

---

# Evaluation de tâches d'estimation sans gold standard

Irène Buvat

Imagerie et Modélisation en Neurobiologie et Cancérologie  
UMR 8165 CNRS - Universités Paris 7 et Paris 11

[buvat@imnc.in2p3.fr](mailto:buvat@imnc.in2p3.fr)

5 mars 2009

# Introduction

---

En présence de gold standard :

- biais
- variabilité
- combinaison des 2 critères (RMSE, etc...)
- corrélation

En l'absence de gold standard ?



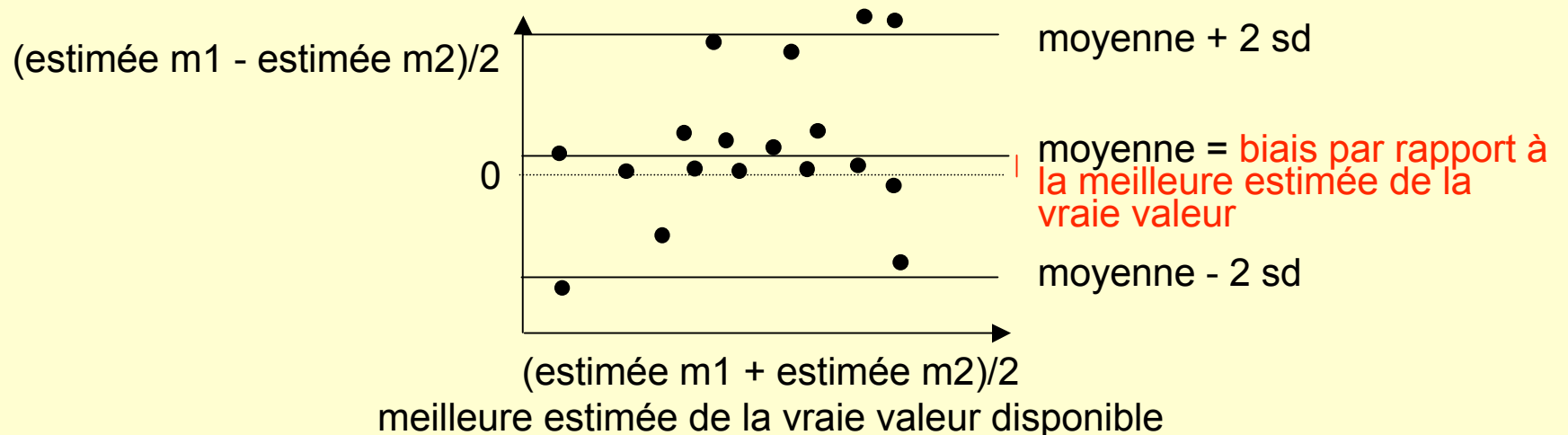
# Plan

---

- Approche Bland-Altman
- Régression dans gold standard
- Evaluation simultanée du gold standard et des performances de méthodes de segmentation d'images (STAPLE)

# Approche de Bland-Altman (1)

## Évaluation de l'accord entre deux méthodes



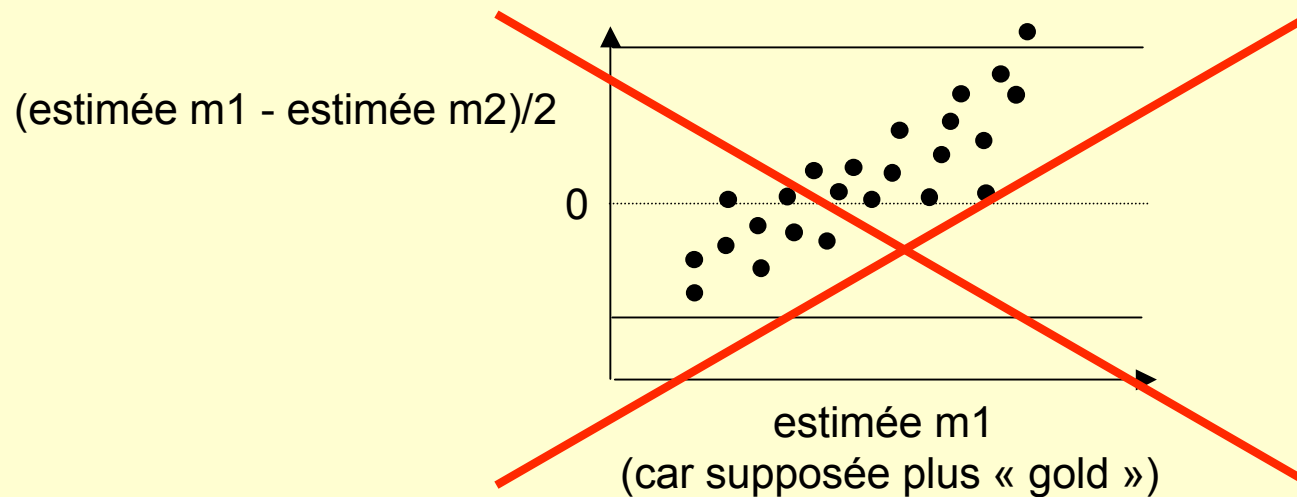
La plupart des différences sont comprises dans l'intervalle [moyenne - 2 sd ; moyenne + 2 sd]

