

Dependence, Correlation and Gaussianity in Independent Component Analysis

Cardoso – Résumé

De Lara, Tilquin, Vidal

1 Définitions préalables et propriétés

1.1 Définitions

Pour une variable aléatoire $Y \in \mathbb{R}^n$, on notera Y_i sa i -ième composante. Jusqu'à la fin, pour X une variable aléatoire à valeur dans \mathbb{R}^n , on notera $P(X)$ sa densité de probabilité et Σ_X sa matrice de covariance.

On posera aussi \mathcal{G} l'ensemble des distributions gaussiennes, \mathcal{P} l'ensemble des distributions « produits » (indiquant une indépendance des composantes). Pour une distribution P , on définira alors $P^{\mathcal{G}}$, $P^{\mathcal{P}}$ et $P^{\mathcal{G} \wedge \mathcal{P}}$ les distributions respectivement gaussiennes, produit et gaussienne produit minimisant la valeur de leur divergence par rapport à P . On verra ces distributions comme des projections sur ces différents espaces.

On définit alors les grandeurs suivantes.

La **divergence de Kullback–Leibler** de la distribution Q par rapport à P :

$$K(P \parallel Q) = \int_{\mathbb{R}^n} P(x) \log \frac{P(x)}{Q(x)} dx. \quad (1)$$

L'**entropie** de Y :

$$H(P) = - \int_{\mathbb{R}^n} P(x) \log P(x) dx. \quad (2)$$

L'**information mutuelle**, que l'on prendra pour mesure d'indépendance :

$$I(Y) = K(P(Y) \parallel \prod_i P_i(Y_i)) = K(P(Y) \parallel P(Y)^{\mathcal{P}}). \quad (3)$$

La **Non-Gaussianité** de Y :

$$G(Y) = K(Y \parallel \mathcal{N}(\mathbb{E}[Y], \Sigma_Y)) = K(P(Y) \parallel P(Y)^{\mathcal{G}}). \quad (4)$$

La **Corrélation** de Y :

$$\begin{aligned} C(Y) &= K(\mathcal{N}(\mathbb{E}[Y], \Sigma_Y) \parallel \mathcal{N}(\mathbb{E}[Y], \text{Diag } \Sigma_Y)) \\ &= K(P(Y)^{\mathcal{G}} \parallel P(Y)^{\mathcal{P} \wedge \mathcal{G}}) \\ &= \frac{1}{2} \log \frac{\det(\text{Diag}(\Sigma_Y))}{\det(\Sigma_Y)}. \end{aligned} \quad (5)$$

Dans la suite, on se permettra, pour éviter des notation trop lourde, d'écrire l'information mutuelle, la non-gaussianité et la corrélation d'une distribution.

1.2 Propriétés

Par propriété, si $K(P \parallel Q) = 0$, les distributions P et Q sont égales sur les espaces de mesures non nulles. Remarquons la propriété suivante, si T est une matrice inversible et μ est un vecteur quelconque, on a :

$$K(P(Y) \parallel P(Z)) = K(P(\mu + TY) \parallel P(\mu + TZ)). \quad (6)$$

Par définition, $I(Y) = 0$, on aura $P(Y) = \prod_i P_i(Y_i)$ et les composantes de Y seront indépendantes. On a la relation, d'où découle l'égalité de la définition de l'indépendance mutuelle

$$K\left(P(Y) \parallel \prod_i Q_i\right) = I(Y) + \sum_i K(P(Y_i) \parallel Q_i).$$

Si T est une matrice inversible et μ est un vecteur quelconque, on a l'égalité suivante :

$$G(P(Y)) = G(P(\mu + TY)). \quad (7)$$

1.3 Relations

On a les relations suivantes, découlant de la définition de la divergence de Kullback–Leibler ou du « théorème de Pythagore » :

$$I(Y) = \sum_i H(P_i) - H(P), \quad (8)$$

$$K(P(Y) \parallel P(Y)^{\mathcal{P} \wedge \mathcal{G}}) = I(Y) + \sum_i G(Y_i),$$

$$K(P(Y) \parallel P(Y)^{\mathcal{P} \wedge \mathcal{G}}) = G(Y) + C(Y),$$

$$\boxed{I(Y) + \sum_i G(Y_i) = G(Y) + C(Y)}. \quad (9)$$

Cette dernière relation nous donne que minimiser la grandeur $I(Y)$ revient à minimiser $C(Y) - \sum_i G(Y_i)$, la Non-gaussianité de Y étant indépendante ici du changement de référentiel que l'on recherche.

