

# IFT 3245

## Simulation et modèles

Fabian Bastin  
DIRO  
Université de Montréal

Automne 2012

# Variables de contrôle (VC)

Idée: exploiter de l'information auxiliaire pour faire une correction à l'estimateur.

On se restreint ici aux VCs linéaires.

Soit  $X$  un estimateur sans biais de  $\mu$  et  $\mathbf{C} = (C^{(1)}, \dots, C^{(q)})^T$  des VCs corrélées avec  $X$ , d'espérance connue  $E[\mathbf{C}] = \boldsymbol{\nu} = (\nu^{(1)}, \dots, \nu^{(q)})^T$ .

L'estimateur avec VC est:

$$X_c = X - \boldsymbol{\beta}^T(\mathbf{C} - \boldsymbol{\nu}) = X - \sum_{\ell=1}^q \beta_\ell(C^{(\ell)} - \nu^{(\ell)}),$$

où  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_q)^T$  (des constantes).

On a  $E[X_c] = E[X] = \mu$ .

# Choix de $\beta$

Comment choisir  $\beta$ ?

Soient  $\Sigma_c = \text{Cov}[\mathbf{C}]$  et  
 $\Sigma_{cX} = (\text{Cov}(X, C^{(1)}), \dots, \text{Cov}(X, C^{(q)}))^T$ .

## Hypothèse (CV1)

$\text{Var}[X] = \sigma^2 < \infty$ ,  $\Sigma_c$  et  $\Sigma_{cX}$  sont finies, et  $\Sigma_c$  est définie positive (et donc inversible).

On a alors

$$\text{Var}[X_c] = \text{Var}[X] + \beta^T \Sigma_c \beta - 2\beta^T \Sigma_{cX}.$$

# Choix de $\beta$

Pour minimiser par rapport à  $\beta$ , il suffit d'annuler le gradient par rapport à  $\beta$ :

$$0 = \nabla_{\beta} \text{Var}[X_c] = 2\boldsymbol{\Sigma}_c\beta - 2\boldsymbol{\Sigma}_{cX}.$$

Le minimum est donc atteint pour

$$\beta = \beta^* = \boldsymbol{\Sigma}_c^{-1}\boldsymbol{\Sigma}_{cX},$$

qui donne la variance minimale

$$\text{Var}[X_c] = (1 - R_{cX}^2)\text{Var}[X] \stackrel{\text{def}}{=} \sigma_c^2,$$

où

$$R_{cX}^2 = \frac{\boldsymbol{\Sigma}_{cX}^T \boldsymbol{\Sigma}_c^{-1} \boldsymbol{\Sigma}_{cX}}{\text{Var}[X]}$$

(le carré du coefficient de corrélation multiple entre  $C$  et  $X$ ) et la variance est réduite par le facteur  $1 - R_{cX}^2 = \sigma_c^2/\sigma^2$ . Mais avec  $\beta \neq \beta^*$ , la variance peut augmenter.

# Types de VCs

- (a) variables internes, basées sur des quantités déjà calculées durant la simulation;
- (b) variables externes, obtenues par des simulations additionnelles;
- (c) VCs implicites obtenues via une moyenne pondérée.  
Soient  $X^{(0)}, \dots, X^{(q)}$  des estimateurs sans biais de  $\mu$ .  
Posons

$$X_c = \sum_{\ell=0}^q \beta_\ell X^{(\ell)} = X^{(0)} - \sum_{\ell=1}^q \beta_\ell (X^{(0)} - X^{(\ell)})$$

où  $\sum_{\ell=0}^q \beta_\ell = 1$ .

On peut interpréter  $C^{(\ell)} = X^{(0)} - X^{(\ell)}$ ,  $\ell = 1, \dots, q$ , comme VC pour  $X = X^{(0)}$ .

# Estimation de $\beta^*$ : propriétés asymptotiques

En pratique, on ne connaît pas  $\beta^* = \Sigma_c^{-1} \Sigma_{cX}$  (parfois  $\Sigma_c$ , mais jamais  $\Sigma_{cX}$ ).

On peut l'estimer, disons par  $\hat{\beta}_n$ , calculé à partir de  $(X_1, \mathbf{C}_1), \dots, (X_n, \mathbf{C}_n)$ .

