

Introduction à la T-estimation

Eddis Amari (1)
22/11/16

Référence: Lucien Birgé - Model selection via testing: an alternative to (penalized) maximum likelihood estimators

But: Développer une méthode générale d'estimation, basée sur des tests, robuste aux erreurs de modèle et adaptée à la sélection de modèle.

- I Estimateur du maximum de vraisemblance
- II Construction des T-estimateurs
- III Application en régression donnée

I Estimateur du maximum de vraisemblance (E.M.V.) (2)

1) Défauts de l'E.M.V.

Au delà du fait que l'E.M.V. n'est pas intrinsèque, ce n'est pas un estimateur universel, au sens où il peut échouer au moins pour deux raisons:

* Comportement pénible de la vraisemblance:

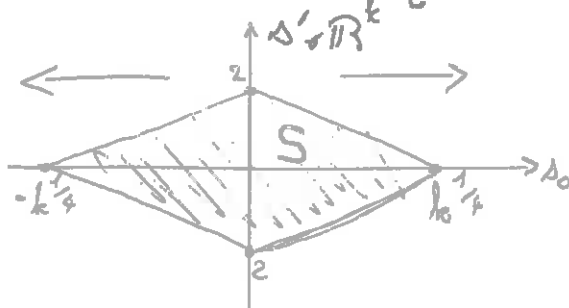
Pour le modèle de translation $f_\Delta(x) = f(x-\Delta)$, $\Delta \in \mathbb{R}$ où $f: \mathbb{R} \rightarrow \mathbb{R}_{>0}$ est une densité telle que $\lim_{x \rightarrow 0} f(x) = +\infty$. Si l'on observe un n -échantillon X_1, \dots, X_n , l'E.M.V. associé n'existe pas car la vraisemblance est non-bornée.

Même en restreignant $\Delta \in K$ compact, le problème est identique.

* Massivité de l'espace des paramètres:

Soit $X = (X_0, \dots, X_k) \sim N(\Delta, I_{k+1})$, avec $\Delta = (\Delta_0, \Delta_1, \dots, \Delta_k) = (\Delta_0, \Delta')$.

On se place dans le modèle $S = \left\{ \Delta \in \mathbb{R}^{k+1} \mid |\Delta_0| \leq k^{\frac{1}{2}}, \|\Delta'\|_2 \leq 2 \left(1 - \frac{|\Delta_0|}{k^{\frac{1}{2}}}\right) \right\}$



Pour $k \geq 128$,

a.g.p. $\hat{\Delta}_{EMV} = \left(0, \frac{2X'}{\|X'\|}\right)$, et $\sup_{\Delta \in S} E_\Delta \|\Delta - \hat{\Delta}_{EMV}\|^2 \geq \frac{3}{4} \sqrt{k+3}$.

Pourtant l'estimateur $\hat{\Delta} = (X_0, 0)$ donne le risque minimum

$$\inf_{\hat{\Delta}} \sup_{\Delta \in S} E_{\Delta} \|\Delta - \hat{\Delta}\|^2 \leq 5.$$

* Erreur dans le modèle :

De plus, même dans des modèles paramétriques compacts, une erreur dans le modèle peut conduire à une non-consistance.

Modèle : les lois uniformes \mathcal{U}_{θ} , $0 < \theta \leq 1$. On observe X_1, \dots, X_n iid de loi \bar{P} ayant pour densité $10 \left[\left(1 - \frac{2}{m}\right) \mathbb{1}_{\left[0, \frac{1}{10}\right]} + \frac{2}{m} \mathbb{1}_{\left[\frac{9}{10}, 1\right]} \right]$ par rapport à la mesure de Lebesgue sur $[0, 1]$.

Le modèle étant faux, on mesure les distances de manière intrinsèque avec Hellinger

$$h^2(P, Q) = \frac{1}{2} \int (\sqrt{dP} - \sqrt{dQ})^2$$

$$\rho(P, Q) = \int \sqrt{dP dQ} = 1 - h^2(P, Q)$$

L'E.M.V. est $\hat{\theta} = \max_{1 \leq i \leq n} X_i = X_{(n)}$. Avec proba. $\geq 1 - \left(1 - \frac{2}{m}\right)^m \geq 1 - e^{-2}$,

$X_{(n)} \geq \frac{9}{10}$, d'où $E_{\bar{P}} h^2(\bar{P}, \mathcal{U}_{X_{(n)}}) > 0.38$.