2 Géométrie

2.1 Théorème de Pythagore

On considère ici les deux variétés \mathcal{P} et \mathcal{G} , comme définies précédemment. Remarquons que ces deux variétés sont toutes les deux des familles exponentielles et ainsi, elles vérifient toutes les deux la propriétés suivantes. Si on y prend deux distributions p et q , alors toutes les distributions du segment exponentiel qu'elles définissent y appartiennent aussi, le segment exponentiel étant défini par

$$w_\alpha(x) = p(x)^{1-\alpha} q(x)^\alpha e^{-\psi(\alpha)},$$

avec $\psi(\alpha)$ un coefficient de normalisation.

On peut montrer, par ailleurs, que la famille exponentielle \mathcal{G} est de dimension $L_{\mathcal{G}} = n + \frac{1}{2}n(n+1)$ dans le sens où on peut trouver une mesure de référence $g(x)$ (non

nécessairement une distribution) et une « base » de L fonctions scalaire $S_l(x)$ telle que toute distribution s'écrive sous la forme :

$$p_\alpha(x) = g(x) \exp\left(\sum_l \alpha_l S_l(x) - \psi(\alpha)\right),$$

pour $\psi(\alpha)$ une fonction de normalisation et $\alpha \in \mathbb{R}^L$.

Remarquons que dans le cas de \mathcal{G} , on peut prendre comme « base » les fonctions $y \mapsto y_i$ et $y \mapsto y_i y_j$ pour i et j dans $\llbracket 1, n \rrbracket$.

Le Théorème de Pythagore s'énonce de la manière suivante :
Si \mathcal{E} est une famille exponentielle et P est une distribution quelconque, il existe alors une unique distribution $P^\mathcal{E}$ de \mathcal{E} vérifiant :

$$\forall Q \in \mathcal{E}, \quad K(P \parallel Q) = K(P \parallel P^\mathcal{E}) + K(P^\mathcal{E} \parallel Q).$$

On peut voir $P^\mathcal{E}$ comme la « projection orthogonale » de P sur la famille \mathcal{E} . Par positivité de la divergence de Kullback-Leibler, on peut voir $P^\mathcal{E}$ comme la distribution de \mathcal{E} minimisant sa divergence par rapport à P .

On peut alors remarquer que l'équation 9 peut être interprété comme suit : on obtient le même résultat en projetant P d'abord sur \mathcal{G} puis sur $\mathcal{P} \wedge \mathcal{G}$ ou d'abord sur \mathcal{P} puis sur $\mathcal{P} \wedge \mathcal{G}$. Remarquons que tous les termes de cette équations sont invariants par translation et changement d'échelle indépendamment sur les différentes coordonnées. L'espace $\mathcal{G} \wedge \mathcal{P}$ étant de dimension $2n$ (chaque coordonnée correspond à une gaussienne avec 2 paramètres), qui correspond exactement à la dimension de l'ensemble des transformations précédentes, tout se ramène exactement au cas d'une seule distribution dans $\mathcal{G} \wedge \mathcal{P}$.

2.2 Structures marginales

On va ici s'intéresser à l'espace « union » des deux variétés \mathcal{G} et \mathcal{P} que l'on notera $\mathcal{G} \vee \mathcal{P}$ qui correspondra à la plus petite famille exponentielle qui contient \mathcal{G} et \mathcal{P} , ce qui revient à considérer exactement toutes les distributions de la forme

$$p(y) = \phi(y) \exp\left(\sum_{l=1}^{L_\mathcal{G}} \alpha_l S_l(y) + \sum_{i=1}^n r_i(y_i) - \psi\right), \quad (10)$$

avec S_l une « base » de \mathcal{G} , r_i des fonctions réelles (puisque l'on a pas de « base » finie pour \mathcal{P}), ψ un coefficient de normalisation et $\phi(y)$ une distribution normale homogène ($\mathcal{N}(0, I_n)$).

