : Part négative des articles contenant au moins un keyword du topic governance moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_health\_hazards** : Part négative des articles contenant au moins un keyword du topic health\_hazards moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_human\_rights** : Part négative des articles contenant au moins un keyword du topic human\_rights moyenné par la similarité moyenne pour un jour donné pour une société donnée.

- **negative\_wordcount\_human\_trafficking** : Part négative des articles contenant au moins un keyword du topic human\_trafficking moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_insider\_trading** : Part négative des articles contenant au moins un keyword du topic insider\_trading moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_kyc** : Part négative des articles contenant au moins un keyword du topic kyc moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_lawsuit** : Part négative des articles contenant au moins un keyword du topic lawsuit moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_money\_laundering** : Part négative des articles contenant au moins un keyword du topic money\_laundering moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_political\_risk** : Part négative des articles contenant au moins un keyword du topic political risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_regulatory\_risk** : Part négative des articles contenant au moins un keyword du topic regulatory risk moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_scandal** : Part négative des articles contenant au moins un keyword du topic scandal moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_sexism** : Part négative des articles contenant au moins un keyword du topic sexism moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_social** : Part négative des articles contenant au moins un keyword du topic social moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_tax\_avoidance** : Part négative des articles contenant au moins un keyword du topic tax\_avoidance moyenné par la similarité moyenne pour un jour donné pour une société donnée.
- **negative\_wordcount\_whistleblower** : Part négative des articles contenant au moins un keyword du topic whistleblower moyenné par la similarité moyenne pour un jour donné pour une société donnée.