

projet traitement numérique

Rédacteur :

Yanis HAMMACI
Réda ARBANE
Akram CHAABNIA

14/04/2023

Contents

1	Aperçu des données	2
1.1	Description du jeu de données	2
1.2	Description des variables	2
1.3	But de l'analyse	2
2	Prétraitement [Annexe : statistique figures.pdf]	3
3	Analyse univariée [Annexe : statistique figures.pdf]	3
3.1	Variables quantitatives	3
3.2	Variables qualitatives	3
4	Analyse Bivariée	3
4.1	Corrélations [Annexe : COR_ACP_figures.pdf]	3
4.1.1	Analyse Globale des Corrélations	4
4.1.2	Détails du Heatmap de Corrélation	4
4.1.3	Corrélations Triées par Valeurs Absolues	4
4.1.4	Conclusion	4
5	Classification	4
5.1	K-means [Annexe : kmeans_figures.pdf]	4
5.1.1	Évaluation	4
5.1.2	Interprétation des résultats	5
5.1.3	Conclusion	6
5.2	Classification Ascendante Hiérarchique (CAH) [Annexe : CAH_figures.pdf]	6
5.2.1	Préparation et Analyse	6
5.2.2	Dendrogramme et Détermination du Nombre Optimal de Clusters	6
5.2.3	Visualisation des Partitions	6
5.2.4	Analyse des Clusters	7
5.2.5	Conclusion	8
6	ACP [Annexe : COR_ACP_figures.pdf]	8
6.1	Réduction de Dimension et Interprétation	8
6.2	Importance des Axes Principaux	8
6.3	Contribution des Variables aux Axes	8
6.4	Qualité de Représentation des Variables sur les deux Axes Factoriel	9
6.5	Visualisation des Clusters de Marques	9
6.6	Conclusion sur l'ACP	9
7	Conclusion	9

1 Aperçu des données

1.1 Description du jeu de données

L'ensemble de données contient des informations sur 261 voitures couvrant 8 variables qui proviennent de trois régions (les USA, l'Europe, le Japon) entre l'année 1971 et 1983.

A data.frame: 6 × 8

	mpg	cylinders	cubicinches	hp	weightlbs	time.to.60	year	brand
	<chr>	<int>	<int>	<int>	<int>	<int>	<int>	<chr>
1	14	8	350	165	4209	12	1972	US.
2	31.9	4	89	71	1925	14	1980	Europe.
3	17	8	302	140	3449	11	1971	US.
4	15	8	400	150	3761	10	1971	US.
5	30.5	4	98	63	2051	17	1978	US.
6	23	8	350	125	3900	17	1980	US.

1.2 Description des variables

1. **MPG (Miles Per Gallon)** : Distance (en miles) pouvant être parcourue avec un gallon d'essence.
2. **Cylinders** : Nombre de cylindres dans le moteur d'une voiture.
3. **Cubic Inches** : Mesure du volume d'un cylindre.
4. **HP (Horsepower)**: La puissance réelle du moteur en chevaux.
5. **Weight (lbs)**: Le poids de la voiture en livres.
6. **Time to 60** : Temps nécessaire pour aller de 0 à 60 mph (miles per hour).
7. **Year** : Année de fabrication de la voiture.
8. **Brand** : Lieu de fabrication de la voiture.

1.3 But de l'analyse

1. **Analyse de l'efficacité énergétique** : explorez les facteurs affectant l'efficacité énergétique (MPG) des voitures.
2. **Analyse des performances du moteur** : étudiez la relation entre les mesures de performances du moteur (HP) et les autres attributs.
3. **Analyse du temps nécessaire pour aller de 0 à 60 mph** : étudier le temps nécessaire pour aller de 0 à 60 mph en fonction des autres attributs.
4. **Analyse temporelle** : examinez l'évolution des attributs de la voiture au fil du temps.
5. **Comparaison de marques** : comparez les voitures de différentes marques en termes de caractéristiques telles que le MPG, la puissance, le poids, etc.

