

Détection quasi-optimale d'informations cachées basée sur un modèle local non-linéaire

Rémi COGRANNE, Cathel ZITZMANN, Lionel FILLATRE, Florent RETRAINT, Igor NIKIFOROV, Philippe CORNU

ICD - LM2S - Université de Technologie de Troyes (UTT) - UMR STMR - CNRS
12, rue Marie Curie - B.P. 2060 - 10010 Troyes cedex - France
prenom.nom@utt.fr

Résumé – Cet article aborde le problème de détection fiable d'informations cachées dans les images naturelles non-compressées. L'objectif est de définir un test dont les probabilités d'erreurs sont analytiquement prédictibles afin de satisfaire une contrainte de faibles taux de fausses alarmes. Dans ce but, la détection d'informations cachées est étudiée dans le cadre de la théorie des tests statistiques. Si la moyenne théorique des pixels est connue, le test optimal au sens de Neyman-Pearson est explicité sous une forme simple et ses performances statistiques sont établies. Dans la pratique, un modèle local non-linéaire des images est développé pour estimer précisément et rapidement cette moyenne. Une linéarisation de ce modèle est proposée et permet d'obtenir un test quasi-optimal, *i.e.* dont la perte d'optimalité est bornée et faible. Enfin, des expérimentations sur 9 000 images montrent la pertinence et la qualité de l'approche par rapport à l'état de l'art.

Abstract – This paper investigates the problem of reliable hidden information detection in uncompressed natural images. The goal is to design a test with analytically predictable probabilities of error to meet a low false alarm rate constraint. For this purpose, the hidden information detection is cast in the framework of hypothesis testing theory. Assuming that the pixels mean is known, the optimal test in Neyman-Pearson sense is explicitly given in a simple form and its statistical performances are established. In practice, a non-linear local model of natural images is developed to estimate these means accurately and quickly. A linearization of this model is proposed and permits to design an almost-optimal test, *i.e.* with a bounded and limited loss of optimality. Eventually, numerical experimentations with 9,000 images show the relevance and the quality of proposed approach compared to the state of art.

1 Introduction

La stéganographie concerne la dissimulation d'informations dans un médium hôte (image, son, vidéo, etc.). La stéganalyse concerne à l'inverse la détection des médias contenant de l'information cachée. Il existe désormais de nombreux logiciels libres de stéganographie, facilement accessibles et utilisables ; il est donc crucial pour les forces de l'ordre de disposer d'outils de détection fiables. Dans ce contexte opérationnel, il est important de pouvoir calculer analytiquement les probabilités d'erreurs (fausse alarme et non-détection) du détecteur afin de fournir un critère quantitatif de décision. En outre, l'obtention immédiate de résultats proscrit l'utilisation de méthodes d'apprentissage supervisé. Aucune méthode actuelle de stéganalyse ne répond à ces contraintes.

Cet article aborde donc le problème de la détection fiable d'informations cachées dans les images. Pour cela, il est nécessaire de modéliser et d'exploiter la redondance entre pixels voisins. Or, les images sont des signaux complexes, non-stationnaires, dont le contenu structuré (moyennes théoriques des pixels) peut nuire à la détection. La méthodologie proposée consiste à utiliser la théorie des tests d'hypothèses et à modéliser le contenu d'une image pour prendre en compte ces moyennes comme des paramètres de nuisance. Les apports sont les suivants :

- Le test optimal au sens Neyman-Pearson est calculé de façon théorique et ses performances sont établies asymptotiquement lorsque le nombre de pixels tend vers l'infini.
- Un modèle local non-linéaire des images est présenté et permet d'estimer efficacement les paramètres de nuisances inconnus en pratique.
- L'utilisation de ces estimations permet l'obtention d'un test quasi-optimal, *i.e.* dont la perte de puissance par rapport au test optimal est bornée et faible.

La méthodologie proposée est, dans cet article, appliquée à la méthode de stéganographie par substitution de LSB car elle demeure la plus couramment étudiée et, bien qu'elle soit la plus simple à mettre en œuvre, elle permet de mettre en évidence les difficultés théoriques majeures.