L'étendue de cet intervalle doit permettre de conclure à l'interchangeabilité des méthodes ou non, **mais PAS au fait que l'une est moins biaisée que l'autre !**

*Bland and Altman. Lancet: 307-10, 1986.*

## Approche de Bland-Altman (2)

Permet de détecter des différences systématiques entre les méthodes

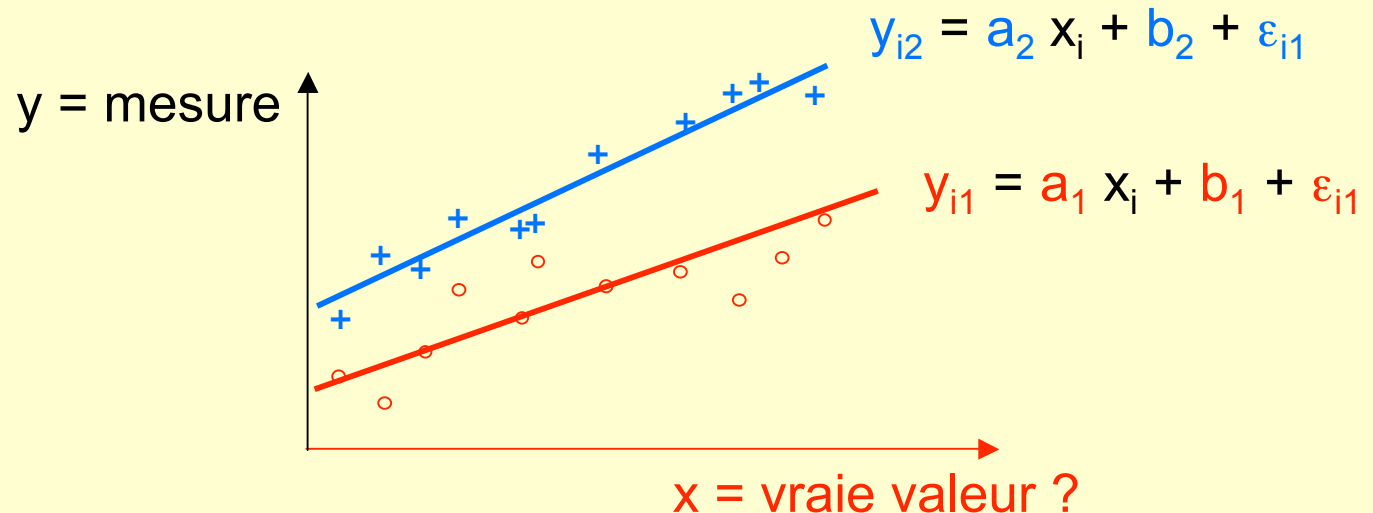


*Bland and Altman. Lancet, 346: 1085-7, 1995.*

# EWAGS (1)

Hypothèses:

1) Au moins 2 méthodes  $m$  d'estimation du même paramètre  $x_i$



2)  $y_{im} = a_m x_i + b_m + \epsilon_{mi}$  pour  $m = 1, 2, \dots$  et  $\epsilon_{mi} \sim \mathcal{N}(0, \sigma_m)$