2 Prétraitement [Annexe : statistique figures.pdf]

le jeu de données ne contient aucune valeur manquante, il peut donc être utilisé tel quel pour la suite de l'analyse (figure 7)

La variable qualitative **brand** est de type char on peut la transformer en type factor car elle contient un nombre fini de valeurs (US, Europe, Japan).

3 Analyse univariée [Annexe : statistique figures.pdf]

3.1 Variables quantitatives

Pour évaluer la dispersion des variables quantitatives, nous avons adopté une approche exhaustive en calculant non seulement la médiane, le minimum, le maximum et les quartiles (Q1 et Q3) à l'aide de la fonction `summary`, mais aussi en déterminant les quartiles à des intervalles de 10 % grâce à la fonction `quantile` (figure 1).

A partir du diagramme en boîte qui nous permet de visualiser ces résultats dans la (figure 2), on peut déduire le suivant :

- **mpg (Miles per gallon)** Les valeurs sont plus dispersées au-dessus de la médiane
- **cylinders** Les valeurs sont majoritairement centrées autour de la médiane, avec une dispersion similaire des valeurs en dessous et au-dessus de la médiane. En examinant les valeurs existantes (3, 4, 5, 6, 8), on constate que 98 % des voitures (256 voitures) ont des valeurs parmi (4, 6, 8), tandis que seulement 2 % des voitures (8 voitures) ont des valeurs parmi (3, 5) (figure 5)
- **cubicinches** Les valeurs sont plus dispersées au-dessus de la médiane, mais il n'y a pas de valeurs aberrantes évidentes.
- **hp(Horsepower)** les valeurs sont plus dispersées au-dessus de la médiane, avec un coefficient d'asymétrie le plus élevé parmi les variables du jeu de données 0.96.
- **weightlbs (Weight in lbs)** La distribution des valeurs est légèrement au-dessus de la médiane.
- **time.to.60** On remarque une distribution plutôt uniforme des valeurs. Quelques valeurs aberrantes sont détectées dans les deux côtés. On identifie 11 valeurs aberrantes (figure 4). Ces valeurs correspondent à des voitures avec des accélérations anormalement rapides (8 secondes pour atteindre 60 mph pour la plus rapide) ou faibles (25 secondes pour la plus lente). Ces observations ne semblent pas résulter d'erreurs de mesure, mais plutôt de voitures qui se situent en dehors de la norme. Donc on peut garder ces valeurs tel qu'elles sont.
- **year** Le coefficient d'asymétrie(skewness) est proche de 0 ce qui indique une distribution symétrique (figure 3).

3.2 Variables qualitatives

La variable **brand** est la seule variable qualitative du jeu de données. On remarque que 62 % des voitures sont américaines, 20 % européenne et 18 % japonaise (figure 6-7).

La dominance des voitures américaines dans notre jeu de données est à prendre en compte durant l'analyse.

4 Analyse Bivariée

4.1 Corrélations [Annexe : COR_ACP figures.pdf]

L'analyse de corrélation est un outil statistique fondamental qui mesure l'intensité et la direction de la relation entre deux variables quantitatives. Elle est essentielle pour comprendre comment les caractéristiques des véhicules influencent les uns les autres, permettant ainsi de déduire des causalités potentielles et des dépendances. Par exemple, identifier comment la puissance du moteur affecte la consommation de carburant ou l'impact du poids sur l'accélération. Ce type d'analyse aide à optimiser les conceptions de véhicules en fonction des attentes en termes de performance et d'efficacité énergétique.

