

Ma recette santé !

Étude expérimentale du transfert analogique

Balkis Bouthaina Dirahoui, Sarah Djouder, Tinhinane Chafai
Encadré par : Marie-Jeanne Lesot



2 juin 2022

Table des matières

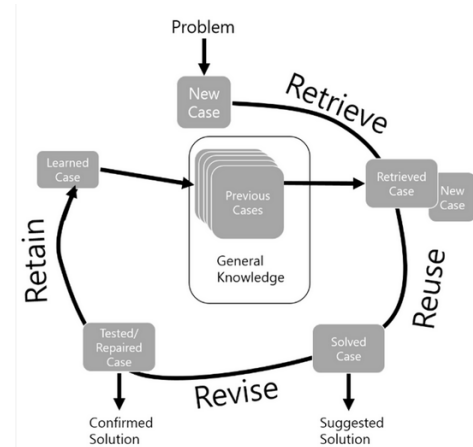
- ❶ Transfert Analogique
- ❷ Implémentation
 - Architecture globale
 - Modélisation
 - CoAT optimisé
- ❸ Expérimentations
 - Validation CCBI
 - Comparaison du temps d'exécution : CCBI vs Knn,SVM
 - Temps d'exécution : CoAT vs CCBI
 - Performance CCBI : PREC, CONF, MAE
- ❹ Application à des données réelles
 - Traitements
 - Classification
 - Régression

Introduction : Transfert Analogique

Définition :

- Méthode d'apprentissage supervisé qui permet de résoudre des problèmes de régression ou de classification
- Semblable au modèle des K plus proches voisins...
- Stockage des données d'apprentissage.
- Utilise ses données pour une nouvelle prédiction.
- Utilisation de ces données pour une nouvelle prédiction
- Travail **principalement** fait à l'inférence.

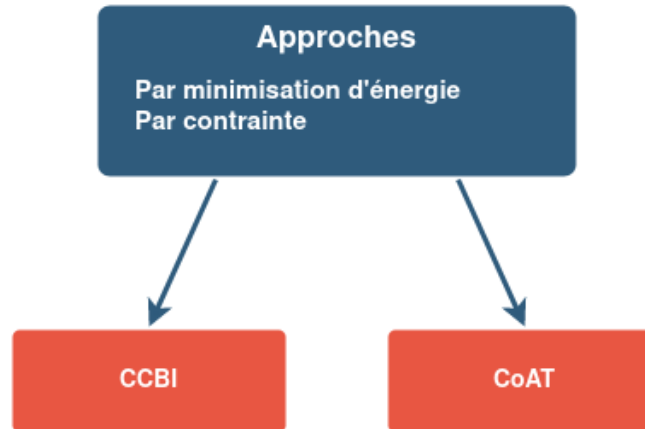
1



Introduction : Transfert Analogique

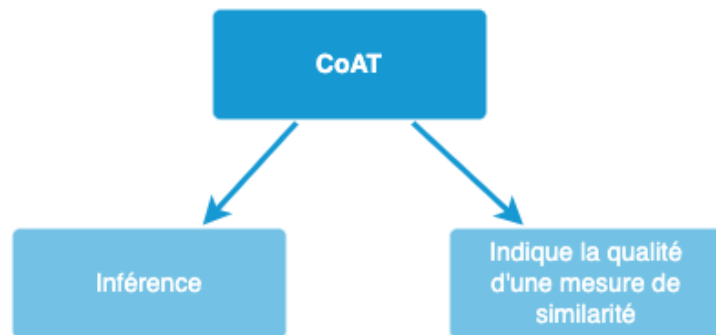
Hypothèse

- Si deux situations sont similaires dans un aspect alors elles sont aussi similaires dans d'autres aspects.



Une approche ordinale : CoAT²

- Traduit le principe d'analogie sous forme **ordinaire**.
- Introduit le concept d'**inversion**.
- Introduit le concept de **complexité**.

