

# Nombre de groupes

Choix du nombre de groupes

---

Véronique Tremblay

# Stratégie pour le choix du nombre de groupes

---

Pour la segmentation de clientèle

1. Regarder le dendogramme
2. Couper à un endroit qui a du sens
3. Analyser les segments
4. Modifier les groupes

Il y a plusieurs autres indicateurs<sup>1</sup>.

La librairie NbClust<sup>2</sup> en contient une trentaine.

---

<sup>1</sup>Surtout si les variables sont continues

<sup>2</sup>Pour variables continues

	KL	CH	Hartigan	CCC	Scott	Marriot
[1,]	6	6	3	6	3	3

	TrCovW	TraceW	Friedman	Rubin	Cindex	DB
[1,]	3	3	5	6	5	7

	Silhouette	Duda	PseudoT2	Beale	Ratkowsky	Ball
[1,]	6	6	6	6	3	3

	PtBiserial	Frey	McClain	Dunn	Hubert	SDindex
[1,]	6	1	2	6	0	5

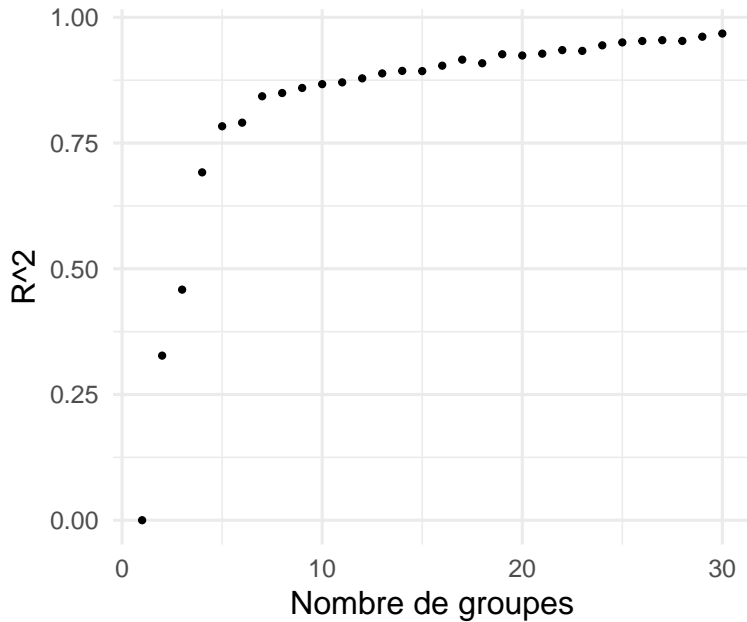
# Les indicateurs basés sur l'inertie

---

$$I_{totale} = I_{intra-groupe} + I_{inter-groupe}$$

- Ces indicateurs sont plus pertinents avec des variables continues
- Prendre garde au poids des variables et à la standardisation

$$\text{Pseudo-}R^2 = \frac{I_{inter-groupe}}{I_{total}}$$



$$CH = \frac{I_{inter-groupe}/(k-1)}{I_{intra-groupe}/(n-k)}$$



On maximise l'indice suivant:

$$D = \frac{\text{Distance minimale entre 2 groupes}}{\text{Distance maximale dans un groupe}}$$

L'indice de Dunn cherche donc à créer des groupes denses et bien séparés.

# L'indice de Silhouette

---

La silhouette de l'observation  $i$  mesure la confiance dans le choix du groupe pour l'observation  $i$ :

$$S(i) = \frac{b_i - a_i}{\max(b_i, a_i)}$$

où

- $a_i$  est la distance moyenne entre l'observation  $i$  et les autres observations de son groupe
- $b_i$  est la distance moyenne entre l'observation  $i$  et les observations du groupe le plus proche de  $i$

On souhaite maximiser la silhouette moyenne des observations.

# L'indice de Silhouette

---

```
s_moy <- silhouette(gr_m_moy,d)
```

```
mean(s_moy[, 'sil_width'])
```

```
## [1] 0.4012718
```

```
s_ppv <- silhouette(gr_m_ppv,d)
```

```
mean(s_ppv[, 'sil_width'])
```

```
## [1] 0.210907
```

```
s_vpd <- silhouette(gr_m_vpd,d)
```

```
mean(s_vpd[, 'sil_width'])
```

```
## [1] 0.39881
```

1. Tester différentes approches
2. Avec différents sous-échantillons
3. Vérifier si les groupes sont stables d'un essai à l'autre

```
table(gr_m_moy, gr_m_ppv)
```

```
##           gr_m_ppv
## gr_m_moy  1  2  3  4  5  6
##           1 11  0  0  0  0  0
##           2  0 11  0  0  0  0
##           3  0  0  3  0  0  0
##           4  0  0 11  0  0  1
##           5  0  6  0  1  0  0
##           6  0  0  0  0  3  0
```

## Quelques conseils pour le choix des groupes

---

- Privilégiez l'interprétabilité et l'utilité des groupes à un critère quelconque.
- Évitez les critères basés sur l'inertie ou la variance pour des groupes de taille et d'étendue inégales.

### Attention

- Est-ce qu'un petit groupe est nécessairement inintéressant?

## Indices liés aux probabilités estimées

- la probabilité a posteriori d'appartenir à chaque classe,  $\hat{\gamma}_{ki}$
- l'indice d'entropie calculé à partir des  $\hat{\gamma}_{ki}$

AIC, BIC, Tests d'hypothèses