4.1.1 Analyse Globale des Corrélations

Le tableau de corrélation (figure 1) illustre les relations entre les différentes caractéristiques techniques des véhicules. Les coefficients de corrélation varient entre -1 et 1, où -1 indique une corrélation négative parfaite, 0 aucune corrélation, et 1 une corrélation positive parfaite. Une observation importante est la forte corrélation négative entre **mpg** et les variables **cylinders**, **cubicinches**, **hp**, et **weightlbs**. Ces coefficients, variant de -0.7767099 à -0.8246487, suggèrent que les voitures plus puissantes et plus lourdes ont tendance à consommer plus de carburant, ce qui est intuitivement logique étant donné la demande énergétique des moteurs plus grands et des véhicules plus massifs.

4.1.2 Détails du Heatmap de Corrélation

(figure 2) utilise un heatmap pour visualiser ces corrélations, où les couleurs passent du bleu (corrélation positive) au rouge (corrélation négative). Cette visualisation met en évidence des blocs de forte corrélation, en particulier parmi les variables liées aux spécifications du moteur (**cylinders**, **cubicinches**, **hp**) et le poids **weightlbs**, tous montrant des corrélations positives supérieures à 0.84. Ces fortes corrélations positives indiquent que les caractéristiques du moteur sont liées entre elles : des moteurs avec plus de cylindres tendent à avoir plus de pouces cubiques et donc plus de puissance.

4.1.3 Corrélations Triées par Valeurs Absolues

(figure 3) présente clairement les corrélations, organisées du plus au moins important, ce qui facilite l'identification rapide des relations les plus fortes et les plus faibles entre les variables. Les données montrent que les dimensions du moteur **cubicinches** et le nombre de cylindres **cylinders** sont presque parfaitement corrélés (0.9512776), ce qui reflète une consistance dans la conception des moteurs où un plus grand nombre de cylindres implique généralement un moteur plus grand. Cette figure permet également de noter des corrélations modérées comme celle entre **time.to.60** (temps pour atteindre 60 mph) et **hp** (-0.7448731), suggérant que des voitures plus puissantes atteignent plus rapidement 60 mph.

4.1.4 Conclusion

Cette analyse des corrélations montre comment différentes caractéristiques des voitures dépendent les unes des autres, nous donnant des informations utiles sur ce qui influence la performance et l'efficacité énergétique des voitures. En plus des relations déjà évoquées, nous voyons d'autres corrélations intéressantes. Par exemple, la corrélation entre les **cubicinches** et le **time.to.60** avec un coefficient de -0.6084126 montre que les moteurs plus grands peuvent rendre les voitures plus rapide. De plus, la relation entre **cylinders** et **year** (coefficient de -0.3222394) pourrait indiquer une tendance à utiliser moins de cylindres dans les moteurs plus récents.

5 Classification

5.1 K-means [Annexe : kmeans_figures.pdf]

Le K-means est un algorithme de clustering non supervisé largement utilisé. Il partitionne un ensemble de données en k groupes en minimisant la variance intra-cluster (entre les individus du même cluster) et maximisant la variance inter-cluster (entre les individus des clusters différents).

pour appliquer cette algorithme à notre jeu de données on doit préciser le nombre de cluster et pour cela on peut utiliser certaines méthodes pour déterminer le nombre optimale de cluster

5.1.1 Évaluation

1. **Critère de silhouette** : Il mesure à quel point les objets d'un même cluster sont similaires entre eux par rapport à ceux des autres clusters, sur une échelle de -1 à 1. Un score élevé (proche de 1) indique une bonne séparation entre les clusters, tandis qu'un score proche de 0 signifie que les clusters se chevauchent ou qu'un objet pourrait appartenir à un autre cluster. Un score négatif suggère qu'un objet a été mal classé.

On calcule la silhouette $s(i)$ pour chaque individu du jeu de données :

- La cohésion $a(i)$ est la distance moyenne entre le point i et tous les autres points du même cluster.
- La séparation $b(i)$ est la distance moyenne entre le point i et tous les points du cluster voisin le plus proche.
- La silhouette $s(i)$ est définie comme $\frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$.

La silhouette moyenne de tous les points est ensuite calculée pour évaluer la qualité globale du clustering

On remarque sur notre jeu de données, on obtient un score de silhouette maximum de 0.62 avec 2 clusters (figure 1), un score qui est proche de 0.6 avec 3 ou 4 clusters.