Posons

$$X_{ce,i} = X_i - \hat{\beta}_n^T (\mathbf{C}_i - \boldsymbol{\nu}),$$

et

$$\bar{X}_{ce,n} = \bar{X}_n - \hat{\beta}_n^T (\bar{\mathbf{C}}_n - \boldsymbol{\nu}).$$

# Estimation de $\beta^*$ : propriétés asymptotiques

## Theorème

Sous l'hypothèse CV1, lorsque  $n \rightarrow \infty$ , si  $\hat{\beta}_n \xrightarrow{D} \beta^*$ , alors

$$\sqrt{n}(\bar{X}_{c,n} - \bar{X}_{ce,n}) \xrightarrow{D} 0,$$

$$S_{ce,n}^2 \stackrel{def}{=} \frac{1}{n} \sum_{i=1}^n (X_{ce,i} - \bar{X}_{ce,n})^2 \xrightarrow{D} \sigma_c^2,$$

$$\frac{\sqrt{n}(\bar{X}_{ce,n} - \mu)}{S_{ce,n}} \xrightarrow{D} \frac{\sqrt{n}(\bar{X}_{c,n} - \mu)}{\sigma_c} \xrightarrow{D} N(0, 1).$$

On peut utiliser ce théorème pour calculer un IC pour  $\mu$ , en supposant que  $\sqrt{n}(\bar{X}_{ce,n} - \mu)/S_{ce,n} \sim N(0, 1)$ .

# Construction de $\hat{\beta}_n$

Méthode de base:

$$\hat{\beta}_n = \hat{\Sigma}_c^{-1} \hat{\Sigma}_{cX}$$

où les éléments de  $\hat{\Sigma}_c$  et  $\hat{\Sigma}_{cX}$  sont

$$\hat{\sigma}_c^{(\ell,k)} = \frac{1}{n-1} \sum_{i=1}^n (C_i^{(\ell)} - \bar{C}_n^{(\ell)})(C_i^{(k)} - \bar{C}_n^{(k)}),$$

$$\hat{\sigma}_{cX}^{(\ell)} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)(C_i^{(\ell)} - \bar{C}_n^{(\ell)}).$$

## Cas multinormal

Dans le cas où  $\begin{pmatrix} X_i \\ \mathbf{C}_i \end{pmatrix} \sim \text{normal}$ , on peut utiliser la théorie de la régression linéaire (avec estimateurs moindres carrés) pour le modèle

$$X = \mu + \beta^T(\mathbf{C} - \nu) + \epsilon$$

où  $\epsilon \sim N(0, \sigma_\epsilon^2)$ .

Si on définit

$$\tilde{S}_{ce,n}^2 = \frac{n}{n-q-1} \left( \frac{1}{n} + \frac{(\bar{\mathbf{C}}_n - \nu)^T \hat{\Sigma}_c^{-1} (\bar{\mathbf{C}}_n - \nu)}{n-1} \right) \sum_{i=1}^n (X_{ce,i} - \bar{X}_{ce,n})^2,$$

où  $\hat{\beta}_n$  utilise les covariances empiriques, on a le théorème ci-après.

## Theorème

Si les  $(X_i, \mathbf{C}_i^T)^T$  sont i.i.d normaux, alors

$$E[\bar{X}_{ce,n}] = \mu \quad (\text{aucun biais}),$$

$$E[\tilde{S}_{ce,n}^2/n] = Var[\bar{X}_{ce,n}] = \frac{n-2}{n-q-2}(1 - R_{cX}^2)Var[\bar{X}_n],$$

$$\sqrt{n}(\bar{X}_{ce,n} - \mu)/\tilde{S}_{ce,n} \sim Student(n-q-1) \quad (\text{loi exacte}).$$

Permet de calculer un IC avec couverture exacte pour  $n$  fini.

Facteur d'inflation de la variance  $(n-2)/(n-q-2) > 1$  dû à l'estimation de  $\beta^*$ . Si on a déjà  $q$  VC, l'ajout d'une nouvelle VC n'est rentable que si la valeur de  $(1 - R_{cX}^2)$  est réduite d'une fraction  $\geq 1/(n-q-2)$ .

# Expériences pilotes pour estimer $\beta^*$ ?

Pour avoir un estimateur sans biais de  $\beta^*$ , on peut faire une expérience pilote de  $n_0$  observations, calculer un estimateur  $\hat{\beta}_0$  et l'utiliser pour les  $n - n_0$  observations restantes. On obtient:

$$\bar{X}_{\text{cp},n} = \frac{1}{n - n_0} \sum_{i=n_0+1}^n (X_i - \hat{\beta}_0^T (\mathbf{C}_i - \boldsymbol{\nu})) \text{ et}$$

$$S_{\text{cp},n}^2 = \frac{1}{(n - n_0 - 1)} \sum_{i=n_0+1}^n (X_i - \hat{\beta}_0^T (\mathbf{C}_i - \boldsymbol{\nu}) - \bar{X}_{\text{cp},n})^2.$$

On a  $E[\bar{X}_{\text{cp},n}] = \mu$  et  $E[S_{\text{cp},n}^2 / (n - n_0)] = \text{Var}[\bar{X}_{\text{cp},n}]$ .

Mais sous l'hypothèse de normalité,

$$\frac{\text{Var}[\bar{X}_{\text{cp},n}]}{\text{Var}[\bar{X}_{\text{ce},n}]} = \frac{n(n - q - 2)(n_0 - 2)}{(n - n_0)(n - 2)(n_0 - q - 2)} > 1.$$

C'est donc inefficace.