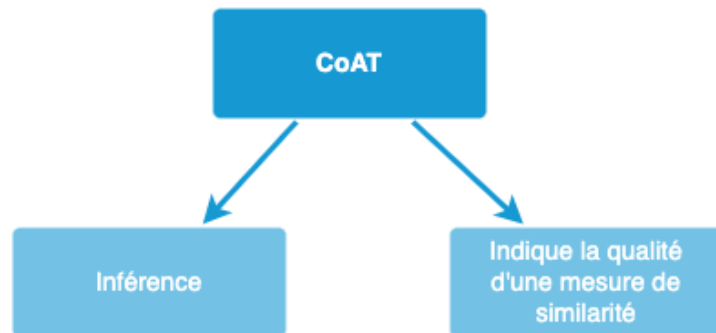


2. Fadi Badra. 2020. A Dataset Complexity Measure for Analogical Transfer. In Proceedings of the 29th International Joint Conference on Artificial Intelligence IJCAI 2020, 1601–1607.

Une approche ordinale : CoAT

- **Notations**

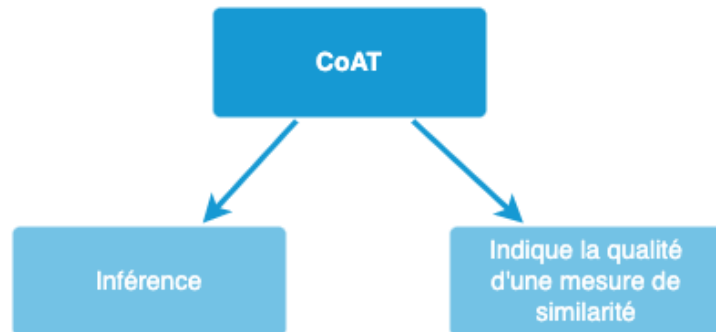
- **triplet** (s_0, s_i, s_j)
- **mesure de similarité** σ_s
- **mesure de similarité** σ_r



Une approche ordinale : CoAT

- **Notion d'inversion**

$$Inv(s_0) = \{s_i s_j | \sigma_s(s_0, s_i) \geq \sigma_s(s_0, s_j) \text{ et } \sigma_r(s_0, s_i) < \sigma_r(s_0, s_j)\}$$



Une approche ordinale : CoAT

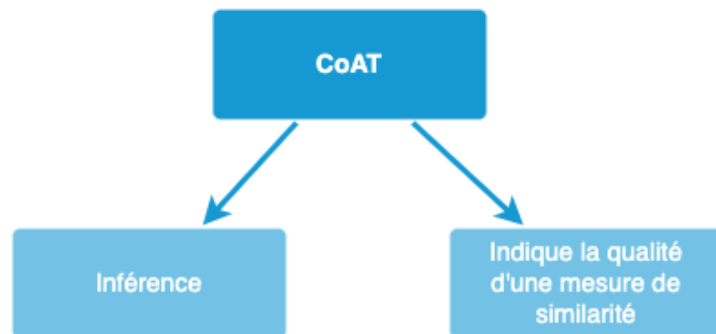
- **Cardinal**

$$\gamma(s_0) = |inv(s_0)|$$

- **Complexité Γ**

$$\Gamma(\sigma_s, \sigma_r, CB) = \sum_{s_0 \in CB} \gamma(s_0)$$

Complexité temporelle en $O(n^3)!!$

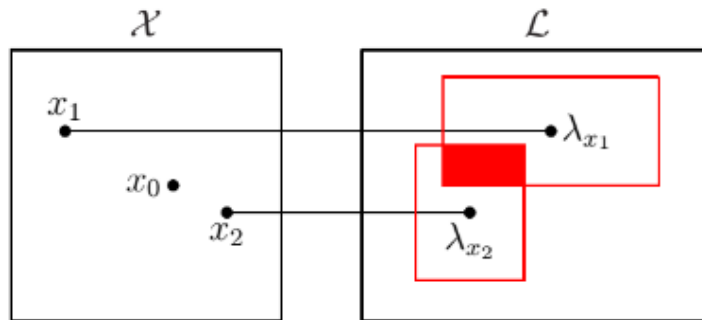


Approche par contrainte : CCBI³

- **Différence** : Pas des triplets $(s_0, s_i, s_j) \rightarrow$ **couples** (x_i, λ_{xi})
- Contrainte de similarité :

$$\forall x, y \in X : \text{sim}X(x, y) \leq \text{sim}L(\lambda_x, \lambda_y)$$

