

Blocks Classification

THOMAS MOULY

AYMERIC LEBOUCHER

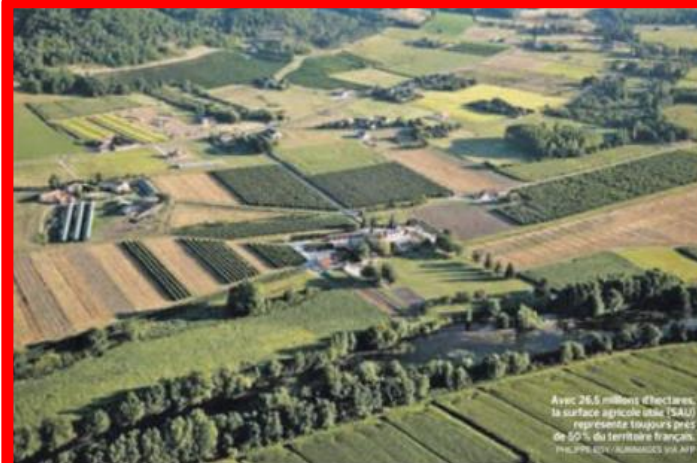
Le contexte de l'étude

- Domaine : Analyse automatique de document
- Filère : Segmentation de bloc
- Problématique : Classifier la nature d'un bloc en fonction de paramètre graphiques, géométriques.
- Chaîne de protocole: Document Analysis -> Bloc segmentation -> **Bloc classification** -> Text extraction
- Intérêt: Extraire l'information textuelle, conditionner un algorithme de data mining
- Documents cibles : Formulaires, chèque bancaire

Image

vendredi 10 décembre 2021 LE FIGARO

30 ÉCONOMIE



Avec 26,5 millions d'hectares, la surface agricole totale (SAU) représente toujours près de 50 % du territoire français. (PHILIPPE L. GOURMAUD VIA AFP)

En dix ans, 20 % des fermes de France ont disparu

Si la surface agricole s'est maintenue, les exploitations ont grandi. Le recensement décennal révèle une profonde mutation du secteur.

OLIVIA DETROYAT @OliviaDet

AGRICULTURE Si la pandémie et les confinements ont remis sur le devant de la scène le rôle des agriculteurs, le profil de la « ferme France » et de ceux qui y travaillent est en profond bouleversement. C'est ce que révèle le grand recensement agricole décennal, réalisé par le ministère de l'Agriculture entre mars et octobre, dont le Figaro dévoile les enseignements.

Le principal constat confirme une tendance engagée il y a cinquante ans. Il y a de moins en moins de fermes en France. Entre 2010 et 2020, l'Hexagone a vu disparaître une exploitation sur cinq. Avec 389 000 fermes recensées en 2020, 109 000 exploitations ont disparu. Cette baisse, continue depuis les années 1970, s'est un peu (-2,3 %) ralentie par rapport aux -3 % annuels enregistrés dans les an-

69

hectares
Superficie moyenne
d'une exploitation,
en hausse de 25 %
sur la décennie

sur la décennie. Le but : tenter de baisser les coûts et d'accroître les marges, pour continuer à rivaliser avec les nouveaux poids lourds mondiaux de l'agriculture. En première ligne les pays de la mer Noire dans les céréales et les élevages de l'autre côté de l'Atlantique.

Mais loin de la course au gigantisme, les campagnes françaises ont aussi accéléré la valorisation de leurs produits, notamment sous l'impulsion des États généraux de l'alimentation de 2017. Une montée en gamme symbolisée par le boom du bio, dont les surfaces ont été multipliées par trois sur la décennie, pour dépasser 12 % de la SAU tricolore. Une croissance que certaines filières estiment s'être faites au détriment d'autres démarches de valorisation (labels, AOP...), qui elles restent globalement stables (27 %) par rapport à 2010.

municants, les salariés non familiaux pesent désormais une personne sur cinq dans les fermes tricolores.

Ce chiffre confirme une autre tendance : l'ouverture croissante de la profession à de nouveaux profils. Dans les lycées agricoles comme chez les reconvertis urbains de tous âges, les « non issus du monde agricole » sont de plus en plus nombreux. Cette bonne nouvelle souligne l'attractivité retrouvée d'un métier exigeant.