Or, un estimateur robuste devrait s'éloigner de la vitesse sans erreur de modèle $\left(\frac{1}{m}\right)$ du même ordre $O\left(\frac{1}{m}\right) = O\left(h^2(\bar{P}, \mathcal{U}_{\left[0, \frac{1}{10}\right]}\right)$.

2) Une ré. interprétation du maximum de vraisemblance

On observe un n -échantillon $X = (X_1, \dots, X_n)$ d'une distribution \bar{P}_s , où $s \in \mathcal{S}$. On prend une paramétrisation $\mathcal{S} \ni s \mapsto \bar{P}_s$ bijective et on munit \mathcal{S} de la métrique $h(s, t) = h(\bar{P}_s, \bar{P}_t)$. On confondra allégrement distributions et paramètres dans la suite.

Supposons que \mathcal{S} est compact. La famille $\{\bar{P}_s\}_{s \in \mathcal{S}}$ est alors dominée, on note $d\bar{P}_s$ la densité de \bar{P}_s .

$$\Lambda_n(t, X) = - \sum_{i=1}^n \log(d\bar{P}_t(X_i))$$

L'E.M.V. est $\hat{s} = \underset{\mathcal{S}}{\text{Argmin}} \Lambda_n(s, X)$ (supposé existant et unique)
 \uparrow
 $\mathcal{S} \subset \mathcal{S}$ fini par exemple.

Pour tous $t, u \in \mathcal{S}$,

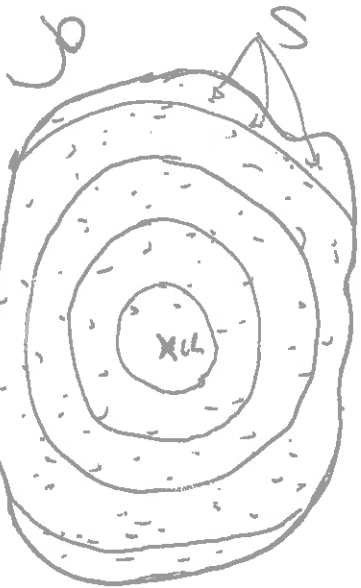
$$\begin{aligned}
\mathbb{P}_t(\Lambda_n(u, X) \leq \Lambda_n(t, X)) &= \mathbb{P}_t\left(\sum_{i=1}^n \log\left(\sqrt{\frac{d\bar{P}_u}{d\bar{P}_t}}(X_i)\right) \geq 0\right) \\
&\stackrel{\text{(Chebiff)}}{\leq} \mathbb{E}_t \exp\left(\sum_{i=1}^n \log\left(\sqrt{\frac{d\bar{P}_u}{d\bar{P}_t}}(X_i)\right)\right) \\
&= \exp\left(m \log(P(\bar{P}_u, \bar{P}_t))\right) \\
&= \log\left[\frac{\int \sqrt{d\bar{P}_u d\bar{P}_t} - 1}{-h^2(d\bar{P}_u, d\bar{P}_t)} + 1\right] \leq -h^2(\bar{P}_u, \bar{P}_t) \\
&\leq \exp(-m h^2(u, t))
\end{aligned}$$

Pour donner une borne uniforme, on somme par colonnes.

$$\text{Pour } K, \eta > 0, S_k = \{u \in S \mid 2^{k/2} K \eta \leq h(\Delta, u) < 2^{(k+1)/2} K \eta\}$$

$$\forall \Delta \in S,$$

$$P_\Delta [h(\Delta, \hat{\Delta}) \geq K \eta] \leq P_\Delta [\exists u \in S, h(\Delta, u) \geq K \eta, \Lambda_n(u, X) \leq \Lambda_n(\Delta, X)]$$



$$\leq \sum_{k \geq 0} P_\Delta [\exists u \in S_k, \Lambda_n(u, X) \leq \Lambda_n(\Delta, X)]$$

$$\leq \sum_{k \geq 0} |S_k| \sup_{u \in S_k} P_\Delta [\Lambda_n(u, X) \leq \Lambda_n(\Delta, X)]$$

$$\leq \sum_{k \geq 0} |S_k| \exp[-2^k m K^2 \eta^2]$$

Pour revenir au modèle complet \mathcal{Y} , on peut réquerir que S est un η -réseau.