2. **méthode du coude** : Elle consiste à tracer la somme des carrés intra-cluster en fonction du nombre de clusters et à identifier le point où l'ajout d'un cluster supplémentaire ne réduit plus significativement cette somme, formant ainsi une "coude" dans le graphique. Ce point est considéré comme le nombre optimal de clusters.

On remarque sur notre jeu de données, le nombre optimal de cluster est 4 car à partir de ce point la différence des WCSS semble se stabiliser (figure 2)

5.1.2 Interprétation des résultats

Les résultats du Kmeans sont disponibles dans (figure 3-4-5)

On remarque que les classes se séparent clairement dans les 3 figures.

On remarque que les variables corrélées ont une distribution similaire de leurs classes, on remarque aussi que la distribution des classes sont en miroir lorsque qu'on prend deux variables anti-corrélées par exemple **mpg** et **cylinders**, **cubicinches**, **hp**, **weightlbs**.

On remarque aussi que la séparation des classes est plus prononcée pour certaines variables comme **weightlbs** comparé à **mpg**.

pour expliquer les classes trouvées par le kmeans on peut utiliser la variable **brand**, on suppose que l'origine de la voiture détermine dans quelle classe elle appartient.

Kmeans avec 2 classes

En examinant les résultats du kmeans (figure 6), on remarque que 92% des voitures européennes (44 voitures) et 100 % des voitures japonaises (51 voitures) sont dans la première classe tandis qu'on trouve que 38 % des voitures américaines (63 voitures). Dans la deuxième classe on trouve 62% des voitures américaines (99 voitures), seulement 8 % des voitures européennes (4 voitures) et 0 voitures japonaises.

On remarque que si une voiture est européenne ou japonaise, il y a de grande chance qu'elle soit dans la classe1, alors que pour les voitures américaines peuvent appartenir aux deux classes, avec plus de chance dans la classe 2.

Kmeans avec 3 classes

En examinant les résultats de la (figure 7), on remarque que 63% des voitures européennes (30 voitures) et 80 % des voitures japonaises (41 voitures) sont dans la première classe tandis qu'on trouve que 13 % des voitures américaines (21 voitures). Dans la deuxième classe on trouve 44% des voitures américaines (72 voitures), seulement 2 % des voitures européennes (1 voiture) et 0 voitures japonaises. Dans la troisième classe on trouve 43 % des voitures américaines (69 voitures), 35 % des voitures européennes (17 voitures) et 20% des voitures japonaises (10 voitures).

On remarque que la majorité des voitures européennes et japonaises, sont dans la classe 1 (un pourcentage supérieure des voitures japonaises). Les voitures américaines se trouvent majoritairement dans la classe 2 et 3. Une minorité de voitures européennes et japonaises se trouvent dans la classe 3 (un pourcentage supérieure des voitures européennes).

Kmeans avec 4 classes

On remarque une répartition similaire des voitures entre les différents groupes, que ce soit pour les k-means précédents avec $k=2$ et $k=3$ (figure 8). La majorité des voitures européenne et japonaise sont dans la classe 4 et 1 tandis que les voitures américaine se trouvent dans toutes les classes avec une présence plus significative dans la classe 3, 2 et 1.

5.1.3 Conclusion

Il est évident qu'il existe une distinction quantitative entre les diverses marques de voitures, avec les voitures américaines et non américaines constituant les deux principales catégories de notre ensemble de données. Les voitures européennes et japonaises présentent des similitudes plus marquées entre elles que celles observées entre les voitures américaines et européenne ou entre les voitures américaine et japonaise. Les voitures américaines sont présentes dans toutes les classes, avec une répartition plus équilibrée, se situant légèrement plus bas dans certaines classes et légèrement plus haut dans d'autres. La distinction entre les voitures européennes et japonaises établie à travers les différentes classes n'est pas clairement, ce qui suggère des similitudes dans leurs caractéristiques.