- Un cas (x_1, λ_{x1}) soit $\text{sim}X(x_1, x_0) = \alpha$ **solution** $\lambda_{x0} \in \alpha$ -voisinage de λ_{x1}

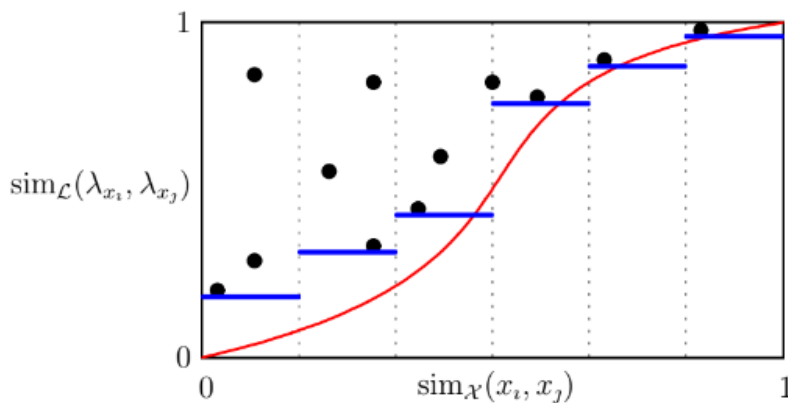


Approche par contrainte : CCBI

- Une contrainte de similarité :

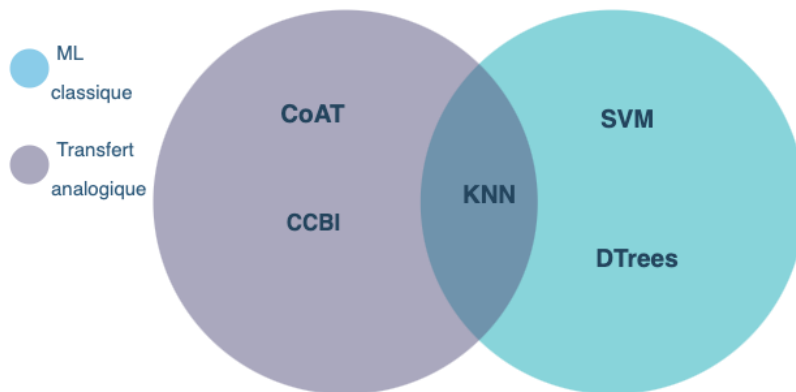
$$\forall x, y \in X : \text{sim}_X(x, y) \leq \text{sim}_L(\lambda_x, \lambda_y)$$

- Apprentissage d'une fonction $h : x \mapsto \sum_{k=1}^n \beta_k \cdot \mathbb{1}_{(A_k)}(x)$



Implémentation : Architecture Globale

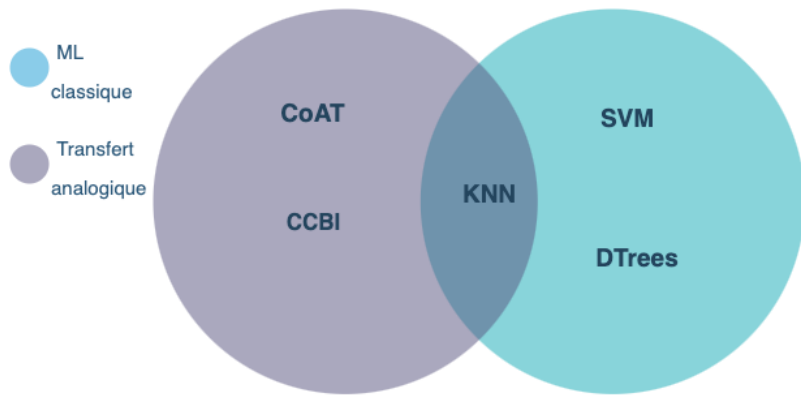
- **Méthodes disponibles** : CoAT,CCBI, KNN,SVM,DTrees...
- **Paramètres de chaque modèle** :
 - k KNN
 - β_k CCBI...
- **Mesures disponibles** :
 - SimX
 - SimL
 - σ_s
 - Distance euclidienne
 - **autre**