Il existe alors $s' \in S$ avec $h(\Delta, s') \leq \eta$, et donc $P_\Delta [h(\Delta, \hat{\Delta}) \geq (K+1)\eta]$

$$\leq P_\Delta [h(s', \hat{\Delta}) \geq K\eta]$$

Mais borner cette quantité est impossible en toute généralité. Autrement dit, la robustesse est parfois impossible, car le test de rapport de vraisemblance entre deux paramètres, même éloignés en Hellinger, fait des erreurs trop importantes.

Changement de point de vue:

(6)

Pour $t \in S$, on note $\mathcal{R}_t = \{u \in S, \Lambda_n(u, X) \leq \Lambda_n(t, X)\}$
 $= \{u \in S, u \text{ est plus vraisemblable que } t\}$,

et $D_X(t) = \sup_{u \in \mathcal{R}_t} h(t, u)$. $D_X(t)$ s'intègre comme un indice de plausibilité: un $D_X(t)$ grand doit laisser penser que le vrai paramètre est loin de t .

Comme pour l'ETV, $\hat{\Delta} = \operatorname{argmin}_{t \in S} \Lambda_n(t, X) \Leftrightarrow D_X(\hat{\Delta}) = 0$,

on obtient:

$$\hat{\Delta} = \operatorname{argmin}_{t \in S} D_X(t).$$

Comme on a toujours $u \in \mathcal{R}_t$ ou $t \in \mathcal{R}_u$, $h(t, u) \leq D_X(t) \vee D_X(u)$.
En particulier, $h(t, \hat{\Delta}) \leq D_X(t) \vee D_X(\hat{\Delta}) \leq D_X(\hat{\Delta})$. Comme précédemment

$$\begin{aligned} \mathbb{P}_{\Delta} (h(\Delta, \hat{\Delta}) \geq (k+1)h) &\leq \mathbb{P}_{\Delta} [h(\Delta', \hat{\Delta}) \geq kh] \\ &\leq \mathbb{P}_{\Delta} [D_X(\Delta') \geq kh] \\ &\leq \mathbb{P}_{\Delta} [\exists u \in \mathcal{R}_{\Delta'} \text{ tel que } h(\Delta', u) \geq kh], \end{aligned}$$

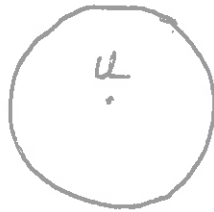
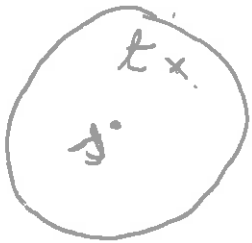
qui est sensiblement la même expression que plus haut.

Qu'a-t-on gagné?

↳ On peut désormais remplacer le test de rapport de vraisemblance (non robuste) entre t et u par une version robuste, au sens où elle vérifierait

$$\mathbb{P}_\Delta \left[\text{le test entre } t \text{ et } u \text{ choisit } u \right] \leq \exp \left[-c m h^2(t, u) \right]$$

$$\text{si } h(t, u) \geq \kappa h(\Delta, t)$$



Notations pour la suite:

- On observe une variable aléatoire $X \in \Xi$ de loi P_x .
- (\mathcal{M}, d) espace des paramètres, métrique (ou presque).
- $\Delta = F(P_x)$ est le paramètre d'intérêt.
- On identifie Δ et P_Δ : $d(\Delta, t) =: d(P_\Delta, P_t)$.
- Par "modèle", on entend un sous-ensemble $(S, S', \bar{S}, \mathcal{P})$ de \mathcal{M} , qui contient Δ ou non.

II Construction des T-estimateurs

(8)

Def: Pour $X \in \Xi$, $t \neq u \in \Pi$, un test entre t et u est un couple de fonctions $\Psi(t, u, X) = 1 - \Psi(u, t, X) \in \{0, 1\}$.

Convention: $\Psi(t, u, X) = 1 \iff$ Accepter u
"H₀" "H₁"

Vue d'ensemble: pour construire un T-estimateur, il faut:

- (i) Un ensemble $S \subset M$ approximant;
- (ii) Un $\varepsilon \geq 0$ et une fonction de poids $\eta: S \rightarrow \mathbb{R}_{>0}$;
- (iii) Une famille de tests entre les points de S .