Cela a permis de maintenir autour de 20 % la proportion d'exploitants de moins de 40 ans. Reste qu'avec près de 60 % des chefs d'exploitation de plus de 50 ans, le mur démographique et le défi du renouvellement des générations se manifestent chaque année avec plus d'acuité. Et il reste, avec la poursuite de la transition écologique de la France agricole, l'enjeu de la prochaine décennie. ■

Un paysage agricole bouleversé



Graphique

Exemple

Texte

Le Dataset

Toutes les observations sont complètes, il ne manque pas de valeurs

Classes: **1,2,3,4,5** respectivement "**text**", "**horiz. line**", "**graphic**", "**vert. line**", "**picture**"

5473 observations

Attributs

Attribut	Signification
Height: entier	Hauteur du block
Length: entier	Largeur du block
Area: entier.	Aire du block
Eccen: continue.	Excentricité du block (largeur / hauteur);
p_black: continue.	Pourcentage de pixels noirs sur le block
p_and: continue.	Pourcentage de pixels noirs sur le block après Run Length Smoothing Algorithm (RLSA)
mean_tr: continue.	Moyenne des transitions Blanc/Noir (blackpix / wb_trans);
blackpix: continue.	Nombre de pixels noirs sur le block original
blackand: continue.	Nombre de pixels noirs sur le block après application du RLSA
wb_trans: continue.	Nombre de transitions blanc/noir du bloc

RLSA ALGORITHM



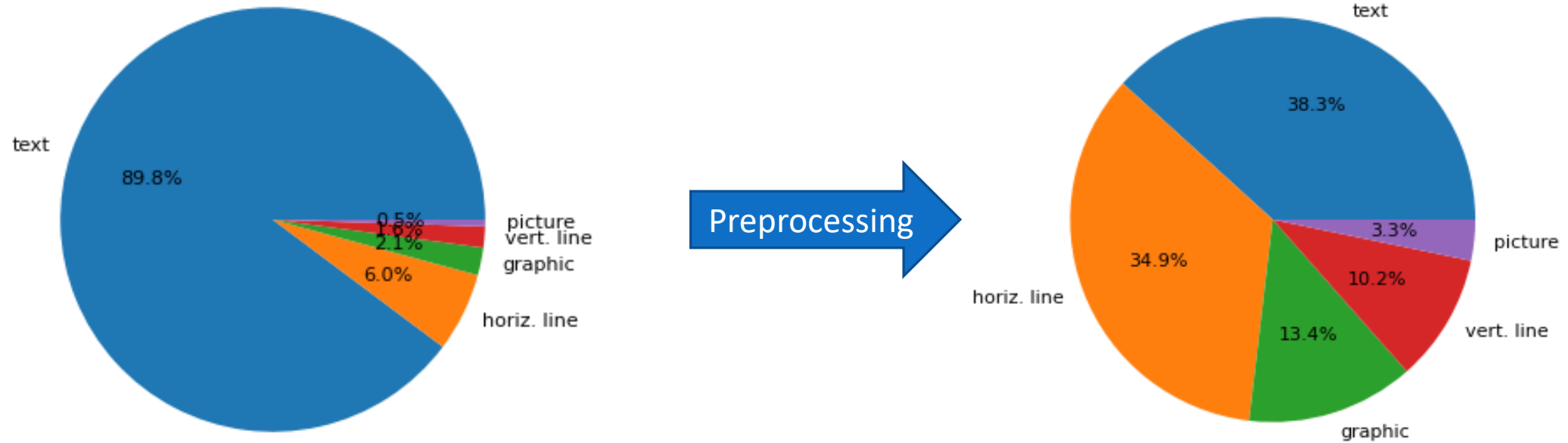
(a)