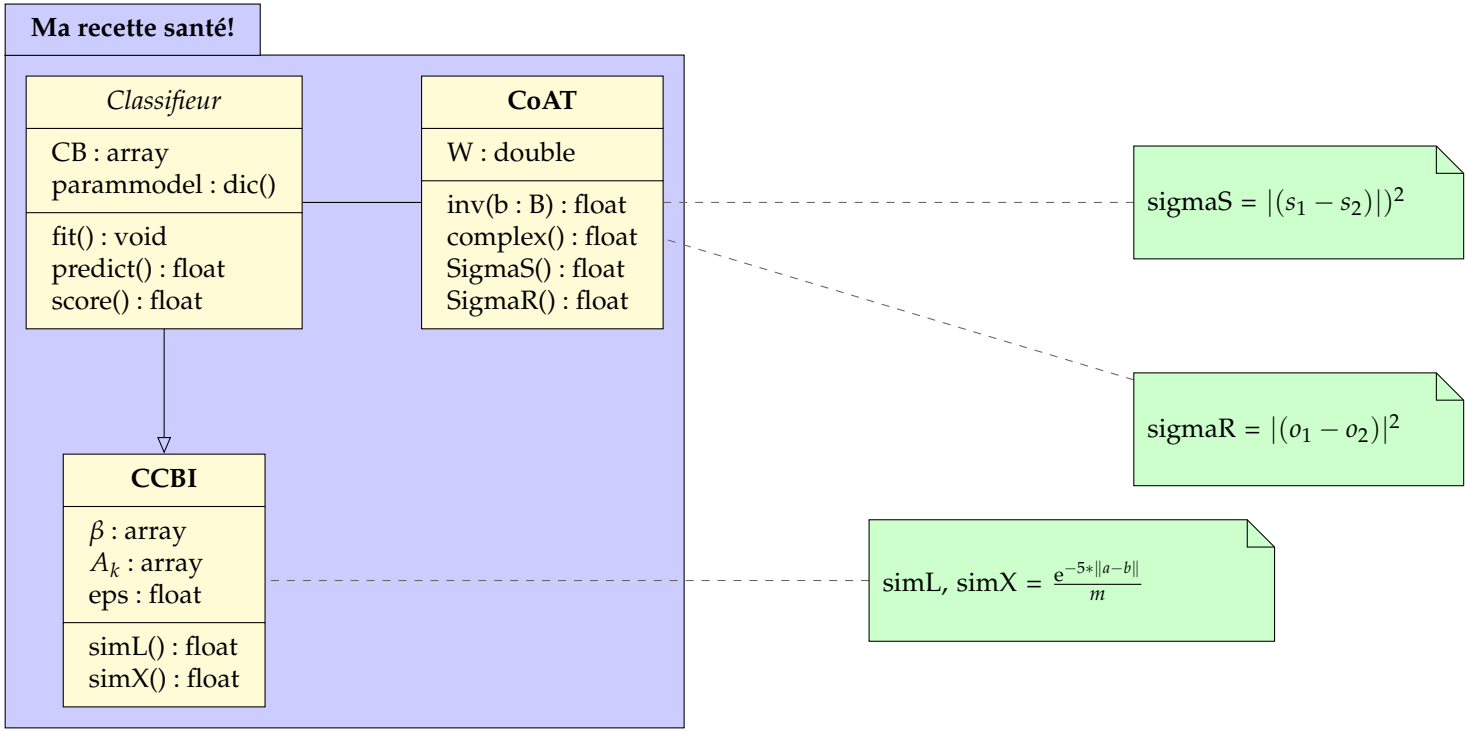
Implémentation : Architecture Globale

Pour la plate-forme de comparaison expérimentale :

- **Évaluation des performances**
 - Temps d'exécution
 - Validation croisée
 - Split test-train
- **Paramètres d'apprentissage :**
 - Nombre d'exemples en test
 - k en validation croisée.



Implémentation : Modélisation



Implémentation : CoAT optimisé

Rappelez vous...