Def: En notant $\forall u \neq t \in S$ $\Psi(t, u, X)$ le test entre t et u , et $R_+ = \{u \in S, \Psi(t, u, X) = 1\}$,

puis
$$D_X(t) = \begin{cases} \sup_{u \in R_+} d(t, u) & \text{si } R_+ \neq \emptyset \\ 0 & \text{si } R_+ = \emptyset \end{cases}$$

on appelle T-estimateur (ou T_ε -estimateur) tout $\hat{\Delta}(X) \in S$ tel que

$$D_X(\hat{\Delta}(X)) \vee \varepsilon \eta(\hat{\Delta}(X)) = \inf_{t \in S} \{D_X(t) \vee \varepsilon \eta(t)\}$$

Lien avec les Π -estimateurs: Si $\gamma(\cdot, X)$ est un contracté sur S ,
et h une fonction poids, soit $\gamma(t, X) = \gamma'(t, X) + \tau h^2(t)$.

Les Π -tests issus de γ' et de pénalité τh sont

$$\Psi(t, u, X) = \begin{cases} 0 & \text{si } \gamma(t, X) > \gamma(u, X) \\ 1 & \text{si } \gamma(t, X) < \gamma(u, X) \end{cases}$$

Remarque: Par définition, $\forall t \in S, d(t, \hat{\Delta}(X)) \leq D_X(t) \vee \varepsilon h(t)$.

En particulier, si $\varepsilon = 0, \forall s \in M,$

$$d(s, \hat{\Delta}) \leq \inf_{t \in S} \{d(s, t) + D_X(t)\},$$

c'est-à-dire que $\hat{\Delta}$ réalise le meilleur compromis entre distance à S et plausibilité.

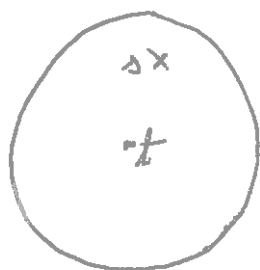
Hypotheses sur la famille de tests:

(10)

(A1) Il existe $\Pi_T \subset \Pi$, $\delta: M \times M_T \rightarrow [0, \infty]$, $a, B > 0$
tels que pour tout $u \neq t \in M_T$ et $x \in \mathbb{R}$, il existe un test $\psi(t, u, X)$
tel que

$$\sup_{\substack{\delta \in M \\ \delta(\delta, t) \leq d(t, u)}} P_{\delta} [\psi(t, u, X) = 1] \leq B \exp[-a(d^2(t, u) + x)]$$

$$\sup_{\substack{\delta \in M \\ \delta(\delta, u) \leq d(t, u)}} P_{\delta} [\psi(u, t, X) = 1] \leq B \exp[-a(d^2(t, u) - x)]$$



Devisin pour $\delta = \kappa \cdot d$

(A2) (Version contracte) $\psi(t, u, X) = 1 \Leftrightarrow \gamma'(t, X) - \gamma'(u, X) > \tau x$
 \hookrightarrow Hypothèse similaire

On prendra ensuite $x = h^2(u) - h^2(t)$

Hypothèses sur les modèles:

On mesurera la massivité des modèles avec la notion de D-modèle ("D" pour Discret et Dimension)

Def (D-modèle) Soient $\eta, D, B' > 0$ et $S' \subset \mathcal{M}$. On dit que S' est un D-modèle pour les paramètres η, D, B' lorsque

$$|S' \cap B_d(t, x\eta)| \leq B' \exp[Dx^2] \quad \forall x \geq 2, t \in \mathcal{M}.$$

De manière équivalente,

$$|S' \cap B_d(t, r)| \leq B' \exp\left[D\left(\frac{r}{\eta} \vee 2\right)^2\right] \quad \forall r > 0, t \in \mathcal{M}.$$

Remarque: le "2" est fixé par convention. L'idée est de regarder S à l'échelle $\approx \eta \times 2\eta$.