5.2 Classification Ascendante Hiérarchique (CAH) [Annexe : CAH_figures.pdf]

5.2.1 Préparation et Analyse

La Classification Ascendante Hiérarchique (CAH) a été appliquée pour regrouper les observations de notre jeu de données en fonction de leurs similarités. Cette méthode commence par considérer chaque observation comme un cluster distinct, puis les fusionne progressivement en se basant sur leur proximité. La distance euclidienne a été utilisée comme mesure de proximité et la méthode de Ward a été choisie pour minimiser la variance intra-cluster.

5.2.2 Dendrogramme et Détermination du Nombre Optimal de Clusters

Le dendrogramme de la CAH (Figure 1) visualise la hiérarchie des fusions des clusters et permet de déterminer le nombre optimal de clusters. L'axe horizontal représente les observations et l'axe vertical représente la distance entre les clusters. En observant le dendrogramme, on peut identifier des regroupements naturels de données à différentes hauteurs de fusion.

L'analyse de l'inertie (Figure 2) vient compléter l'analyse du dendrogramme pour affiner la détermination du nombre optimal de clusters. Le graphique de l'inertie montre l'évolution de l'inertie totale à mesure que le nombre de clusters augmente. L'inertie représente la variabilité intra-classe, c'est-à-dire la dispersion des points de données à l'intérieur de chaque cluster.

En observant le graphique de l'inertie, on remarque une diminution importante de l'inertie lorsque le nombre de clusters augmente de 1 à 4. Cela indique que ces regroupements successifs permettent de capturer une part significative de la variabilité des données. Au-delà de 4 clusters, la diminution de l'inertie devient beaucoup plus lente, suggérant que les regroupements supplémentaires n'apportent pas d'amélioration majeure à la structure des données.

Le point d'inflexion observé sur le graphique de l'inertie, souvent appelé "coude", correspond approximativement à 4 clusters. Ce coude indique que c'est à ce niveau que l'ajout de clusters supplémentaires devient moins bénéfique en termes de réduction de la variabilité intra-classe.

En combinant l'analyse du dendrogramme et du graphique de l'inertie, on peut conclure que le nombre optimal de clusters pour cette analyse est de 4.

5.2.3 Visualisation des Partitions

Les partitions des données en 2, 3 et 4 clusters sont illustrées dans le dendrogramme (Figure 3). Les différentes couleurs représentent les différents clusters. Cette visualisation permet d'observer comment les données sont regroupées à différents niveaux de granularité.

5.2.4 Analyse des Clusters

L'analyse des clusters révèle des regroupements distincts au sein des données, caractérisés par des caractéristiques communes des observations dans chaque groupe.

Nombre de clusters: L'analyse a été réalisée avec 2, 3 et 4 clusters. Le choix du nombre optimal dépend de l'objectif de l'analyse et de l'interprétation souhaitée des résultats.

Répartition des individus dans les classes:

2 classes:

- La classe 1 regroupe 100% des individus européens et japonais.
- La classe 2 regroupe 38% des individus américains, 8% des individus européens et 0% des individus japonais.

3 classes:

- La classe 1 regroupe 63% des individus européens et 80% des individus japonais.
- La classe 2 regroupe 44% des individus américains, 2% des individus européens et 0% des individus japonais.
- La classe 3 regroupe 43% des individus américains, 35% des individus européens et 20% des individus japonais.

4 classes:

- La classe 1 regroupe 0% des individus européens, 0% des individus japonais et 23% des individus américains.
- La classe 2 regroupe 23% des individus européens, 43% des individus japonais et 14% des individus américains.
- La classe 3 regroupe 71% des individus européens, 47% des individus japonais et 9% des individus américains.
- La classe 4 regroupe 6% des individus européens, 10% des individus japonais et 54% des individus américains.