- **Minimiser** la complexité :

$$\Gamma(\sigma_s, \sigma_r, CB) = \sum_{s_0 \in CB} \gamma(s_0)$$

Implémentation : CoAT optimisé

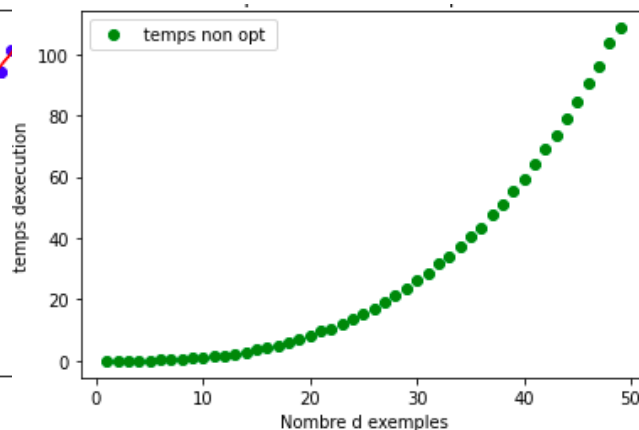
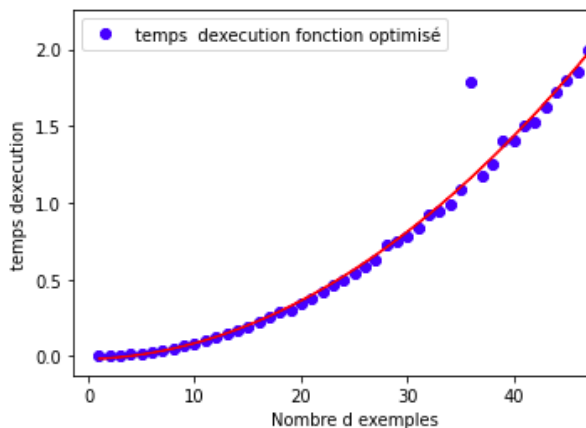
Et si ? :

- En inférence : calculer pour tout le dataset ?
- Calculer **seulement** $|inv(s_0)|$.
- Optimiser le temps de calcul..

\Rightarrow Peut être réduit de $O(n^3)$ à $O(n^2)$

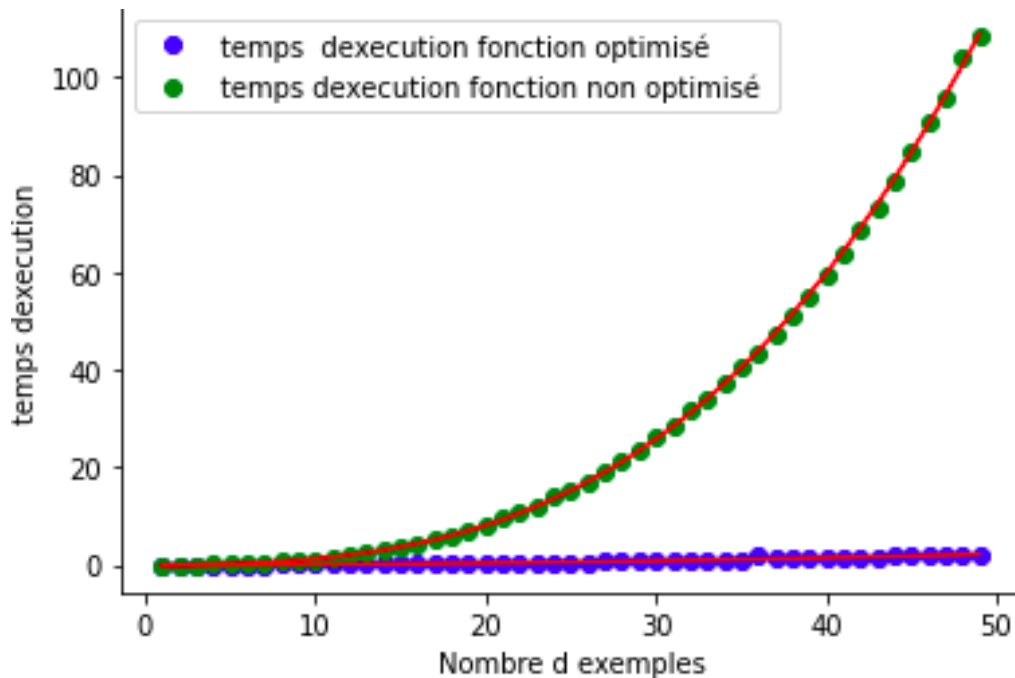
Expérimentation : Performances de la version optimisée