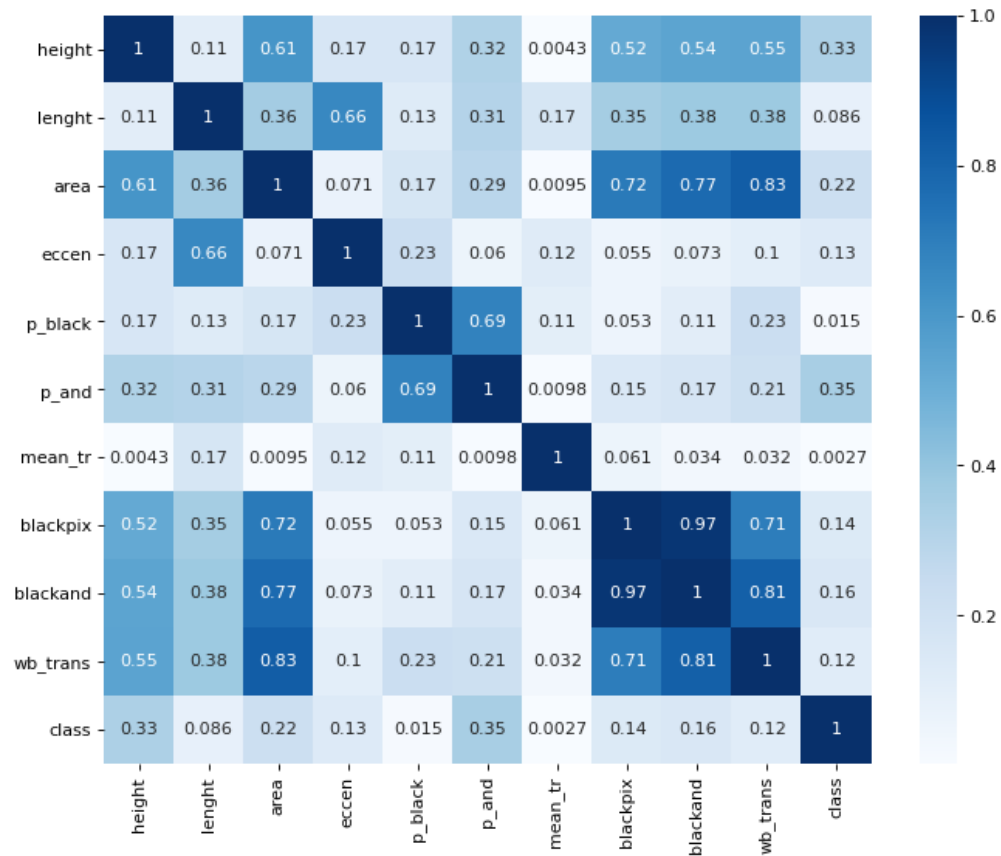
(b)

Datavisualisation

Trop forte présence d'observations de classe texte dans le dataset



Datavisualisation

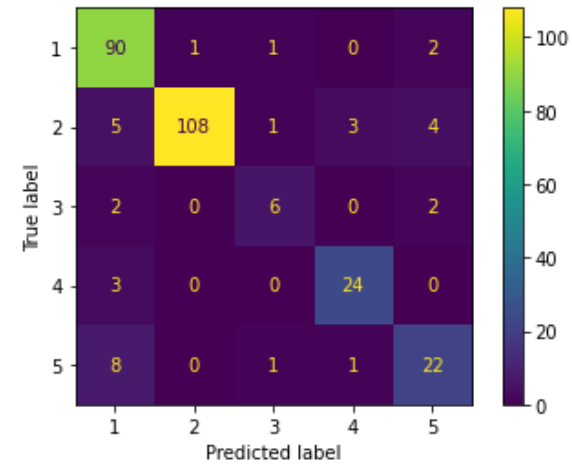
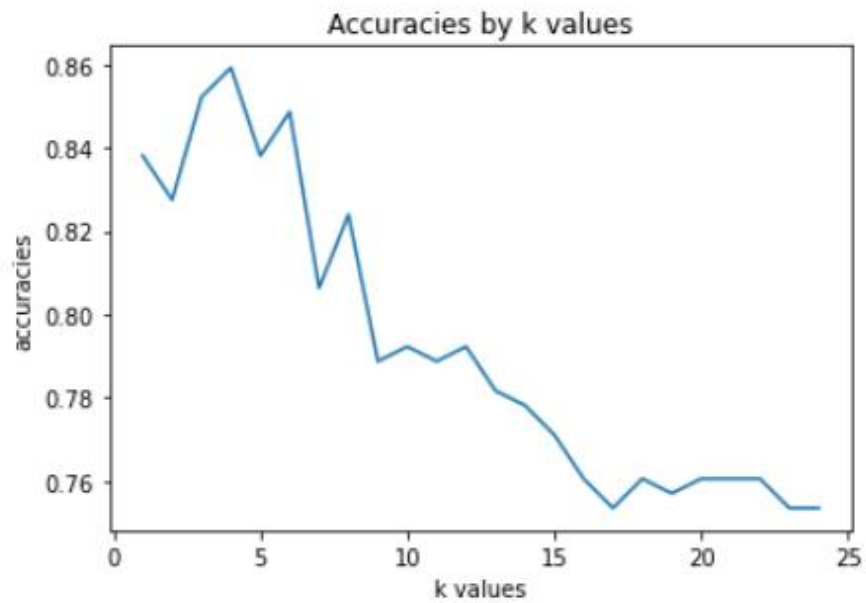