- Voir section 6 pour le lien avec la dimension métrique et la dimension métrique intérieure. Et autres...
- le carré " x^2 " peut paraître contre-intuitif si l'on se réfère à \mathbb{R}^D où il est inutile. Malgré tout, des " $\log x$ " apparaissent naturellement dans la construction de modèles approximatifs. Ne servirait-ce que pour des raisons techniques liées à l'entropie quand elle est utilisée. Cf section 6: Il semble être plus intéressant d'être optimal en "D" et de perdre sur le rayon des boules.

T-estimateurs basés sur un seul D-modèle :

Pas de pénalité via : $\varepsilon = 0$.

Théorème (3): (Simplifié) Sous l'hypothèse (A1), soit $S \subset M_T$ un D-modèle tel que $D \geq \frac{1}{2}$ et $2ah^2 \geq 3D$.

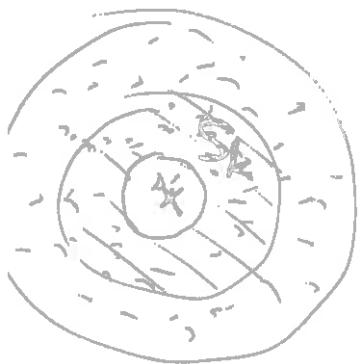
Pour $s \in \Pi$, si $\delta = \underset{>0}{\chi} d$ et $1 \leq q \leq \frac{16ah^2}{3D} \wedge 17$, on a Artificiel

$$\mathbb{E}_\Delta d^q(s, \hat{s}) \leq C_{q, B, B', \chi} \left(d(s, S) \vee \frac{4h}{\chi} \right)^q$$

Dem: Si $s' = \pi_S(s)$, $d(s, \hat{s}) \leq d(s, s') + \mathcal{D}_\chi(s')$
 $= d(s, S) + \mathcal{D}_\chi(s')$.

On borne $\mathbb{P}_\Delta[\mathcal{D}_\chi(s') > y]$ en sommant sur des couronnes centrées en s' , notées S_k :

$$\mathbb{P}_\Delta[\mathcal{D}_\chi(s') > y] = \mathbb{P}_\Delta[\exists t \in S, d(t, s') > y \text{ et } \Psi(s', t, X) = 1]$$



$$\leq \sum_{k \geq 0} \mathbb{P}_\Delta[\exists t \in S_k, \Psi(s', t, X) = 1]$$

$$\leq \sum_{k \geq 0} |S_k| \sup_{s \in S_k} \mathbb{P}_\Delta[\Psi(s', t, X) = 1]$$

Borné via
D-modèle

Borné via
Hypothèse (A1)

□

remarque: le pendant "contraire" fonctionne tout pareil.

T-estimateurs basés sur plusieurs D-modèles:

(13)

Ici, S est l'union de différents D-modèles.

(A3) $S = \bigcup_{m \in \mathcal{M}} S_m$ est l'union finie ou dénombrable de D-modèles S_m de paramètres η_m, D_m, B' , avec $D_m \geq \frac{1}{2}$.

Pour $t \in S$, on note $\eta(t) = \inf \{ \eta_m \mid m \in \mathcal{M} \text{ et } t \in S_m \}$
= le meilleur η_m disponible pour t .

On suppose que $\sum_{m \in \mathcal{M}} \exp(-a \eta_m^2 / 21) = \Sigma < \infty$

Théorème (5) (Simplifié) On suppose (A1) et (A3) avec $S \subset \mathcal{M}_T$,

$$a \eta_m^2 \geq 21 D_m / 5 \quad \forall m \in \mathcal{M}$$

Si $\varepsilon > 0$, $\delta = \kappa \varepsilon$ et $1 \leq q \leq 79$, $\forall s \in \mathcal{M}$,

$$\mathbb{E}_s \left[d_2^q(s, \hat{s}) \right] \leq C_{q, B, B', \Sigma, \kappa} \times \inf_{m \in \mathcal{M}} \left\{ d(s, S_m) \vee \frac{4 \eta_m}{\kappa} \right\}^q$$

Remarque: le pendant "contraste pénalisé" est analogue.

