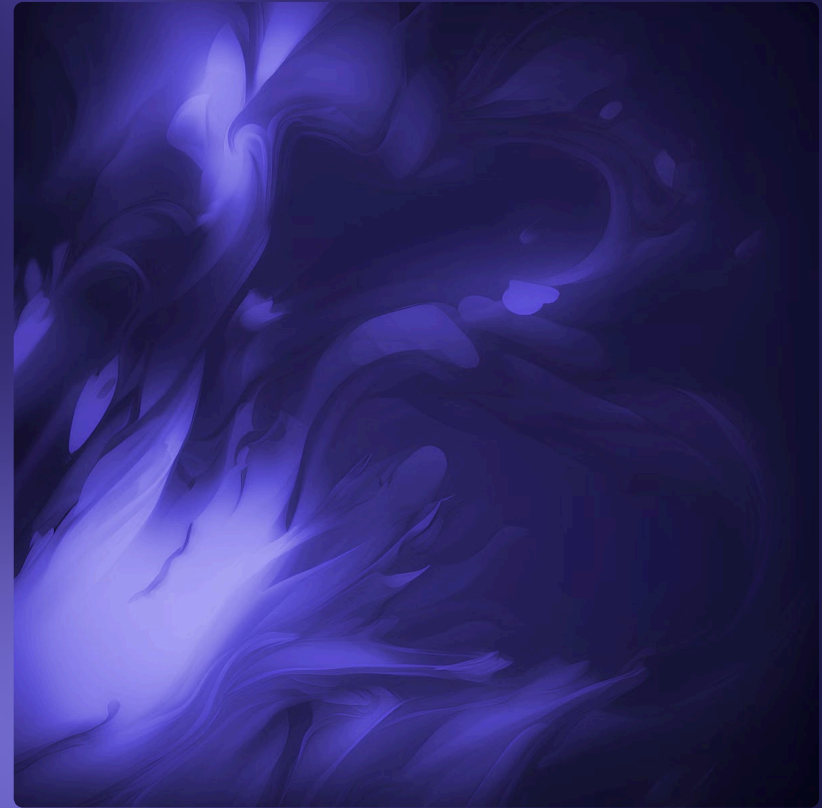


Préparation des données pour l'apprentissage automatique

La préparation des données est une étape cruciale dans tout projet d'apprentissage automatique. Bien que souvent négligée, elle peut faire la différence entre un modèle performant et un modèle médiocre. Dans cette présentation, nous explorerons les techniques clés de préparation des données pour optimiser les performances de vos modèles prédictifs.

 by etixi hacking



Pourquoi la préparation des données est-elle importante ?

Qualité des données

La qualité des données d'entrée est primordiale pour la construction d'un modèle fiable. La préparation des données permet de s'assurer que les données sont cohérentes, précises et complètes.

Performances du modèle

Les performances de votre modèle dépendent directement de la qualité des données utilisées pour l'entraîner. Une préparation adéquate des données optimisera les résultats de votre modèle.

Efficacité de l'apprentissage

Une préparation des données bien menée accélère et facilite le processus d'apprentissage automatique, permettant d'obtenir des résultats plus rapidement.

Techniques de préparation des données

Nettoyage des données

Le nettoyage des données implique d'identifier et de traiter les données manquantes, aberrantes ou en double. Cela inclut l'imputation des valeurs manquantes, la suppression des doublons et le traitement des valeurs aberrantes.

Sélection de fonctionnalités

La sélection de fonctionnalités permet d'identifier les variables les plus pertinentes pour votre modèle. Cela peut inclure des méthodes comme l'importance des fonctionnalités et la sélection de fonctionnalités récursive (RFE).

Transformations de données

Les transformations de données, telles que la normalisation, la standardisation et l'encodage, permettent de préparer vos données pour une meilleure performance du modèle.

Nettoyage des données

Suppression des doublons

Identifiez et supprimez les lignes en double dans votre jeu de données pour éviter les biais et améliorer la qualité des résultats.

Traitement des valeurs aberrantes

Détectez et traitez les valeurs aberrantes qui pourraient fausser les résultats de votre modèle. Vous pouvez les supprimer ou les remplacer par des valeurs plus représentatives.

1

2

3

Gestion des valeurs manquantes

Imputez les valeurs manquantes de manière appropriée, en utilisant des techniques comme la moyenne, la médiane ou des méthodes plus avancées comme la régression.

Sélection de fonctionnalités

1

Importance des fonctionnalités

Évaluez l'importance relative de chaque fonctionnalité pour votre modèle à l'aide de métriques comme le coefficient de corrélation ou le gain d'information.

2

Sélection de fonctionnalités réursive (RFE)

Utilisez la sélection de fonctionnalités réursive pour identifier de manière itérative les fonctionnalités les plus pertinentes pour votre modèle.

3

Sélection de fonctionnalités pour la régression

Pour les problèmes de régression, vous pouvez utiliser des méthodes comme la régression Ridge ou Lasso pour sélectionner les fonctionnalités les plus importantes.

Transformations de données



Normalisation

Normalisez vos variables pour les ramener à la même échelle, facilitant ainsi l'apprentissage de votre modèle.



Standardisation

Standardisez vos variables en les centrant sur leur moyenne et en les réduisant à l'unité, ce qui améliore la stabilité numérique.



Encodage One-Hot

Transformez vos variables catégorielles en vecteurs binaires pour les rendre compatibles avec la plupart des algorithmes d'apprentissage automatique.



Transformations de puissance

Appliquez des transformations de puissance, comme la transformation de Box-Cox, pour corriger les problèmes de non-normalité des distributions.

Réduction de dimensionnalité

1

Analyse en composantes principales (PCA)

Utilisez la PCA pour réduire le nombre de fonctionnalités tout en conservant la majorité de l'information contenue dans les données.

2

Analyse discriminante linéaire (LDA)

Appliquez la LDA pour projeter vos données dans un espace de plus faible dimensionnalité, tout en maximisant la séparabilité entre les classes.

3

Autres méthodes

D'autres techniques comme l'Analyse en composantes indépendantes (ICA) ou l'Apprentissage de caractéristiques (Representation Learning) peuvent également être utilisées pour réduire la dimensionnalité.

Mieux vaut prévenir que guérir

| | |
|-------------------------|---|
| Comprenez vos données | Analysez en détail la structure, les statistiques et les relations de vos données. |
| Documentez le processus | Consignez chaque étape de la préparation des données pour faciliter la reproductibilité et le suivi. |
| Validez les résultats | Vérifiez que les transformations appliquées ont bien l'effet escompté sur les performances de votre modèle. |
| Soyez vigilants | Restez attentifs aux biais et aux hypothèses sous-jacentes qui pourraient affecter la qualité de vos données. |