

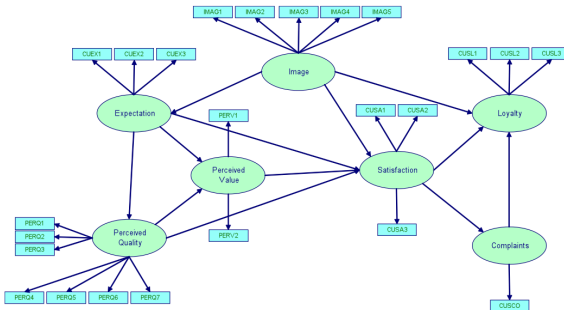
Algorithme et critères du PLS Path Modeling

Giorgio Russolillo
giorgio.russolillo@cnam.fr

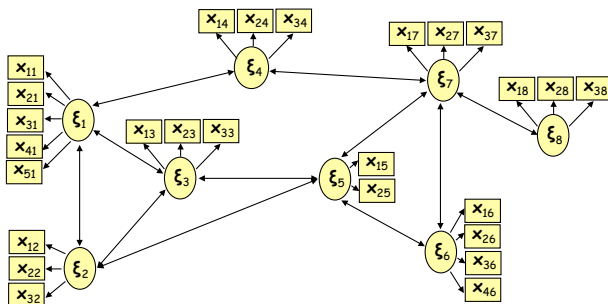
Laboratoire CEDRIC, Conservatoire National des Arts et Métiers, France

PLS Path Modeling

Le PLS-PM permet la modélisation d'un réseau de relations statistiques entre variables latentes indirectement mesurés au moyen d'ensembles d'indicateurs (les variables manifestes)



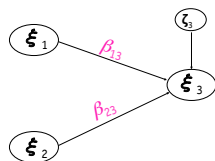
PLS Path Modeling: notations



- P variables manifestes (MV) observées sur n individus $\Rightarrow x_{pq}$ MV générique
- Q blocs de variables manifestes. Le bloc q est composé par P_q variables manifestes $\Rightarrow \sum_{q=1}^Q P_q = P$
- Q variables latentes (LV) $\Rightarrow \xi_q$ LV générique

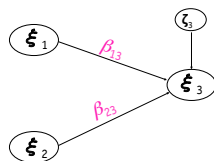
Le modèle structurel

Le **modèle structurel** décrit les relations entre les **variables latentes**.



Le modèle structurel

Le **modèle structurel** décrit les relations entre les **variables latentes**.



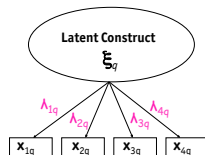
Pour chaque LV endogène le modèle structurel peut être écrit comme :

$$\xi_j = \sum_{q: \xi_q \rightarrow \xi_j} \beta_{qj} \xi_q + \zeta_j$$

where β_{qj} est le path-coefficient qui relie la q -ième LV à la j -ième LV endogène

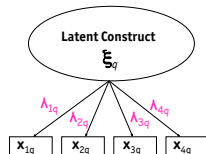
Le modèle de mesure

Le **modèle de mesure** décrit les relations entre **les variables manifestes** et la **variable latente correspondante**.



Le modèle de mesure

Le **modèle de mesure** décrit les relations entre les **variables manifestes** et la **variable latente** correspondante.



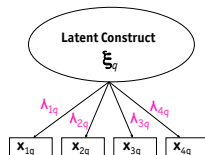
Pour chaque variable manifeste le modèle de mesure peut être écrit comme :

$$x_{pq} = \lambda_{pq}\xi_q + \epsilon_{pq}$$

où λ_{pq} est le loading (coefficient de régression de la q -ième LV sur la p -ième MV)

Le modèle de mesure

Le **modèle de mesure** décrit les relations entre les **variables manifestes** et la **variable latente** correspondante.