2 Formulation du problème

Soit le vecteur $\mathbf{z}=(z_1, \dots, z_M)^T$ représentant une image de M pixels en niveaux de gris : $z_m \in \mathcal{Z}=\{0, \dots, 2^b-1\}$. Chacun des pixels résulte de la quantification $z_m = \lfloor y_m \rfloor$ où $\lfloor \cdot \rfloor$ est l'opération de quantification linéaire avec un pas unitaire (arrondi entier) et $y_m \in \mathbb{R}_+$ est l'intensité du m -ième pixel enregistré par le capteur photosensible. La valeur y_m peut s'écrire

$y_m = \theta_m + \xi_m$ où θ_m est l'espérance mathématique de y_m (contenu structuré) et ξ_m représente le bruit d'acquisition $\xi_m \sim \mathcal{N}(0, \sigma_m)$, cf.[1]. La distribution de probabilité du pixel z_m d'une image de couverture $P_{\theta_m} = \{p_{\theta_m}[0], \dots, p_{\theta_m}[2^b - 1]\}$ est donnée $\forall k \in \mathcal{Z}$ par :

$$p_{\theta_m}[k] = \Phi\left(\frac{k + 1/2 - \theta_m}{\sigma_m}\right) - \Phi\left(\frac{k - 1/2 - \theta_m}{\sigma_m}\right),$$

avec $\Phi(x) = \int_{-\infty}^x \phi(u) du$ et $\phi(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du$.

L'impact de la stéganographie dépend du nombre de bits d'informations insérés par pixel appelé, en stéganographie, le taux d'insertion et noté R dans cet article. Sous les hypothèses classiques en stéganalyse [2, 3], la distribution de probabilité du pixel z_m après insertion avec un taux d'insertion R est donnée par [2] : $Q_{\theta_m}^R = \{q_{\theta_m}^R[0], \dots, q_{\theta_m}^R[2^b - 1]\}$ où $\forall k \in \mathcal{Z}$:

$$q_{\theta_m}^R[k] = \left(1 - \frac{R}{2}\right) q_k(\theta_m) + \frac{R}{2} q_{(k-\bar{k})}(\theta_m)$$

et \bar{k} symbolise l'entier k dont le LSB a été inversé [3], i.e., $\bar{k} = k + 1$ si k est pair et $\bar{k} = k - 1$ si k est impair.

Le problème de détection de d'informations cachées, ou stéganalyse, dans une image consiste donc à choisir entre les hypothèses :

$$\begin{cases} \mathcal{H}_0 = \{z_m \sim P_{\theta_m}, \forall m = 1, \dots, M\} \\ \mathcal{H}_1 = \{z_m \sim Q_{\theta_m}^R, \forall m = 1, \dots, M, \forall R \in \mathcal{R} \subset]0; 1]\} \end{cases} \quad (1)$$

en tenant compte du fait que dans le contexte de cet article, il est nécessaire de maîtriser la probabilité de fausse alarme. Dans le cas théorique où la distribution P_{θ_m} des pixels avant insertion est connue la principale difficulté provient du fait que l'hypothèse alternative \mathcal{H}_1 est composite, puisque le taux d'insertion n'est pas connu. Une solution idéale serait alors de trouver un test qui soit uniformément le plus puissant (UPP) dans la classe $\mathcal{K}_\alpha = \{\delta : \sup_{R \in \mathcal{R}} \Pr[\delta(\mathbf{z}) = \mathcal{H}_1 | \mathcal{H}_0] \leq \alpha\}$ des tests dont la probabilité de fausse alarme est majorée par α .

En pratique les paramètres de distribution θ_m ne sont pas connus. Ces derniers sont des paramètres de nuisances, ils interviennent dans la définition des hypothèses statistiques (1) mais ne présentent pas d'intérêt pour décider en faveur d'une des hypothèses. La difficulté supplémentaire est alors de construire un test permettant une prise de décision indépendamment des paramètres de nuisance θ_m .

3 Modèle d'image naturelle

Afin de modéliser simplement le contenu de l'image, \mathbf{z} est découpée en K segments statistiquement indépendants de $N > 0$ pixels contigus $\mathbf{z}_k = (z_{k,1}, \dots, z_{k,N})^T$ avec $KN = M$. Le modèle d'image utilisé dans cet article suppose que la scène est constituée d'objets distincts sur lesquels l'intensité lumineuse varie de façon continue. Ainsi l'intensité lumineuse θ_k du segment de la scène (portion de largeur négligeable) dont est issu le vecteur \mathbf{z}_k est continue par morceaux. En supposant, dans

un souci de clarté, que chaque segment possède au plus une discontinuité dont la position (si elle est présente) est notée t_k , alors θ_k peut s'écrire :

$$\theta_k(x) = \theta_k^c(x) + u_k \mathbf{1}(x - t_k)$$

où $\mathbf{1}$ est la fonction indicatrice $\mathbf{1}(u) = 1$ si $u > 0$ et $\mathbf{1}(u) = 0$ si $u \leq 0$ et θ_k^c est une fonction continue que l'on suppose approximée fidèlement par un polynôme de degré $p-1$. Lors de l'acquisition, le système optique modifie la scène imagée par un flou qui peut être approximé sur l'ensemble des segments par un filtre Gaussien de réponse $h(x) = \varsigma^{-1} \phi(x/\varsigma)$ où ς est un paramètre constant si l'on admet un système optique stationnaire.