Corrélation à class:

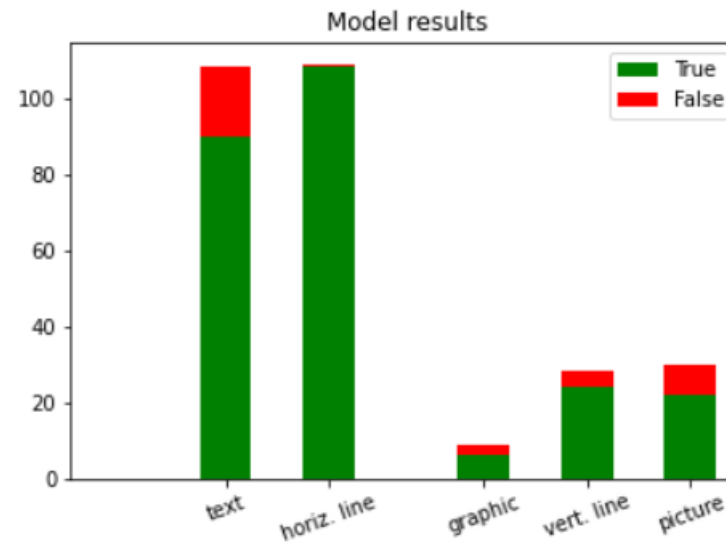
- Height
- Area
- P_and

KNN

- Recherche du k optimal pour notre model
- Meilleur k= 4



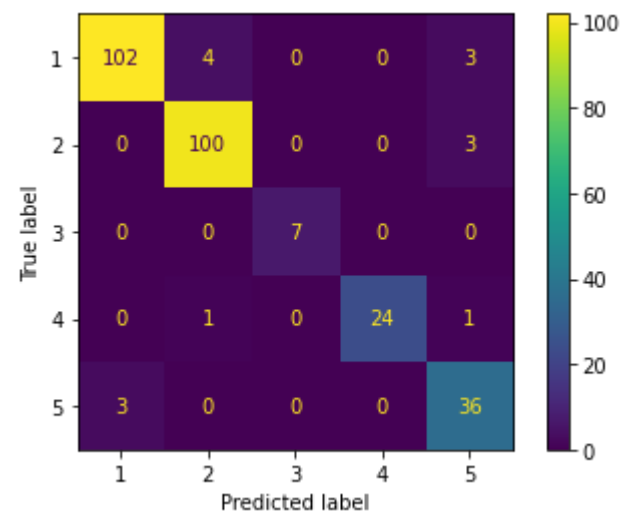
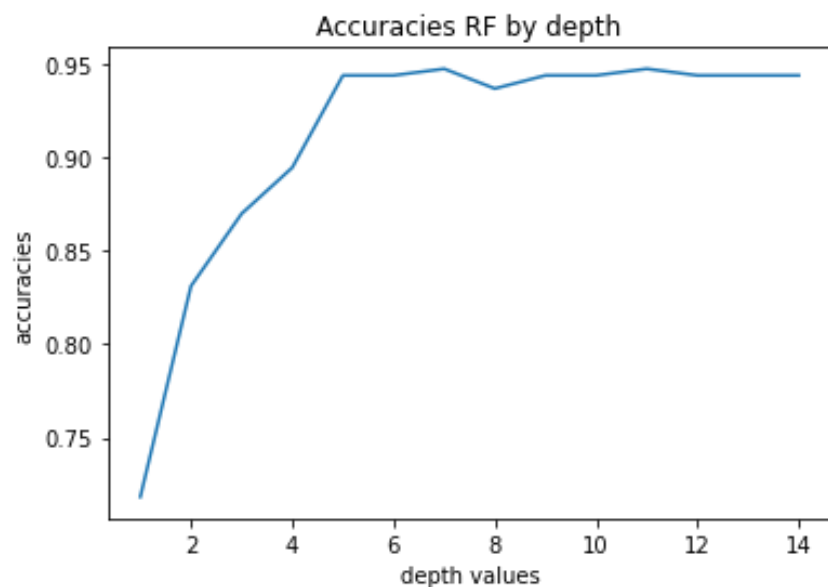
text : 83.3 %
horiz. line : 99.1 %
graphic : 66.7 %
vert. line : 85.7 %
picture : 73.3 %