- Non-optimisée (vert) en $O(n^3)$
- Optimisée (bleu) en $O(n^2)$



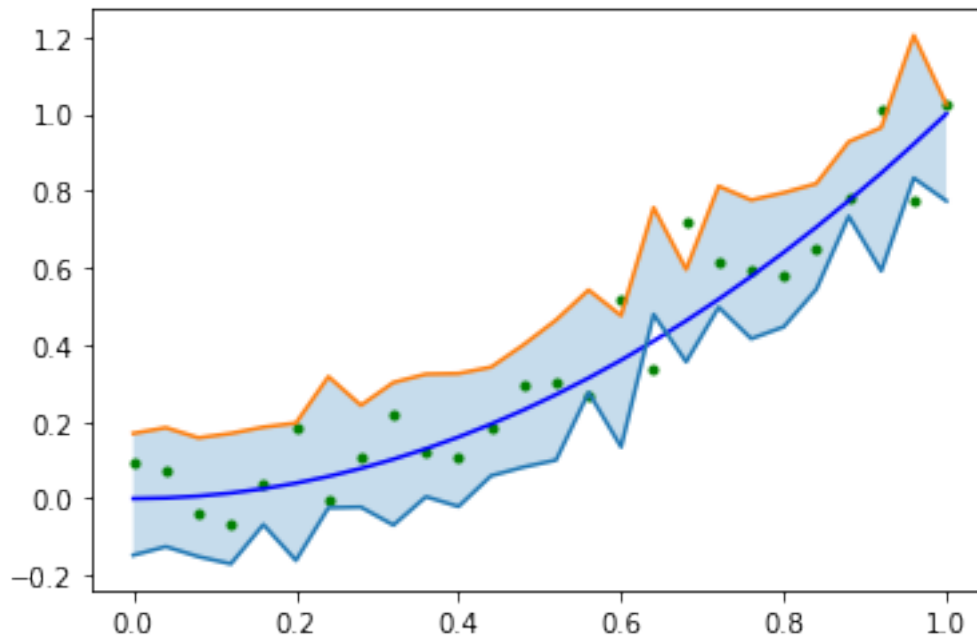
Expérimentation : Performances de la version optimisée

- Non-optimisée (vert) en $O(n^3)$
- Optimisée (bleu) en $O(n^2)$



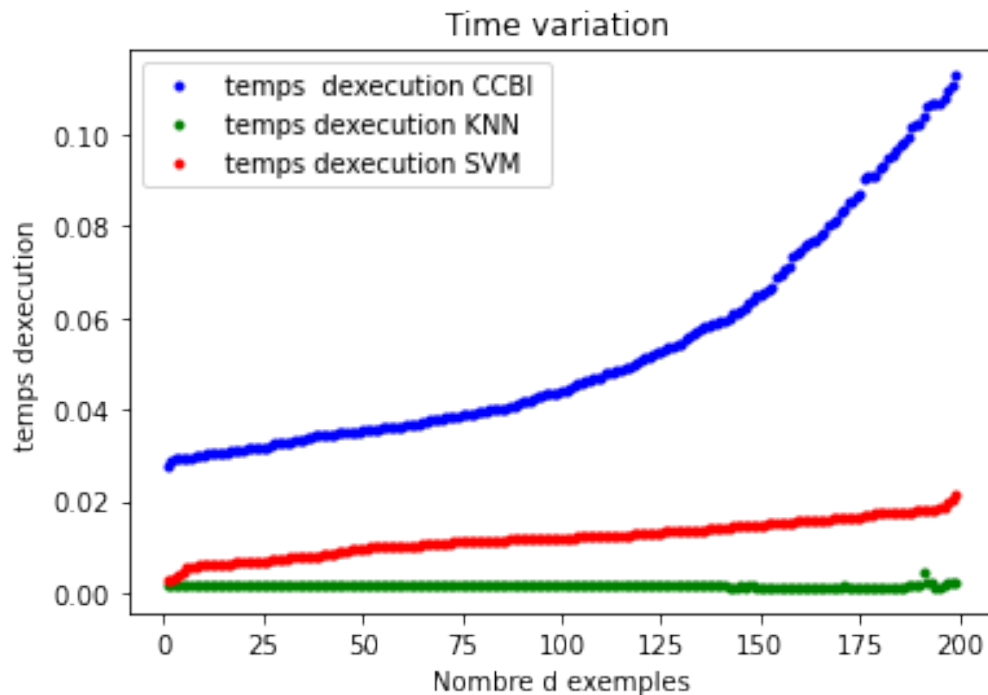
Validation des résultats CCBI

- Ligne bleue : $x \rightarrow x^2$
- Points verts : Exemple d'apprentissage bruités $x \rightarrow x^2$
- Intervalle prédit (région entre les lignes brisées)