Pour chaque variable manifeste le modèle de mesure peut être écrit comme:

$$x_{pq} = \lambda_{pq}\xi_q + \epsilon_{pq}$$

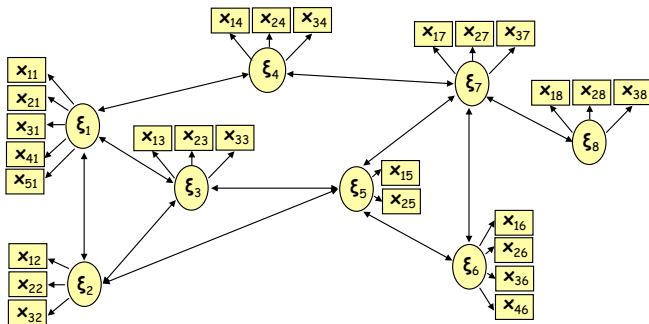
où λ_{pq} est le loading (coefficient de régression de la q -ième LV sur la p -ième MV)

L'approche PLS est basée sur les composantes: la LV est estimée sous forme d'une combinaison linéaire des MV de son bloc ([weight relation](#)):

$$\hat{\xi}_q = \sum_{p=1}^{P_q} w_{pq}x_{pq}$$

PLS Path Modeling: l'algorithme

- 1 Obtenir Q vecteurs de poids w_q travers la boucle itérative PLS
- 2 Calculer les scores des variables latentes
- 3 Calculer les loadings et les path coefficients



PLS Path Modeling: l'algorithme

Le but de la boucle PLS est de définir un **système de pondérations** appliquer chaque bloc de MV en vue d'estimer les LV correspondantes, selon la relation:

$$\hat{\xi}_q = \sum_{p=1}^{P_q} w_{pq} X_{pq}$$

PLS Path Modeling: l'algorithme

Le but de la boucle PLS est de définir un **système de pondérations** appliquer chaque bloc de MV en vue d'estimer les LV correspondantes, selon la relation:

$$\hat{\xi}_q = \sum_{p=1}^{P_q} w_{pq} X_{pq}$$

L'algorithme itératif est basé sur deux étapes principales; dans chaque étape la variable latente est approximée par une **proxy**.

- **Estimation externe**

⇒ t_q = somme pondérée des ses VM

- **Estimation interne**

⇒ z_q = somme pondérée des estimations externes des LV reliées

L'estimation externe

estimation externe

La proxy t_q est une somme pondérée des ses VM

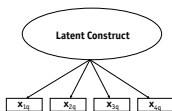
$$t_q \propto \pm \sum_{p=1}^{P_q} w_{pq} x_{pq} = \pm X_q w_q$$

L'estimation externe

estimation externe

La proxy t_q est une somme pondérée des ses VM

$$t_q \propto \pm \sum_{p=1}^{P_q} w_{pq} x_{pq} = \pm X_q w_q$$



Mode A (flèches orientées vers l'extérieur - reflective):

$$w_{pq} = X_q^T z_q / (z_q^T z_q)$$

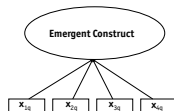
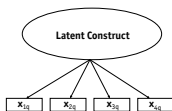
- Les indicateurs devraient **covier**
- Plusieurs régressions simples
- Variance expliquée (**AVE**, communality)

L'estimation externe

estimation externe

La proxy t_q est une somme pondérée des ses VM

$$t_q \propto \pm \sum_{p=1}^{P_q} w_{pq} x_{pq} = \pm X_q w_q$$



Mode A (flèches orientées vers l'extérieur - **reflective**): **Mode B** (flèches orientées vers l'intérieur - **formative**):

$$w_{pq} = X_q^T z_q / (z_q^T z_q)$$

- Les indicateurs devraient **covier**
- Plusieurs régressions simples
- Variance expliquée (**AVE**, communality)

$$w_{pq} = (X_q^T X_q)^{-1} X_q^T z_q$$

- Les indicateurs ne devraient pas **covier**
- Une régression multiple (multicollinéarité?)
- Structural Predictions (R^2 plus élevées pour les variables latentes endogènes)