Après échantillonnage, la moyenne théorique θ_k du segment \mathbf{z}_k s'écrit $\theta_k = \mathbf{H}\mathbf{s}_k + \mathbf{F}(\eta_k)u_k$ où \mathbf{s}_k sont les coefficients du polynômes, u_k est l'intensité de la discontinuité et $\eta_k = (t_k, \varsigma)^T$. La matrice \mathbf{H} de taille $N \times p$ est une base polynomiale de degré $p-1$ et le vecteur $\mathbf{F}(\eta_k)$ est le profil de la discontinuité filtrée (si elle est présente).

Supposons que l'on dispose d'une méthode permettant d'obtenir une estimation $\hat{\eta}_k$ de η_k (voir [4] pour un exemple d'estimateur). Pour estimer précisément et rapidement le paramètre de nuisance θ_k malgré la présence de la non-linéarité par rapport à η_k l'approche proposée dans [5, 6] est adaptée pour linéariser le modèle par :

$$\theta_k = \dot{\mathbf{G}}_k \mathbf{v}_k + o(\vartheta^2) \quad (2)$$

où $\mathbf{v}_k = (\mathbf{s}_k \mid u_k \mid u_k(\eta_k - \hat{\eta}_k))^T$, $\dot{\mathbf{G}}_k = (\mathbf{H} \mid \mathbf{F}(\hat{\eta}_k) \mid \dot{\mathbf{F}}(\hat{\eta}_k))$, $\dot{\mathbf{F}}(\hat{\eta}_k)$ est la matrice Jacobienne de $\mathbf{F}(\hat{\eta}_k)$ de taille $N \times 2$ et la notation $x = o(y)$ signifie que $x/y \rightarrow 0$ lorsque $y \rightarrow 0$. En supposant que la variance du bruit est constante par segment, la relation (2) permet alors d'obtenir une estimation linéaire de θ_k par maximum de vraisemblance :

$$\hat{\theta} = \dot{\mathbf{G}}_k (\dot{\mathbf{G}}_k^T \dot{\mathbf{G}}_k)^{-1} \dot{\mathbf{G}}_k^T.$$

4 Stéganalyse quasi-optimale

Dans le cas où θ_k , σ_k et R sont connus pour tout k , (1) se réduit à un problème de décision entre deux hypothèses simples. Il est alors connu que le test du Rapport de Vraisemblance (RV) est optimal au sens de Neyman-Pearson. On s'intéresse ici au cas délicat de faibles taux d'insertion et de forts bruits. Dans ces conditions, un développement en série de Taylor montre que le test du RV, noté δ , est donné par [7] :

$$\delta(\mathbf{z}) = \begin{cases} \mathcal{H}_0 & \text{si } \Lambda(\mathbf{z}) < \tau_\alpha, \\ \mathcal{H}_1 & \text{si } \Lambda(\mathbf{z}) \geq \tau_\alpha, \end{cases} \quad (3)$$

où

$$\Lambda(\mathbf{z}) = \sum_{k=1}^K \sum_{n=1}^N w_{k,n} (z_{k,n} - \bar{z}_{k,n}) (z_{k,n} - \theta_{k,n}) + o\left(\frac{R^2}{\bar{\sigma}^2}\right) \quad (4)$$

$$w_{k,n} = \frac{\bar{\sigma}}{\sigma_{k,n}^2 \sqrt{K \cdot N}} \quad \text{et} \quad \frac{1}{\bar{\sigma}^2} = \frac{1}{KN} \sum_{k=1}^K \sum_{n=1}^N \frac{1}{\sigma_{k,n}^2}.$$