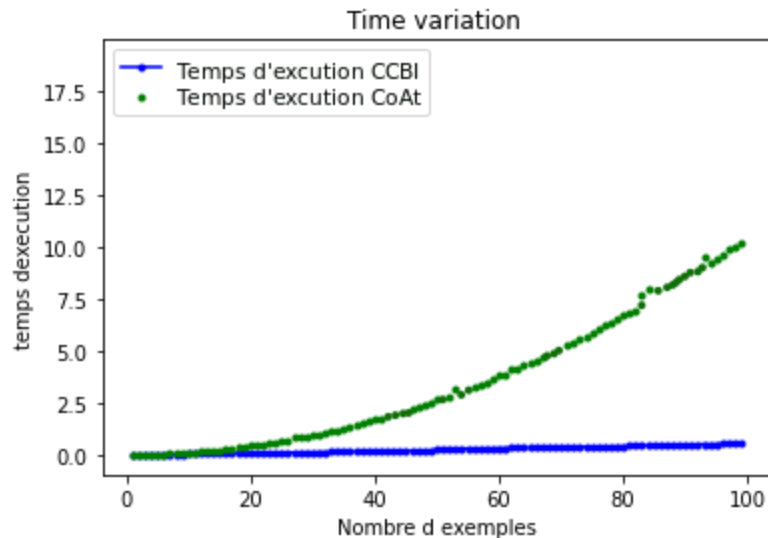
Temps d'exécution : CCBI vs Knn, SVM

- Ligne bleue : $x \rightarrow x^2$
- Temps d'exécution par rapport au nombre d'exemples



Temps d'exécution : CoAT vs CCBI

- Problème de régression : $x \rightarrow x^2$
- **Différence** : double boucle 'inv()' utilisé dans **CoAT**
- Ligne verte : CoAT
- Ligne bleue : CCBI



Performance CoAT/CCBI vs SVM, Knn, Dtrees

Modèle	SVM	Knn	Dtrees	CoAT
Score (Accuracy)	0.98	0.98	0.90	0.98

Performances des modèles sur les données IRIS

Modèle	SVM	Knn	Dtrees	CCBI
Score	0.98	0.99	0.98	0.94

Performances des modèles sur les données x^2

- Ici la mesure utilisée pour CCBI est la confiance (CONFIDENCE)
- CONF : Taux de bonne classification $y \in \text{intervalle}$.

Performance : CONF, PREC, MAE

- Validation de l’algorithme CCBI sur des données synthétiques
- **Mpg (resp. Cpu)** : Données sur les véhicules (resp. Données sur les ordinateurs) de la base UCI ⁴
- **CONF** : Taux de bonne classification
 $y \in \text{intervalle}$.
- **PREC** : Taille de l’intervalle prédit.
- **MAE** : Mean Absolute Error.

