

Rapport Projet Science des données

1. Présentation de la problématique du projet et des données

On dispose d'un jeu de données correspondant à des propriétés chimiques de vins, provenant d'un jeu de données disponible librement sur le dépôt de données de l'*University of California Irvine* 1. L'objectif de ce projet est de pouvoir prédire la qualité d'un vin à partir des données mises à disposition.

Les données sont définies dans un fichier ".csv", possédant les valeurs de 12 propriétés chimiques différentes caractérisant chaque vin, et séparées par le caractère ','.

Ces propriétés chimiques sont respectivement :

- "fixed acidity"
- "volatile acidity"
- "citric acid"
- "residual sugar"
- "chlorides"
- "free sulfur dioxide"
- "total sulfur dioxide"
- "density"
- "pH"
- "sulphates"
- "alcohol"
- "quality"

2. Analyse des données du problème

Tout d'abord, étant donné que nous cherchons à prédire la qualité d'un vin à partir des données mises à disposition, la classe à prédire est donc la "quality" du vin. Elle ne peut prendre que deux valeurs possibles (-1 ou 1), ce qui nous permet donc d'affirmer qu'il s'agit d'un problème de classification binaire.

Ensuite, il est important d'analyser les données pour identifier et corriger les données manquantes ou aberrantes. Une donnée manquante est représentée dans le fichier CSV par deux guillemets consécutifs, et par la valeur NaN dans un langage informatique comme python dans notre cas, comme le montre l'exemple ci-dessous :

```
54 | 6.6,"0.5","","2.1","0.068","6.0","14.0","0.9955","3.39","0.64","9.4","1"
```

fixed acidity	6.6000
volatile acidity	0.5000
citric acid	NaN
residual sugar	2.1000
chlorides	0.0680
free sulfur dioxide	6.0000
total sulfur dioxide	14.0000
density	0.9955
pH	3.3900
sulphates	0.6400
alcohol	9.4000
quality	1.0000

En ce qui concerne les données aberrantes, nous n'avons pas les connaissances nécessaires en oenologie pour affirmer qu'une valeur, supérieure à 3 écarts-types de la moyenne, soit réellement une valeur aberrante ou une valeur réelle mais hors-normes.

Dans le cadre de ce projet, nous avons pris le choix de conserver ces valeurs. En revanche, nous savons que le pH est une propriété chimique prenant des valeurs de 0 à 14, c'est pourquoi nous pouvons éliminer toutes données dont le pH n'est pas contenu dans cet interval.

Les histogrammes par attributs, disponibles en Annexe 1, permettent de visualiser les distributions des données restantes et ainsi de repérer quelles valeurs sont les plus exprimées par attribut.

On sait que deux attributs sont significativement corrélés si la valeur absolue du coefficient de corrélation > 0.5 . Or cela peut poser problème chez certains algorithmes de classification qui peuvent être sensible à la corrélation, et donc produire un modèle moins performant. On choisit d'afficher les coefficients de corrélation significatifs entre deux attributs pour avoir une idée de quels attributs pourraient être supprimés pour possiblement obtenir de meilleurs résultats :

Corrélations absolues significatives entre les différents attributs

	Fixed acidity	Volatile acidity	Citric acidity	Free sulfur dioxide	Total sulfur dioxide	Density	pH
Fixed acidity			0.671739			0.667610	0.684950
Volatile acidity			0.554347				
Citric acidity	0.671739	0.554347					0.540597
Free sulfur dioxide					0.668295		
Total sulfur dioxide				0.668295			
Density	0.667610						
pH	0.684950		0.540597				

Pour avoir un meilleur aperçu du type de corrélation qu'il existe entre ces attributs, vous pouvez trouver en Annexe 2 les graphiques de chacune de ces corrélations testées. On remarque alors que ces corrélations sont de types linéaires.

Avec ces résultats, on peut se demander si supprimer l'attribut "*fixed acidity*" pourrait potentiellement améliorer le modèle produit pour certains classifieurs parce que cet attribut est fortement corrélé à certains autres attributs. On pourrait également essayer de supprimer l'attribut "*free sulfur dioxide*" ou "*total sulfur dioxide*" car un seul suffirait à être pertinent. Nous verrons si cela est vrai par la suite lors de l'analyse des résultats et de l'influence des différents facteurs.