Ici encore la notation $\bar{z}_{k,n}$ représente l'entier $z_{k,n}$ dont le LSB à été inversé. Le RV $\Lambda(\mathbf{z})$ ne dépendant pas de R , ce test peut donc être utilisé si R est inconnu. Dans la pratique les valeurs $\theta_{k,n}$ ne sont pas connues ; il est proposé de les remplacer par $\hat{\theta}_{k,n}$. Cela correspond ici au test du RV Généralisé (RVG), noté $\hat{\delta}$, dont la fonction de décision est donné par [7]:

$$\hat{\Lambda}(\mathbf{z}) = \sum_{k=1}^K \sum_{n=1}^N \hat{w}_k(z_{k,n} - \bar{z}_{k,n})(z_{k,n} - \hat{\theta}_{k,n}) + o\left(\frac{R^2}{\bar{\sigma}^2}\right) \quad (5)$$

où désormais

$$\hat{w}_k = \frac{\bar{\sigma}}{\sigma_k^2 \sqrt{K(N-p-d_k)}}$$

et $d_k \in \{0, 1\}$ représente le nombre de discontinuités du k -ième segment.

Les performances du test du RV (3) sont données par le théorème suivant.

Théorème 1. *Si $\bar{\sigma} \gg 1$ et $R \approx 0$ alors en négligeant le terme en $o(R^2/\bar{\sigma}^2)$ dans l'expression du RV (4), il découle du théorème de la limite centrale de Lindeberg que :*

$$\begin{cases} \Lambda(\mathbf{z}) \rightsquigarrow \mathcal{N}(0; 1) \\ \Lambda(\mathbf{z}) \rightsquigarrow \mathcal{N}\left(\frac{R}{2\bar{\sigma}}\sqrt{KN}; 1\right) \end{cases} \quad (6)$$

où \rightsquigarrow représente la convergence en distribution lorsque $KN \rightarrow \infty$.

Corollaire 1. *En fixant le seuil de décision $\tau_\alpha = \Phi^{-1}(1 - \alpha)$ il découle du théorème 1 que le test δ appartient à la classe \mathcal{K}_α et que sa puissance est donnée par :*

$$\beta_{max} = 1 - \Phi\left(\tau_\alpha - \frac{R\sqrt{KN}}{2\bar{\sigma}}\right) \quad \text{lorsque } K \rightarrow \infty. \quad (7)$$

De façon analogue, lorsque les paramètres $\theta_{k,n}$ ne sont pas connus, les performances du test du RVG $\hat{\delta}$ sont données par le théorème suivant. Dans un souci de simplicité, le cas d'une unique discontinuité par segment est traité bien que la validité de ce dernier peut-être étendu pour les cas de multiples discontinuités.

Théorème 2. *Supposons que $\forall k \in \{1, \dots, K\}$, $d_k \leq 1$ et que l'estimation $\hat{\eta}_k$ est telle que $\|\eta_k - \hat{\eta}_k\|_1 \leq \vartheta$. Si $\bar{\sigma} \gg 1$ et $R \approx 0$ alors en négligeant le terme en $o(R^2/\bar{\sigma}^2)$ dans l'expression du RV (5), il découle du théorème de la limite centrale de Lindeberg que :*

$$\begin{cases} \hat{\Lambda}(\mathbf{z}) \rightsquigarrow \mathcal{N}(0; 1 + b) \\ \hat{\Lambda}(\mathbf{z}) \rightsquigarrow \mathcal{N}\left(\frac{R}{2\bar{\sigma}}\sqrt{\kappa}; 1 + b\right) \end{cases} \quad (8)$$

où $\kappa = K(N-p) - 3 \sum_{k=1}^K d_k < K \cdot N$ et b est un biais dû aux erreurs d'estimation de $\hat{\eta}_k$ vérifiant :

$$b \leq b_{max} = \frac{o(\vartheta^2)}{\bar{\sigma}^2(N-p-3)}.$$

Corollaire 2. *En fixant le seuil de décision $\hat{\tau}_\alpha = \Phi^{-1}(1 - \alpha)\sqrt{1 + b_{max}}$, il découle du théorème 2 que le test $\hat{\delta} \in \mathcal{K}_\alpha$ et que sa puissance est bornée, lorsque $K \rightarrow \infty$, par :*

$$1 - \Phi\left(\frac{\hat{\tau}_\alpha - \frac{R\sqrt{\kappa}}{2\bar{\sigma}}}{1 + b_{max}}\right) \leq \hat{\beta} \leq 1 - \Phi\left(\hat{\tau}_\alpha - \frac{R\sqrt{\kappa}}{2\bar{\sigma}}\right) \quad (9)$$