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$

1 Schéma Centroïde

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$

- 1 Schéma Centroïde
- 2 Schéma Factoriel

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$

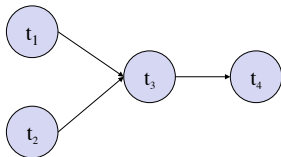
- 1 Schéma Centroïde
- 2 Schéma Factoriel
- 3 Schéma Structurel (ou Path weighting)

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$



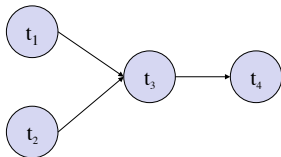
- ❶ Schéma Centroïde
- ❷ Schéma Factoriel
- ❸ Schéma Structurel (ou Path weighting)

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$



❶ Schéma Centroïde

$$z_3 = e_{31}t_1 + e_{32}t_2 + e_{34}t_4 \text{ where } e_{qq'} = \text{sign}(\text{cor}(t_q, t_{q'}))$$

❷ Schéma Factoriel

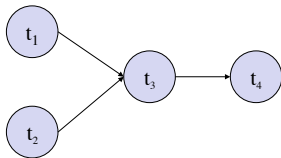
❸ Schéma Structurel (ou Path weighting)

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$



❶ Schéma Centroïde

$$z_3 = e_{31}t_1 + e_{32}t_2 + e_{34}t_4 \text{ where } e_{qq'} = \text{sign}(\text{cor}(t_q, t_{q'}))$$

❷ Schéma Factoriel

$$z_3 = \text{cor}(t_3, t_1)t_1 + \text{cor}(t_3, t_2)t_2 + \text{cor}(t_3, t_4)t_4$$

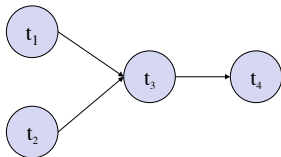
❸ Schéma Structurel (ou Path weighting)

L'estimation interne

Inner Estimation

La proxy z_q est une somme pondérée des estimations externes des LV reliées

$$z_q \propto \sum_{q'=1}^{Q'} e_{qq'} t_{q'}$$



❶ Schéma Centroïde

$$z_3 = e_{31}t_1 + e_{32}t_2 + e_{34}t_4 \text{ where } e_{qq'} = \text{sign}(\text{cor}(t_q, t_{q'}))$$

❷ Schéma Factoriel

$$z_3 = \text{cor}(t_3, t_1)t_1 + \text{cor}(t_3, t_2)t_2 + \text{cor}(t_3, t_4)t_4$$

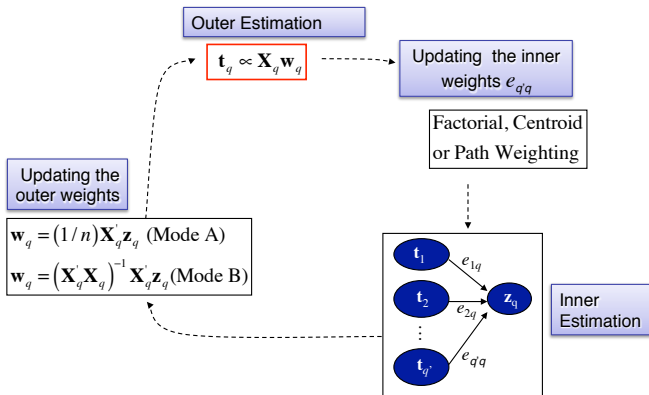
❸ Schéma Structurel (ou Path weighting)

$$z_3 = \beta_{31}t_1 + \beta_{32}t_2 + \text{cor}(t_3, t_4)t_4$$

où les betas sont les coefficients de régression du modèle:

$$t_3 = \beta_{31}t_1 + \beta_{32}t_2 + \zeta_3$$

PLS-PM Iteration



Critères d'optimisation de l'algorithme

Glang (1988) et Mathes (1993) ont montré que l'équation stationnaire du PLS-PM Mode B résout le critère:

PLS-PM Mode B

$$\arg \max_{\forall w_q} \sum_q c_{qq'} g(\text{cov}(X_q w_q, X_{q'} w_{q'})) \quad \text{s.t. } \|X_q w_q\| = 1$$

$$\begin{cases} c_{qq'} = 1 & \text{si } X_q \text{ and } X_{q'} \text{ sont reliés} \\ c_{qq'} = 0 & \text{sinon} \end{cases} \quad \begin{cases} g() = \text{carré (Schéma Factoriel)} \\ g() = \text{valeur absolue (Schéma Centroïde)} \end{cases}$$

Hanafi (2007) a prouvé que l'algorithme PLS dans ce cas est monotone convergente à ces critères.

Critères d'optimisation de l'algorithme

Dans sa thèse, Kramer (2007)

- 1 montre que l'algorithme PLS, quand le Mode A est utilisé pour tous les blocs, n'est pas basé sur une équation stationnaire liée à l'optimisation d'une fonction dérivable deux fois.
- 2 étend les résultats de Hanafi à un mode A modifié (connu dans la littérature comme **new mode A**), dans lequel une contrainte de normalisation est mise sur les poids extérieurs plutôt que sur les scores des variables latentes

Full New Mode A PLS-PM

$$\arg \max_{\forall w_q} \sum_q c_{qq} g(\text{cov}(X_q w_q, X_{q'} w_{q'})) \quad \text{s.t. } \|w_q\| = 1$$

Critères d'optimisation de l'algorithme

Dans sa thèse, Kramer (2007)

- 1 montre que l'algorithme PLS, quand le Mode A est utilisé pour tous les blocs, n'est pas basé sur une équation stationnaire liée à l'optimisation d'une fonction dérivable deux fois.
- 2 étend les résultats de Hanafi à un mode A modifié (connu dans la littérature comme **new mode A**), dans lequel une contrainte de normalisation est mise sur les poids extérieurs plutôt que sur les scores des variables latentes

Full New Mode A PLS-PM

$$\arg \max_{\forall w_q} \sum_q c_{qq} g(\text{cov}(X_q w_q, X_{q'} w_{q'})) \quad \text{s.t. } \|w_q\| = 1$$

Il a été montré de façon empirique que le mode A se rapproche de ce critère.

Critères d'optimisation de l'algorithme

Tenenhaus et Tenenhaus (2011) ont proposé ce critère qui unifie les deux critères précédents:

critère général

$$\begin{aligned} \arg \max_{\forall X_q, w_q} \sum_{q \neq q'} c_{qq'} g \{ \text{cov}(X_q w_q, X_{q'} w_{q'}) \} \\ \text{s.t. } \tau_q \|w_q\| + (1 - \tau_q) \text{var}(x_{pq} w_q) = 1 \end{aligned}$$

$$\begin{cases} c_{qq'} = 1 & \text{si } X_q \text{ et } X_{q'} \text{ sont reliés} \\ c_{qq'} = 0 & \text{sinon} \end{cases} \quad \begin{cases} g\{\cdot\} = \text{carré (Schéma Factoriel)} \\ g\{\cdot\} = \text{valeur absolue (Schéma Centroïde)} \end{cases} \quad \begin{cases} \tau_q = 1 & \text{si Mode A} \\ \tau_q = 0 & \text{si Mode B} \end{cases}$$

Conclusions

- Le PLS-PM est utilisé comme une méthode alternative pour estimer un SEM
- il est (presque) jamais sous-identifié (règles habituelles de régression)
- c'est une technique **distribution-free**
- En tant que technique d'estimation, PLS promet que la **consistence "at large"**.
- La **prédiction** (entendue comme reproduction de données) est considérée plus importante que l'estimation des paramètres
- Les **scores** des variables latentes sont directement obtenus.
- Les propriétés de l'algorithme dépendent du schéma et du mode utilisés.