Certains classifieurs sont sensibles à la proportion d'individus de chaque classe et peuvent produire un modèle de mauvaise qualité si celle-ci est déséquilibrée. Après préparation des données, c'est-à-dire suppression des données manquantes et aberrantes, on obtient les proportions suivantes : proportion de 1 = 54% et proportion de -1 = 46%. On peut donc en conclure que les classes sont assez bien équilibrées.

Enfin, la normalisation des données peuvent être nécessaire chez certains classifieurs pour réaliser des modèles performant à partir de leurs modèles mathématiques sous-jacent. Les données restantes ont donc subi une normalisation.

Finalement, on se retrouve avec 1566 individus différents.

3. Description du protocole expérimental mis en place

On dispose, après préparation des données, de données normalisées. On souhaite réaliser une cross-validation afin de permettre à nos classifieurs de pouvoir s'entraîner et produire un modèle à partir de données d'entraînements, puis de les évaluer à l'aide de métriques grâce aux données de tests restantes. Il est donc primordiale pour une bonne comparaison d'utiliser les mêmes données d'entraînements et de tests entre classifieurs.

Pour effectuer la cross-validation, nous avons tout d'abord divisé nos données en différents jeux (ou plis) de tests et d'entraînements en utilisant la méthode du Stratified KFold à 10 itérations. Cela permet d'éviter un sur-apprentissage des classifieurs et de respecter la stratification, c'est à dire d'assurer une distribution des classes dans les plis fidèle au jeu de données complet. On se retrouve alors avec les données suivantes :

- X_train : la liste de chaque plis d'entraînements, utilisé par le classifieur pour créer son modèle.
- y_train : la liste des valeurs de qualité réelle correspondant respectivement à chacun des plis d'entraînements, utilisé par le classifieur pour créer son modèle.
- X_test : la liste de chaque plis de tests, utilisée par le classifieur pour prédire une qualité à partir du modèle créé.
- y_test : la liste des valeurs de qualité réelle correspondant aux plis de tests, utilisée pour tester et évaluer le classifieur par comparaison avec les qualités prédites.

Dans le cadre de ce projet, nous allons tester les classifieurs suivants :

- Régression logistique (Logistic regression)
- Machines à vecteurs supports (SVC avec `max_iter = 10000`)
- Analyse discriminante linéaire (Linear discriminant analysis)
- Analyse discriminante quadratique (Quadratic discriminant analysis)
- K-plus proches voisins (K-nearest neighbors classification)
- Arbres de décision (Decision trees)
- Perceptron (avec `epochs=500`, `learning_rate=0.01`)

Comme expliqué précédemment, on commence par entraîner chaque classifieur avec chacun des plis de `X_train` et `y_train` créés pour produire un modèle par cross-validation. Ce modèle est ensuite testé à l'aide des plis de `X_test` afin d'obtenir des qualités prédites. Ces qualités sont comparées aux vraies qualités correspondantes (`y_test`) en calculant les métriques suivantes :

- Précision
- Recall
- F1 score
- Accuracy

Puisqu'on calcule ces métriques pour chaque jeu de tests, on réalise la moyenne des résultats obtenus pour avoir des résultats de métriques plus précis.

4. Résultats expérimentaux obtenus

Les résultats obtenus seront présentés sous forme de tableaux regroupant chaque classifieur et chaque métrique correspondante. Le plus grand résultat pour chaque métrique sera mis en valeur par la coloration de la cellule correspondante en jaune. Chaque tableau représente une situation différente. On retrouvera donc, dans l'ordre :

- Résultats obtenus avec les données aberrantes retirées et sans normalisation
- Résultats obtenus avec les données aberrantes retirées et normalisation
- Résultats obtenus avec les données aberrantes et les valeurs supérieures à 3 écarts-types de la moyenne retirées et normalisation
- Résultats obtenus avec les données aberrantes, "fixed acidity" retirés et normalisation
- Résultats obtenus avec les données aberrantes, "fixed acidity" et "free sulfur dioxide" retirés et normalisation