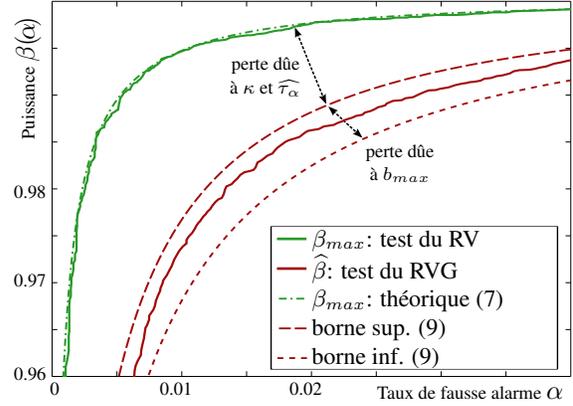


Figure 1: Comparaison des fonctions de puissance sur des données simulées.

Le résultat (7) offre une borne supérieure de la puissance de tout test statistique dont le taux de fausse alarme est majoré par α . Le résultat (9) montre que la perte d'optimalité du test du RVG est bornée et minimale si ϑ et $(KN - \kappa)$ sont suffisamment petits.

Enfin, il est intéressant de rapprocher les rapports de vraisemblance (4)-(5) aux méthodes *Weighted-Stego* (WS), initialement proposée dans [3], reposant sur la statistique suivante [9] :

$$\sum_{k=1}^K \sum_{n=1}^N w_{k,n}(z_{k,n} - \bar{z}_{k,n})(z_{k,n} - \hat{c}_{k,n})$$

où le poids $w_{k,n}$ est choisi pour que les zones texturées aient une importance moindre et $\hat{c}_{k,n}$ est une estimation de la valeur du pixel avant insertion $c_{k,n}$ obtenue par un filtrage autoregressif.

5 Résultats numériques

Les résultats numériques présentés dans cette section ont vocation à montrer la qualité de la méthodologie théorique proposée ainsi que sa pertinence en pratique.

La figure 1 présente les résultats d'une simulation Monte-Carlo des tests du RV et du RVG. La simulation a été répétée 25000 fois avec des images artificielles constituées de 400 segments de 32 pixels. Chacun de ces segments est généré avec une unique discontinuité de paramètres $u_k=96$, $\varsigma = 1.75$, $\vartheta = 1$ et un polynôme d'ordre $p=4$ de coefficients aléatoires. Le taux d'insertion est de $R = 0.47$ et un bruit Gaussien stationnaire est ajouté avec $\bar{\sigma} = 5.43$. On constate que les résultats empiriques coïncident très bien avec les résultats théoriques annoncés dans la section 4.

La figure 2 présente une comparaison expérimentale des performances de plusieurs détecteurs sur 9000 images de 128×128 pixels en niveaux de gris codés sur 8 bits issues de la base BOSS[8]. Chacune des images a été analysée avant et après insertion d'un message binaire de 820 bits (soit $R \approx 0.05$) constitué de $\{0, 1\}$ uniformément distribués. De nombreux test peuvent potentiellement être choisis pour la comparaison. Le choix s'est porté sur le SPA/LSM qui est parmi les plus performants

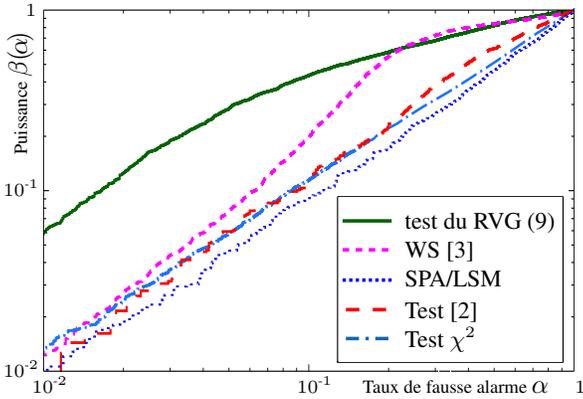


Figure 2: Puissance β des tests en fonction de la probabilité de fausse alarme α .

des détecteurs structuraux [9]. Le test du χ^2 et le test de [2] ont été choisis car ils font parties des rares détecteurs reposant sur des tests statistiques. Enfin, le WS [3] est un compétiteur naturel au regard de la similarité de son expression avec le test du RV (3). On constate que le test du RVG proposé a une puissance bien plus grande que tous les autres détecteurs pour des taux de fausses alarmes suffisamment faibles pour permettre une utilisation pratique, typiquement $\alpha = 10^{-2}$. Pour comprendre la perte de puissance, pour des faibles taux de fausses alarmes α , des test présentés sur la figure 2, la probabilité de distributions empiriques de $\hat{\Lambda}$ et du WS sont présentés sur la figure 3.