Interprétation des résultats:

2 classes:

- Une distinction claire se dessine entre les voitures européennes et japonaises d'une part, et les américaines d'autre part.
- Les voitures européennes et japonaises sont regroupées dans la même classe, ce qui suggère qu'elles présentent des caractéristiques similaires.
- Les voitures américaines sont plus dispersées entre les deux classes, ce qui indique qu'elles présentent une plus grande diversité en termes de caractéristiques.

3 classes:

- Cette analyse permet d'affiner la distinction entre les différentes origines de voitures.
- La classe 1 regroupe principalement les voitures européennes et japonaises, ce qui confirme qu'elles présentent des caractéristiques similaires.
- La classe 2 regroupe principalement les voitures américaines ayant des caractéristiques plus proches de celles des voitures européennes et japonaises.
- La classe 3 regroupe principalement les voitures américaines ayant des caractéristiques plus distinctes de celles des voitures européennes et japonaises.

4 classes:

- Cette analyse permet d'obtenir une segmentation plus fine des différentes origines de voitures.
- Les classes 1 et 3 regroupent principalement les voitures européennes et japonaises, avec une distinction plus fine en fonction de leurs caractéristiques spécifiques.
- La classe 2 regroupe les voitures américaines ayant des caractéristiques plus proches de celles des voitures européennes et japonaises.
- La classe 4 regroupe les voitures américaines ayant des caractéristiques plus distinctes de celles des voitures européennes et japonaises.

5.2.5 Conclusion

L'analyse par CAH met en évidence une distinction claire entre les différentes origines de voitures en termes de caractéristiques. Les européennes et japonaises présentent des similarités plus marquées que les américaines, qui se distinguent par une plus grande diversité.

Le choix du nombre optimal de clusters dépend des objectifs de l'analyse et de l'interprétation souhaitée. Si l'objectif est de simplement distinguer les origines, une analyse avec 2 classes peut suffire. Si l'objectif est d'affiner la distinction, une analyse avec 3 ou 4 classes peut être plus appropriée.

Il est important de noter que la CAH est une méthode exploratoire et que les résultats doivent être interprétés avec prudence. Des analyses complémentaires, telles que des analyses statistiques, pourraient être menées pour confirmer les conclusions tirées.

6 ACP [Annexe : COR_ACP_figures.pdf]

L'Analyse en Composantes Principales (ACP) nous permet de simplifier la complexité des données en réduisant le nombre de variables. Nous partons des corrélations entre les variables pour les résumer dans des axes principaux qui capturent l'essence des informations, nous aidant à identifier les grandes tendances.

6.1 Réduction de Dimension et Interprétation

L'ACP nous permet de simplifier la structure de nos données en transformant nos sept variables quantitatives en axes principaux. Comme l'illustre la (Figure 5), ces axes fournissent une visualisation des relations entre les variables. Les variables **hp**, **cylinders**, **cubicinches** et **weightlbs** sont étroitement liées et pointent dans la même direction, soulignant leur association avec la puissance et la taille des véhicules.

En revanche, **mpg** et **time.to.60** sont proches l'une de l'autre, indiquant une corrélation positive entre ces deux mesures. Cependant, elles sont en opposition avec le groupe de variables liées à la puissance et à la taille du véhicule, indiquant une relation négative.

La variable **year**, positionnée perpendiculairement aux autres, ne montre pas de corrélation directe avec les autres mesures, ce qui signifie que l'évolution temporelle des véhicules est relativement indépendante des autres caractéristiques mesurées ici.

6.2 Importance des Axes Principaux

Les deux premiers axes de l'ACP capturent une portion significative de l'information totale. Le premier axe représente **72.4%** de l'inertie totale, soulignant son importance pour comprendre les caractéristiques dominantes des voitures, telles que la puissance et le poids. Le deuxième axe, représentant **13%**, se concentre davantage sur l'évolution temporelle des véhicules, capturant les tendances liées aux années de fabrication. Ensemble, ces deux axes restituent environ **85%** de l'inertie totale, permettant une interprétation globale et fiable des données principales sans s'attarder sur les **15%** restants qui pourraient concerner des variations moins explicatives.