Par le théorème de Pythagore, on trouve que pour toute distribution Q de \mathcal{G} ou de \mathcal{P} , on a :

$$K(P \parallel Q) = K(P \parallel P^{\mathcal{G} \vee \mathcal{P}}) + K(P^{\mathcal{G} \vee \mathcal{P}} \parallel Q). \quad (11)$$

Ainsi, la divergence minimum sur \mathcal{G} par rapport à P et par rapport à $P^{\mathcal{G} \vee \mathcal{P}}$ est atteinte au même point $P^\mathcal{G}$, ce qui revient à dire que :

$$(P^{\mathcal{G} \vee \mathcal{P}})^\mathcal{G} = P^\mathcal{G},$$

de même pour $P^\mathcal{P}$.

Ainsi, en reprenant l'équation 11 et en l'utilisant avec $Q = P^\mathcal{P}$ et avec $Q = P^\mathcal{G}$, on trouve les relations suivantes :

$$\begin{aligned} I(P) &= K(P \parallel P^{\mathcal{G} \vee \mathcal{P}}) + I(P^{\mathcal{G} \vee \mathcal{P}}) \\ G(P) &= K(P \parallel P^{\mathcal{G} \vee \mathcal{P}}) + G(P^{\mathcal{G} \vee \mathcal{P}}) \end{aligned}$$

En supposant que les valeurs des divergences sont suffisamment petites, on peut supposer que la figure formé par $P^{\mathcal{G} \vee \mathcal{P}}$, $P^{\mathcal{G}}$, $P^{\mathcal{P}}$ et $P^{\mathcal{G} \wedge \mathcal{P}}$ et un rectangle. L'égalité des longueurs nous donne alors :

$$I(P^{\mathcal{P} \vee \mathcal{G}}) \simeq C(P) \quad \text{et} \quad G(P^{\mathcal{P} \vee \mathcal{G}}) \simeq \sum_i G(P_i).$$

Cela donne, avec les équations précédentes :

$$\begin{aligned} I(P) &\simeq K(P \parallel P^{\mathcal{G} \vee \mathcal{P}}) + C(P) \\ G(P) &\simeq K(P \parallel P^{\mathcal{G} \vee \mathcal{P}}) + \sum_i G(P_i) \end{aligned}$$

La distribution $P^{\mathcal{G} \vee \mathcal{P}}$ peut être interprété comme la distribution la plus simple approchant la distribution P , dans le sens où elle capture la structure marginale de P et sa structure de premier et second ordre (voir équation 10).

3 Cumulant et géométrie locale

Pour rappel, les cumulants κ_n d'une variable aléatoire X sont définis avec la fonction génératrice des cumulants :

$$g(t) = \log \mathbb{E}[e^{tX}] = \sum_{n=1}^{\infty} \frac{\kappa_n}{n!} t^n \quad (12)$$

On s'intéresse ici aux distributions au voisinage de $P^{\mathcal{G} \wedge \mathcal{P}}$, ce qui correspond aux distributions faiblement corrélées et faiblement non-gaussiennes. On assimilera les variétés \mathcal{P} et \mathcal{G} à leur plan tangent et la divergence de Kullback-Leibler à une mesure quadratique.

3.1 Construction des plans tangents

Pour deux distributions $p(x)$ et $n(x)$, on définit la fonction

$$e_p(x) = \frac{p(x)}{n(x)} - 1,$$

qui sera alors d'espérance nulle selon $n(x)$: $\mathbb{E}_{X \sim n}[e_p(X)] = 0$.

On cherche alors à identifier les distributions p proches de n aux « petites » fonctions d'espérance nulle selon n . Si on considère une distribution proche de n , la fonction e_p sera *petite* et d'espérance nulle. Réciproquement, pour une fonction e_p données, on considérera alors $p(x) = n(x)(e_p(x) + 1)$ qui sera bien une distribution proche de $n(x)$.

Ainsi, on peut identifier l'espace vectoriel des variables aléatoires d'espérance nulle et de variance finie au plan tangent à la variété des distributions au point n .

Si on considère deux distributions p et q proches de n , on peut exprimer K_n l'expansion au second ordre suivant e_p et e_q de $K(p \parallel q)$ de la manière suivante :

$$K_n(p \parallel q) = \frac{1}{2} \mathbb{E}_{X \sim n}[(e_p(X) - e_q(X))^2] \quad (13)$$

3.2 Expansion de Gram-Charlier

3.3 Base d'Hermite

3.4 Approximation de la divergence par des petits cumulants

3.5 Décomposition en quatre parties

3.6 Divergences et objectifs