Mesure	CONF	PREC	MAE
mpg	0.94	0.06	0.01
cpu	0.98	0.35	0.15

4. <https://archive.ics.uci.edu/ml/datasets.php>

Application à des données réelles

- Données recettes de cuisine issues du site BBC Good Food
- Sous format Json

```
{
  "page": {
    "article": {
      "author": "Jane Hornby",
      "description": "Give traditional roast chicken a Chinese-style twist with this zesty recipe",
      "id": "101919",
      "tags": [],
      "recipe": {
        "collections": [
          "Broccoli",
          "Winter roasts",
          "Chinese chicken",
          "Chicken leg"
        ],
        "cooking_time": 2700,
        "prep_time": 600,
        "serves": 4,
        "keywords": [
          "Chinese",
          "Lunch",
          "Family",
          "Roast",
          "One pot",
          "Supper",
          "Super",
          "Dinner",
          "Diner",
          "BBC Good Food magazine July",
          "Broccoli",
          "Chicken leg",
          "Chicken legs",
          "Clear honey",
          "Garlic clove",
          "Garlic cloves",
          "Ginger",
          "Lemon",
          "Lemons",
          "Light soy sauce",
          "Sesame oil",
          "Sesame seed",
          "Sesame seeds"
        ],
        "ratings": 87,
        "nutrition_info": [
          "Added sugar 6g",
          "Carbohydrate 7g",
          "Kcal 340 calories",
          "Protein 28g",
          "Salt 2.35g",
          "Saturated fat 7g",
          "Fat 23g"
        ],
        "ingredients": [
          "chicken leg",
          "lemon",
          "broccoli",
          "garlic clove",
          "ginger",
          "light soy sauce",
          "clear honey",
          "sesame seed",
          "sesame oil"
        ],
        "courses": [
          "Main course",
          "Supper",
          "Lunch"
        ],
        "cuisine": "Chinese",
        "diet_types": [],
        "skill_level": "Easy",
        "post_dates": "1309474800",
        "channel": "Recipe",
        "title": "Lemon, broccoli & sesame roast chicken"
      }
    }
  }
}
```

Extrait d'une recette du Dataset Json

- Environ 3000 recettes avec 30 attributs différents pour chacune

Traitements

- Implémentation d'un parser :

- 1 Expressions régulières
- 2 Extractions des valeurs nutritionnelles .

- Ingrédients : sac de mots

- 1 Tokenization
- 2 Stemming
- 3 TF-IDF

- Binarisation des valeurs catégorielles



Jeu de Données : Caractéristiques du jeu de données

Caractéristiques du jeu de donnée

- ➊ Plus 1000 recettes
- ➋ 49 attributs différents (ingrédients,sucre,fibres,...)
- ➌ 5 tâches de classification
- ➍ 2 tâches de régression

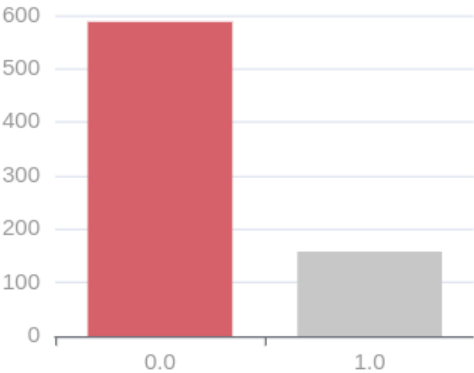
Classification : La recette est-elle saine ou non ?

Protocole Expérimental

- **Équilibrer les classes**
 - UnderSampling
- **Séparer les données**
 - train , test

A Healthy

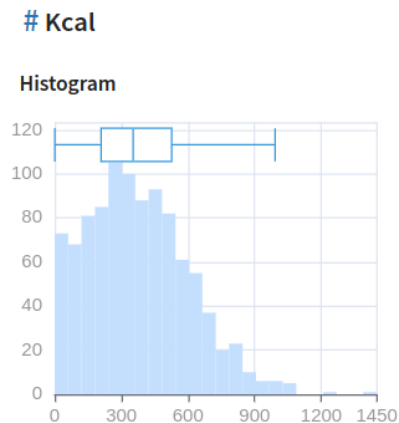
Histogram



Modèle	SVM	Knn	Dtrees	CoAT
Score (Accuracy)	0.72	0.82	0.78	0.80

Performances des modèles en classification sur les données BBC Good Food

Régression : quel est l'apport calorique de la recette ?



Distribution des calories

Modèle	LR	CCBI	CoAT
R2-Score	0.82	0.80	0.91

Performances des modèles en régression sur les données BBC Good Food

Conclusion

❶ Réalisations :

- Implémentation d'approches du transfert analogique (CCBI + CoAT)
- Optimisation théorique et pratique de CoAT
- Développement d'une plateforme de comparasion expérimentale
- Construction d'un jeu de données réelles
- Application à des données réelles

Conclusion

❶ Réalisations :

- Implémentation d'approches du transfert analogique (CCBI + CoAT)
- Optimisation théorique et pratique de CoAT
- Développement d'une plateforme de comparasion expérimentale
- Construction d'un jeu de données réelles

❷ Bilan :

- L'idée du transfert analogique est séduisante mais :
- Cout de Calcul prohibitif
- Pertinent pour des petits jeux de données.