1/Résultats obtenus avec les données aberrantes retirées et sans normalisation

	Précision	Recall	F1 score	Accuracy
Régression logistique	0.76460768193 351	0.73623249299 71988	0.74352164399 62973	0.73438673852 68659
SVC	Ne converge pas			
Analyse discriminante linéaire	0.76673903410 90918	0.73621848739 4958	0.74348279340 40067	0.73566062387 71844
Analyse discriminante quadratique	0.72667250790 39789	0.77187675070 02802	0.74198753456 15331	0.71453944145 02695
K-plus proches voisins	0.62265375287 56054	0.62340336134 45378	0.61587844369 56447	0.59185448309 65212
Arbres de décision	0.66661124735 90428	0.66025210084 03361	0.65939608981 45832	0.64114813000 16331
Perceptron	0.72911990663 92877	0.70422969187 67508	0.68781335439 3197	0.68329658664 05357

2/Résultats obtenus avec les données aberrantes retirées et
normalisation

	Précision	Recall	F1 score	Accuracy
Régression logistique	0.76288334900 83109	0.74455182072 82913	0.74670589503 18171	0.73630164951 821
SVC	0.76509800654 91467	0.73978991596 63865	0.74514176404 1775	0.73630164951 82101
Analyse discriminante linéaire	0.76673903410 90918	0.73621848739 4958	0.74348279340 40067	0.73566062387 71844
Analyse discriminante quadratique	0.72667250790 39789	0.77187675070 02802	0.74198753456 15331	0.71453944145 02695
K-plus proches voisins	0.69133920659 64394	0.73754901960 78431	0.70576079528 12984	0.67691082802 54778
Arbres de décision	0.66263576750 71856	0.65787114845 93838	0.65629968550 41525	0.63731830801 89448
Perceptron	0.77457142370 49759	0.58780112044 81793	0.65853838958 81509	0.68644863628 94007

3/Résultats obtenus avec les données aberrantes et les valeurs supérieures à 3 écarts-types de la moyenne retirées et normalisation

	Précision	Recall	F1 score	Accuracy
Régression logistique	0.75801275748 40508	0.75269730269 73028	0.74488929162 60922	0.73113857973 01291
SVC	0.76748858101 17255	0.74750249750 24975	0.74686035219 9227	0.73746675859 35192
Analyse discriminante linéaire	0.77164419732 52365	0.74620379620 37963	0.74682197495 37753	0.73887028464 49325
Analyse discriminante quadratique	0.73991469008 46458	0.72928737928 73793	0.72082343627 53386	0.70726386289 76657
K-plus proches voisins	0.70043513398 55133	0.72682317682 31769	0.70287272527 71974	0.67846449325 32257
Arbres de décision	0.66276167489 21992	0.64783549783 54978	0.64841540325 23388	0.63208903772 28406
Perceptron	0.70561971286 33885	0.60760905760 90575	0.64440657291 52629	0.65102432778 4891

4/Résultats obtenus avec les données aberrantes, "fixed acidity" retirées et normalisation

	Précision	Recall	F1 score	Accuracy
Régression logistique	0.76665947167 59618	0.74812324929 97199	0.75053608404 01764	0.74078066307 36566
SVC	0.77218199355 88243	0.74218487394 95798	0.74992311271 84574	0.74204638249 22423
Analyse discriminante linéaire	0.77729995796 73613	0.74218487394 95798	0.75124526069 45953	0.74459415319 28792
Analyse discriminante quadratique	0.72254524229 02949	0.78138655462 18486	0.74168664256 66147	0.71264494528 82574
K-plus proches voisins	0.70203068524 50935	0.74822128851 54061	0.71690714865 12034	0.68774293646 90512
Arbres de décision	0.67126613556 13829	0.67568627450 98039	0.66955809792 25966	0.64949371223 25657
Perceptron	0.76332862886 31962	0.61037815126 05042	0.67059201865 3906	0.69092764984 48473

5/Résultats obtenus avec les données aberrantes, “fixed acidity” et “free sulfur dioxide” retirées et normalisation

	Précision	Recall	F1 score	Accuracy
Régression logistique	0.75680723483 58513	0.74808123249 29971	0.74705474117 84442	0.73565654091 1318
SVC	0.76161615033 61452	0.74333333333 33333	0.74673423200 51349	0.73693042626 16365
Analyse discriminante linéaire	0.76577689301 61184	0.73976190476 19047	0.74691891541 72445	0.73820839457 78213
Analyse discriminante quadratique	0.72856926783 34504	0.78260504201 68067	0.74892577924 40236	0.72221541727 91115
K-plus proches voisins	0.68855060972 31298	0.73626050420 16807	0.70392908200 97283	0.67495917034 1336
Arbres de décision	0.67086486044 06326	0.66854341736 69468	0.66599247540 84171	0.64821166095 05144
Perceptron	0.79408083792 54957	0.60672268907 56302	0.67940627678 17362	0.70751674016 00522

5. Conclusion et analyse critique des résultats obtenus

Après analyse des résultats obtenus, on remarque tout d'abord que la normalisation des données est importante pour un meilleur fonctionnement des classifieurs comme le SVC notamment.