6.3 Contribution des Variables aux Axes

La (Figure 6) montre comment chaque variable contribue aux axes de l'ACP. Les caractéristiques techniques du moteur dominent le premier axe, alors que l'année (year) affecte principalement le deuxième axe, révélant des informations sur l'âge des véhicules.

6.4 Qualité de Représentation des Variables sur les deux Axes Factoriel

(Figure 7) montre comment chaque variable est bien représentée sur les deux premiers axes de l'ACP. Si une variable a un Cos2 élevé, cela signifie que les axes captent bien ses informations.

- Les variables **cylindres**, **cubicinches**, **horsepower (hp)**, **weightlbs (poids)** et **mpg (miles per gallon)** ont des Cos2 très hauts, montrant qu'elles sont bien représentées par le premier axe.
- La variable **time.to.60 (temps pour atteindre 60 mph)** a un Cos2 plus faible, ce qui signifie qu'elle n'est pas aussi bien représentée par les deux premiers axes.
- La variable **year** a le Cos2 le plus élevé.

(Figure 6) montre la contribution de chaque variable aux dimensions de l'ACP. Les variables **hp**, **weightlbs**, **cubicinches**, **horsepower (hp)** et **mpg (miles per gallon)** dominent le premier axe, alors que **year** influence surtout le deuxième axe, montrant l'évolution des voitures avec le temps.

6.5 Visualisation des Clusters de Marques

(Figure 8) illustre clairement comment les marques de voitures américaines, européennes et japonaises sont regroupées dans l'espace factoriel de l'ACP, montrant leurs différences distinctes.

- Une grande partie des voitures américaines, sont situées à droite du cercle de corrélation, ce qui indique des valeurs élevées pour les variables **cylindres**, **cubicinches**, **hp** et **weightlbs**. En revanche, elles présentent des valeurs plus basses pour les variables **time.to.60** et **mpg**.
- Les marques européennes et japonaises, regroupées à gauche, montrent des similitudes entre elles. Ces véhicules ont généralement des valeurs inférieures pour les variables **cylindres**, **cubicinches**, **hp** et **weightlbs**, mais présentent des valeurs plus hautes pour les variables **time.to.60** et **mpg**.
- une proportion des voitures américaines est positionnée à proximité des voitures européennes et japonaises, suggérant ainsi des similitudes.

6.6 Conclusion sur l'ACP

Les résultats montrent que les caractéristiques des voitures européennes et japonaises sont assez similaires, notamment en termes de poids plus léger, d'un nombre de cylindres inférieur, et de puissance réduite, contrairement aux voitures américaines qui sont généralement plus lourdes, avec un nombre de cylindres supérieur et plus de puissance. De plus, les voitures européennes et japonaises consomment généralement moins de carburant et mettent plus de temps à atteindre 60 miles par heure par rapport aux voitures américaines, ce qui reflète des différences significatives dans la conception et la performance.

7 Conclusion

Ce projet a permis d'explorer en profondeur les caractéristiques et les tendances des véhicules de différentes régions, en utilisant des méthodes d'analyse de données avancées telles que l'Analyse en Composantes Principales (ACP), la Classification Ascendante Hiérarchique (CAH) et le k-moyenne. Nos analyses révèlent des distinctions claires entre les véhicules américains et ceux provenant du Japon et de l'Europe, ce qui souligne des approches variées en matière de design automobile, de performance, et d'efficacité énergétique.

Les voitures américaines sont généralement caractérisées par une plus grande puissance et un poids supérieur, ce qui affecte leur consommation de carburant et leur rapidité d'accélération. À l'opposé, les voitures européennes et japonaises tendent à favoriser l'efficacité énergétique et la performance en termes de consommation de carburant, ce qui les rend plus adaptées aux préoccupations actuelles en matière de durabilité environnementale.