III Application en regression bornée

(11)

On observe un n -échantillon $(X_1, Y_1), \dots, (X_n, Y_n)$ où

$$\left\{ \begin{array}{l} \cdot X \in \mathcal{X} \\ \cdot Y_i = s(X_i) + \xi \\ \cdot Y, s(X) \in [0, 1] \end{array} \right. , \left\{ \begin{array}{l} X \sim \mu \text{ inconnue} \\ E_\mu[\xi | X] = 0 \end{array} \right.$$

On note, pour $s, t \in \{ \mathcal{X} \rightarrow [0, 1] \text{ mesurables} \} = \mathcal{M}$

$$d_2^2(t, s) = E_\mu \|s(X) - t(X)\|^2.$$

On rappelle que l'on a besoin $\left\{ \begin{array}{l} \text{d'une famille de tests} \\ \text{de D-modèles} \end{array} \right.$

Proposition (5) Pour tout $s, t, u \in \mathcal{M}$, $\eta, \gamma \in \mathbb{R}$, si $\gamma = 4 \|s - t\|^2 - \|t - u\|^2 / 4$,

$$\mathbb{P}_\mu \left[\chi'(t, X) - \chi'(u, X) \geq m \eta \right] \leq \exp \left[- \frac{3m}{100} \left(\|t - u\|^2 + \frac{98(\eta - \gamma)}{25} \right) \right],$$

où $\chi'(t, X) = \sum_{i=1}^n (Y_i - t(X_i))^2$

Dem Bernstein à $Z_i = (Y_i - t(X_i))^2 - (Y_i - u(X_i))^2$ qui vérifie

$$\left\{ \begin{array}{l} |Z_i| \leq 1 \\ |Z_i| \leq 2 |t(X_i) - u(X_i)| \end{array} \right.$$

À propos de la construction de D-modèles

Comme la loi de la covariable est inconnue, l'espace métrique à discrétiser est inconnu: $d_2^2(s, t) = E_{\mu=?} |t(x) - \mu(x)|^2$.

Par conséquent, la façon sûre de discrétiser est d'utiliser la métrique plus forte L^∞ .

Pour $\{m = (j, k) \in \mathbb{N}^2 \mid 3 \leq |m| \leq \frac{M}{10}\} = \mathcal{M} \ni m$

où $|m| = j(k+1)$,

considérons l'espace \mathcal{P}_m des polynômes par morceaux sur $[0, \frac{1}{j}]$, ..., $[\frac{k-1}{j}, 1]$ de degrés $\leq k$. \mathcal{P}_m est un espace vectoriel de dimension $|m|$: on construit un D-modèle $S_m \subset \mathcal{P}_m$ de paramètres:

$$h_m = \sqrt{\frac{18|m|}{n} \log\left(\frac{n}{|m|}\right)}; D_m = \frac{|m|}{4} \log\left(1 + \frac{2}{h_m}\right); B' = 1$$

Troncature à $[0, 1]$ et existence d'un réseau en dimension finie sur les compacts
 $\rightarrow d_2(s, S_m) \leq d_{\infty}(s, \mathcal{P}_m) + h_m$

Corollaire (8) On rappelle que si $h_m \geq \frac{140 D_m}{n}$ et $\sum \exp(-n h_m^2 / 100) = \sum < \infty$,

L'unique minimiseur de $\sum_{i=1}^n (X_i - t(X_i))^2 + \frac{25n}{9B} h^2(t)$ vérifie

$$E_0 \|D - \hat{D}\|^q \leq C_{q, B', \Sigma} \left\{ \left(\prod_{m \in \mathcal{M}} \|s - t\| \vee h_m \right) \right\}^q \text{ pour } 1 \leq q \leq 79.$$

Dans notre cas, les polynômes locaux fournissent:

$$\mathbb{E}_s \|\Delta - \hat{\Delta}\|^q \leq c_q \inf_{m \in \mathcal{M}} \left\{ d_\infty(s, \mathcal{P}_m) + \sqrt{\frac{|B| |m|}{n}} \log\left(\frac{n}{|m|}\right) \right\}^q.$$

Si $s \in \mathcal{L}^\beta(\mathbb{R})$ et dans une classe de Hölder, on a

$$d_\infty(s, \mathcal{P}_m) \leq C_\beta R j^{-\beta},$$

donc après optimisation en $m = (j, k) \in \mathcal{M}$, on obtient:

$$\mathbb{E}_s \|\Delta - \hat{\Delta}\|^2 \leq C_\beta R^{\frac{2}{2\beta+1}} \left(\frac{\log n}{n} \right)^{\frac{2\beta}{2\beta+1}}.$$

On perd le "log n" ici, mais $\hat{\Delta}$ est robuste