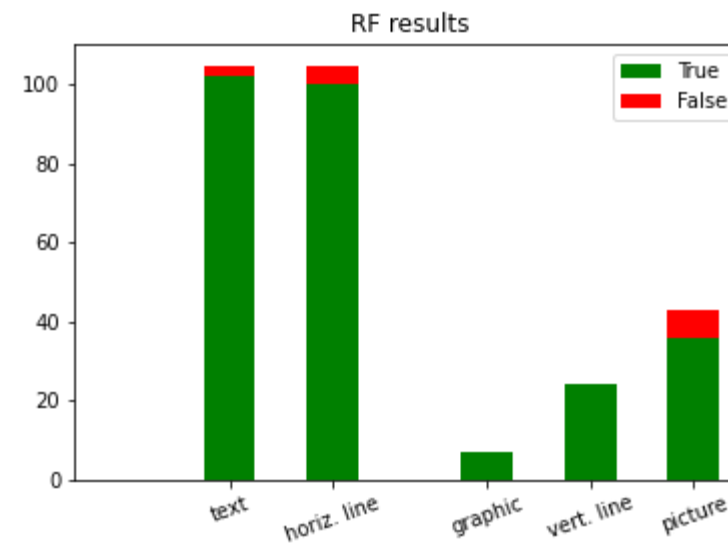
88.03 %

Random Forest

- Recherche de paramètre pour trouver le model optimal
- Meilleur depth= 7 ou 11



text : 97.1%
horiz. line : 95.1 %
graphic : 100 %
vert. line : 100 %
picture : 83.7 %

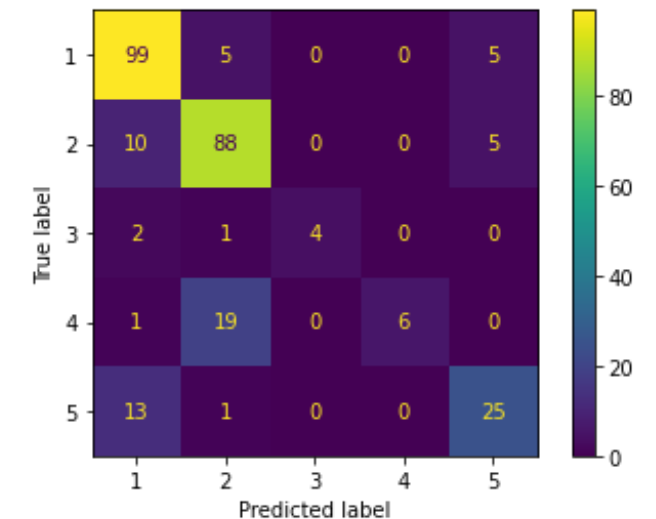


94,72 %

Logistic Regression



text : 72,2%
horiz. line : 77,2 %
graphic : 100 %
vert. line : 100 %
picture : 71,39 %



Accuracy: 78,17 %

API REST

- Input : fichier .csv ou .data
 - Output : Bloc classifiés + % de texte dans le document
 - *aller plus loin pour un potentiel client : une analyse de document "end to end"*
- Appliquer un algorithme de segmentation à partir d'un fichier pdf, jpg...
- Analyse automatique des 9 paramètres qui alimente le modèle
- Extraire l'information textuelle