La figure 3 illustre l'importance du modèle du contenu des images. Le WS reposant sur un modèle autoregressif simple qui ne permet pas une approximation précise des images dont le contenu est complexe ; cela se traduit par des queues de distributions lourdes rendant difficile le respect d'une contrainte d'un faible taux de fausse alarme. A l'opposé, le modèle d'image proposé tient compte explicitement des discontinuités ce qui évite l'apparition de valeurs de $\hat{\Lambda}$ extrêmes.

6 Conclusion

Cet article est une première passerelle permettant l'application de la théorie de la décision statistique à la stéganalyse. Un premier test théorique est proposé pour le cas idéal où les paramètres

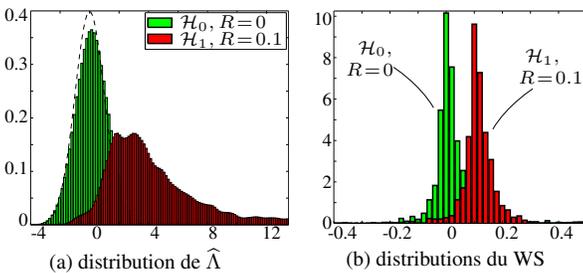


Figure 3: Distribution des statistiques de décisions sous \mathcal{H}_0 et \mathcal{H}_1 .

de distribution de tous les pixels sont connus. Les performances de ce test sont explicitées analytiquement ce qui permet de garantir le respect d'une contrainte de fausses alarmes.

Dans le cas pratique de l'analyse d'une image inconnue, un modèle local non-linéaire du contenu est proposé. Une méthode d'estimation simple et de bonne qualité du contenu est proposée par linéarisation du modèle. La perte de puissance due à l'estimation et à cette procédure de linéarisation est bornée. Cela permet de respecter une contrainte de fausses alarmes.

Les résultats numériques présentés sur des données simulées comme sur une base d'images naturelles montrent la pertinence de la méthodologie proposée. Le test pratique décrit dans cet article offre, pour des faibles taux de fausses alarmes, une puissance de détection bien supérieure aux tests proposés dans la littérature.

References

- [1] A. Foi, & al. "Practical poissonian-gaussian noise modeling and fitting for single-image raw-data," *Image Processing, IEEE Transactions on*, vol. 17, no. 10, pp. 1737–1754, Oct. 2008.
- [2] O. Dabeer, & al. "Detection of hiding in the least significant bit," *Signal Processing, IEEE Transactions on*, vol. 52, no. 10, pp. 3046 – 3058, oct. 2004.
- [3] J. Fridrich and M. Goljan, "On estimation of secret message length in LSB steganography in spatial domain," in *Security, Steganography, and Watermarking VI*, ser. Proc. SPIE, vol. 5306, 2004, pp. 23–34.
- [4] C. Bruni, A. De Santi, G. Koch, and C. Sinisgalli, "Identification of discontinuities in blurred noisy signals," *Circuit and systems-I, IEEE Trans. on*, vol. 44, no. 5, pp. 422 – 433, May 1997.
- [5] L. Fillatre, I. Nikiforov, and F. Reiraint, " ϵ -optimal non-bayesian anomaly detection for parametric tomography," *Image Processing, IEEE Transactions on*, vol. 17, no. 11, pp. 1985 –1999, nov. 2008.
- [6] R. Cogranne, C. Zitzmann, L. Fillatre, I. Nikiforov, F. Reiraint, and P. Cornu, "A cover image model for reliable steganalysis," in *Proc. of the 13th Information Hiding Conference*, 2011.
- [7] R. Cogranne, C. Zitzmann, L. Fillatre, F. Reiraint, I. Nikiforov, and P. Cornu, "Statistical decision by using quantized observations," in *Proc. of the IEEE International Symposium on Information Theory*, 2011.
- [8] BOSS contest, "Break our steganographic system," 2010. [Online]. Available: <http://boss.gipsa-lab.grenoble-inp.fr>
- [9] A. Ker, "A fusion of maximum likelihood and structural steganalysis," in *Proc. of 9th Information Hiding conference*, 2007, pp. 204–219.