En revanche, la suppression des valeurs à plus de 3 écarts-types apportent des légères améliorations de métriques pour certains classifieurs, mais ce n'est pas le cas pour tous.

Selon la préparation des données et la métrique privilégiée, le choix du classifieur ne sera pas forcément le même, comme présenté dans le tableau ci-dessous.

Meilleurs classifieurs selon les différentes situations et métriques

	Précision	Recall	F1 score	Accuracy
Situation 1	Analyse discriminante linéaire	Régression logistique	Régression logistique	Analyse discriminante linéaire
Situation 2	Perceptron	Analyse discriminante quadratique	Régression logistique	SVC
Situation 3	Analyse discriminante linéaire	Régression logistique	SVC	Analyse discriminante linéaire
Situation 4	Analyse discriminante linéaire	Analyse discriminante quadratique	Analyse discriminante linéaire	Analyse discriminante linéaire
Situation 5	Perceptron	Analyse discriminante quadratique	Analyse discriminante quadratique	Analyse discriminante linéaire

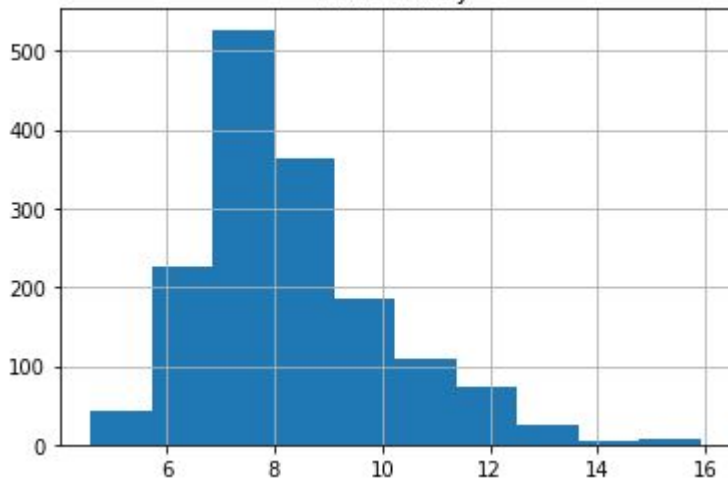
En général, on remarque que le classifieur à “*Analyse discriminante linéaire*” est le plus performant dans la majorité des situations, mais la situation 4 semble être préférable aux autres car ce classifieur atteint le maximum de F1 score et accuracy, et la deuxième meilleure précision toutes situations confondues.

L'accuracy obtenue à la situation 4 est 0.7445941531928792. Si nous voulons avoir des classifieurs avec une accuracy accrue, il nous faudrait donc un échantillon plus large, ainsi les classifieurs pourront avoir plus de données d'entraînement et de test pour produire des modèles plus fiables.

