

## Article choisi: EXPLAINING AND HARNESSING ADVERSARIAL EXAMPLES

### DÉPOTS :

- Git :
  - o Code source
  - o Readme
  - o Notebook (demo execution des expériences de l'article)
- Rapport 6 à 8 pages en ANGLAIS
- Vidéo 5 min de présentation du projet

### BAREME :

- Rapport /40
- Code /30
- Soutenance /30

<b>Problème étudié</b>	Les réseaux de neurones se trompent sur des entrées légèrement perturbées.
<b>Ce qu'on pensait avant</b>	Que c'était à cause de leur non-linéarité ou de l'overfitting.
<b>Nouvelle explication</b>	En fait, c'est leur <b>nature trop linéaire</b> (à grande échelle) qui cause ces erreurs.
<b>Conséquence</b>	On peut générer facilement des exemples adversariaux en exploitant cette linéarité.
<b>Application utile</b>	En utilisant ces exemples pour l'entraînement, on améliore la robustesse du modèle.

Pas compliqué, ressources internet

Objectifs : Produire des exemples qui amènent à une mauvaise classification

➔ Gros focus sur **évaluation et résultats**

2 équipes :

- Réimplémenter l'article
  - o Générer des exemples adversarios
- Se concentrer sur l'évaluation
  - o Entraîner le modèle
  - o Récup un modèle et voir comment il marche

2sem :

- récup base de donnée et présenter, cmbn de classes, ce qu'on retient, résultats qualitatifs (brouillage qui marche et qui marche pas)
- récup ttes infos sur internet à ce sujet

Ds 1 mois :

- majorités des entraînements commencés/finis
- Majorité de l'architecture faites

On peut utiliser code ou IA génératives mais à indiquer

Pas de courbe époque =  $f(\text{loss})$  on s'en balec (jsp ce que ça veut dire mais oké), reproduire les résultats de l'article

Pour aller plus loin : Regarder comment les modèles se comportent avec d'autres bases de données