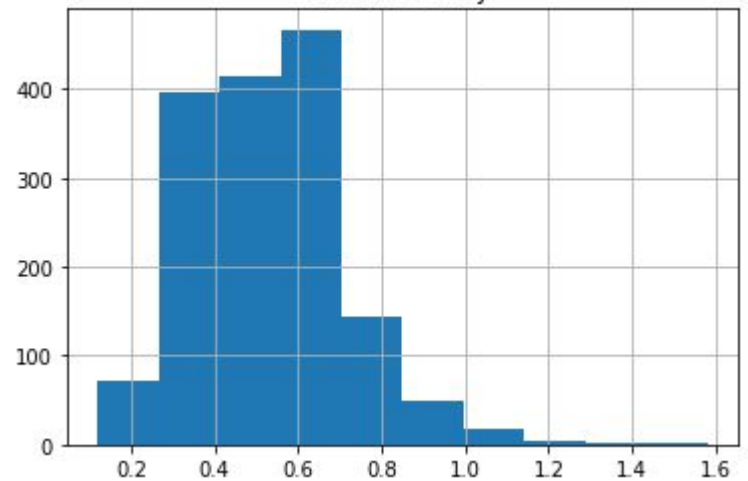
ANNEXES

1/Histogrammes représentant la distribution des attributs

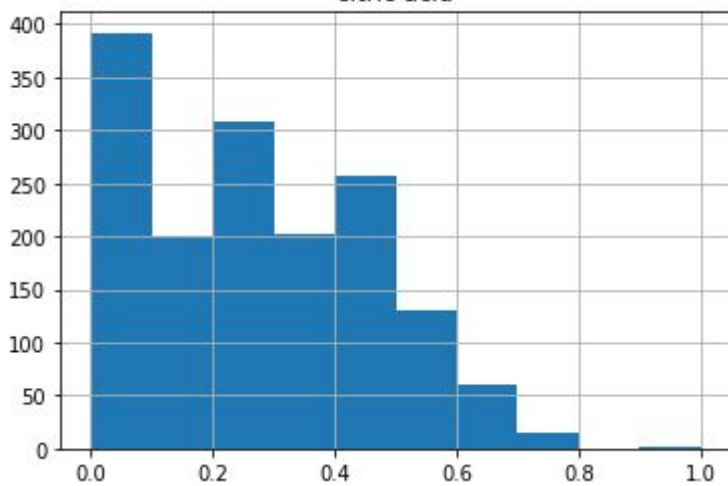
fixed acidity



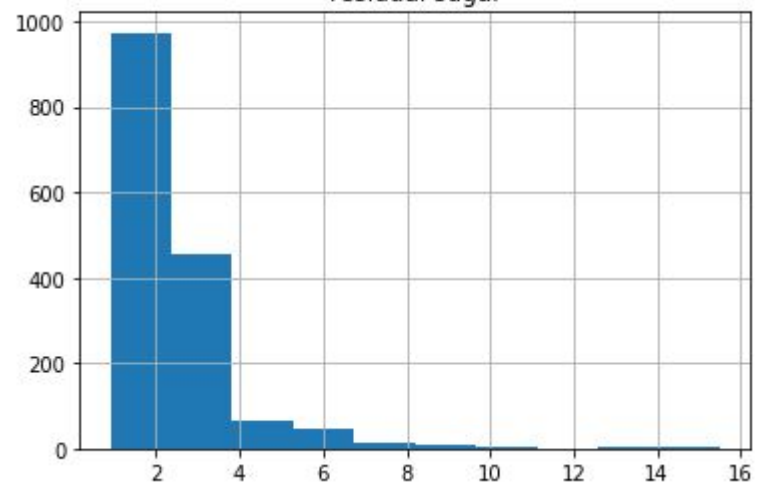
volatile acidity



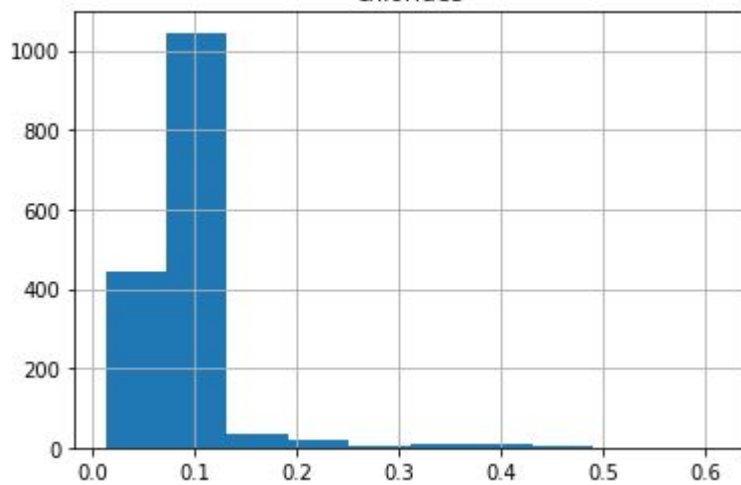
citric acid



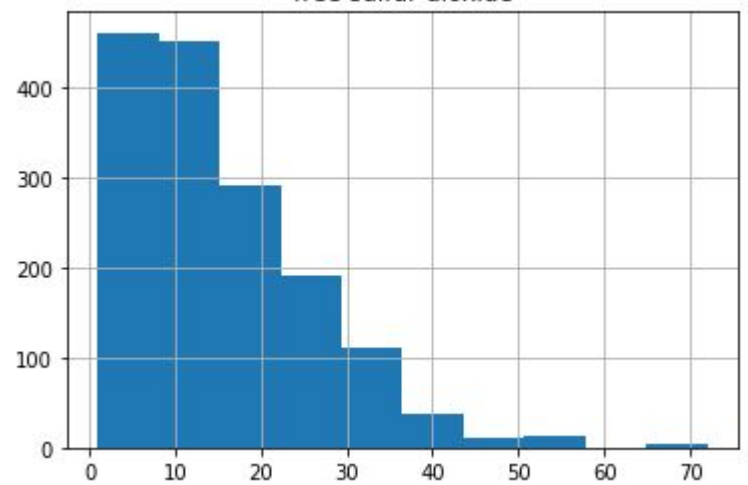
residual sugar



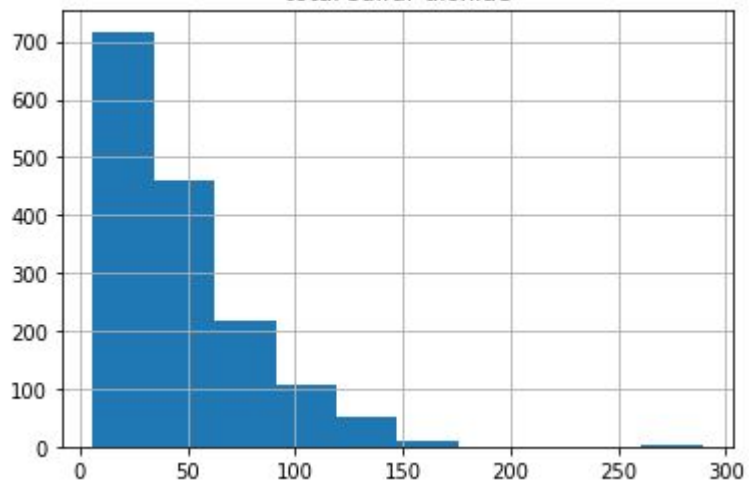
chlorides



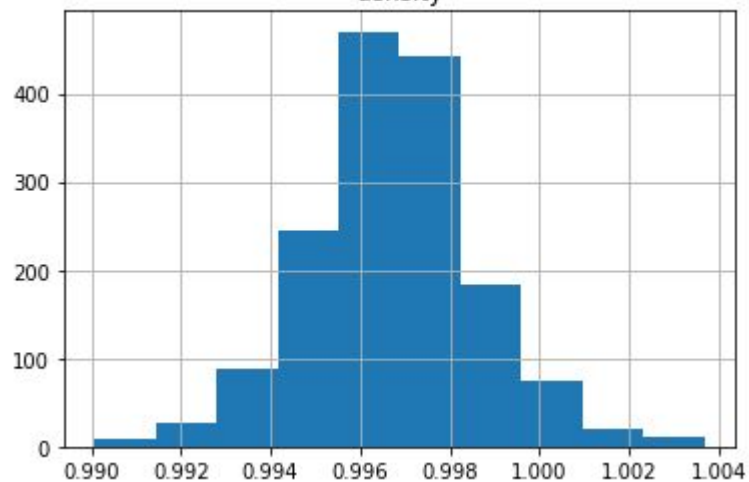
free sulfur dioxide



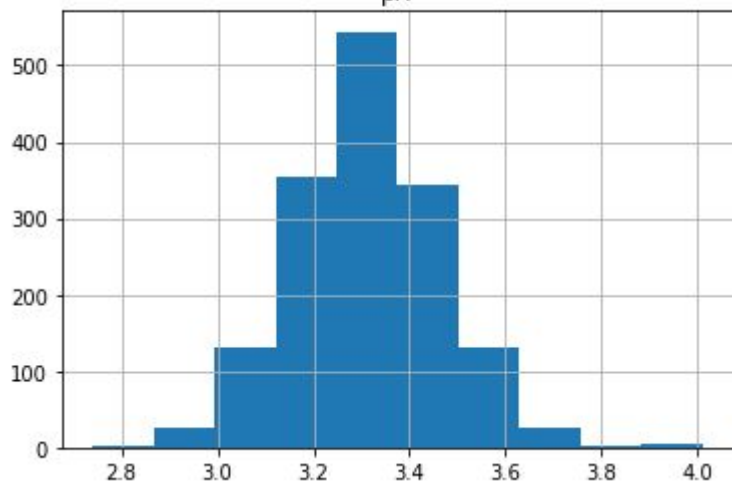
total sulfur dioxide



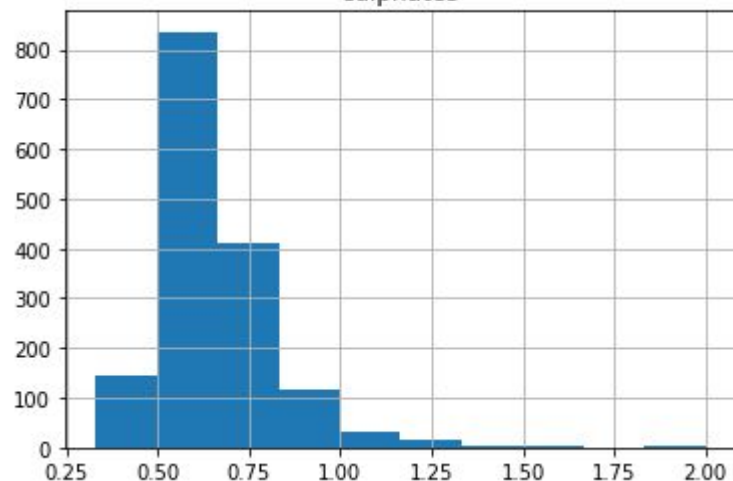
density



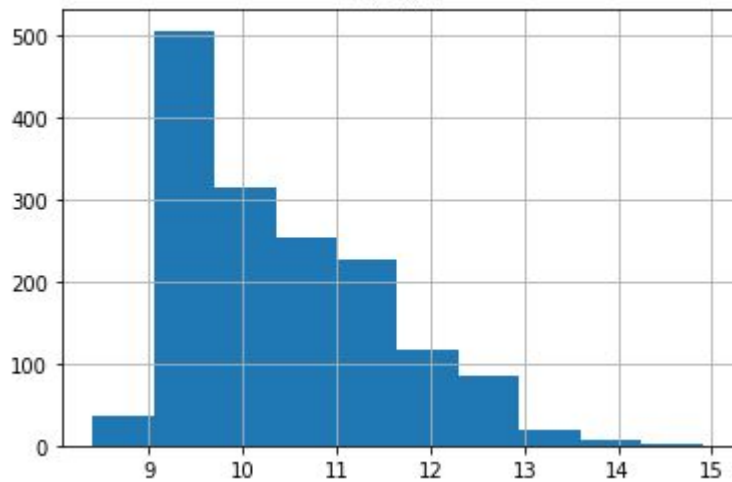
pH



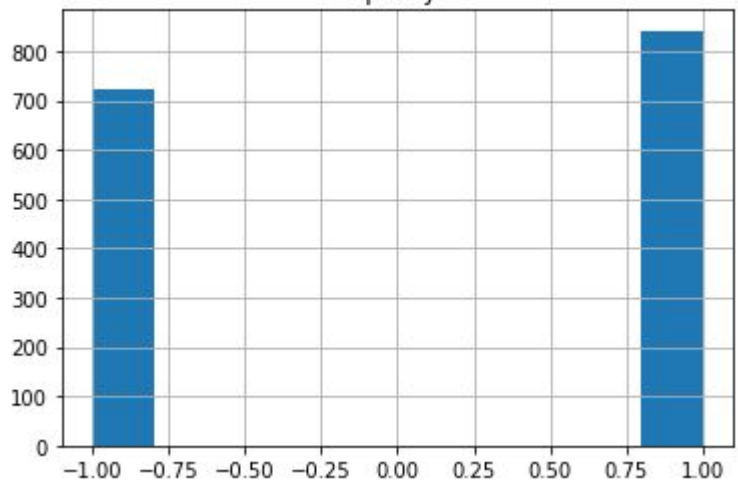
sulphates



alcohol



quality



2/Graphiques des corrélations entre les différents attributs