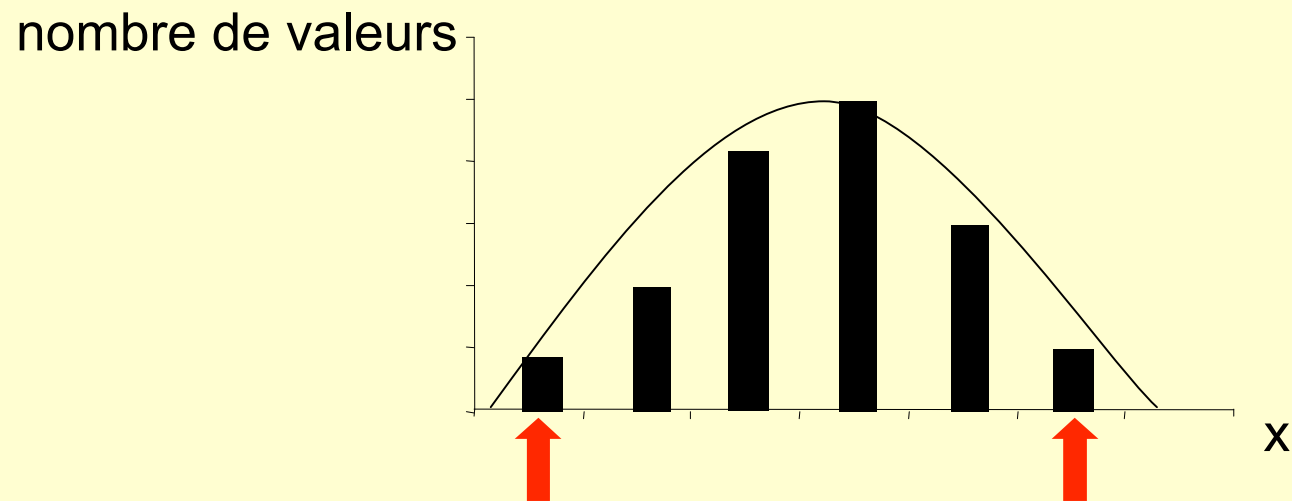
Hoppin, Kupinski, Kastis, Clarkson, Barrett. *IEEE Trans Med Imaging* 21: 441-449, 2002

Kupinski, Hoppin, Clarkson, Barrett, Kastis. *Acad Radiol* 9: 290-297, 2002

## EWAGS (2)

Hypothèses :

- 3)  $x_i$  suit une loi beta ou une distribution normale tronquée  $\theta$  définie par 2 paramètres inconnus  $\pi_1$  et  $\pi_2$ .  
Les valeurs **min** et **max** de la distribution sont approximativement connues



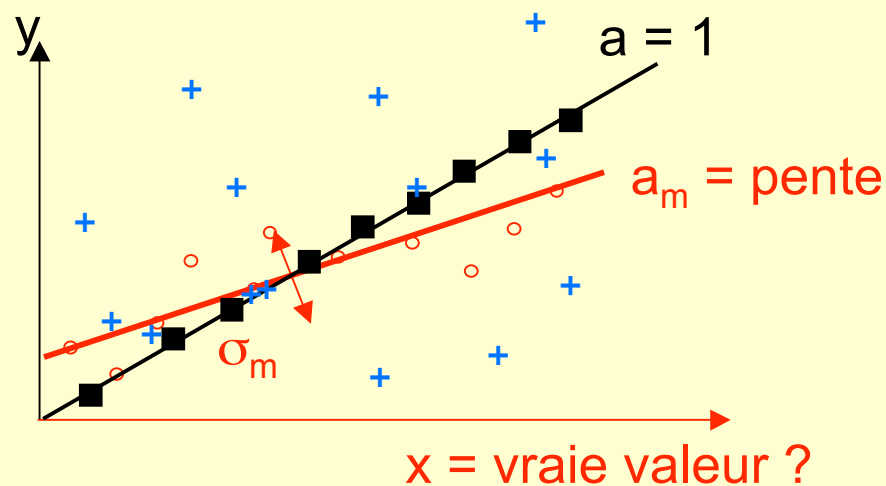
# EWAGS (3)

Méthodes :

1) Maximisation de la log-vraisemblance des paramètres du modèle

$$\mathcal{L}(\{a_m, b_m, \sigma_m\}_{\text{pour tt } m} \mid \{y_{im}\})$$

2) Calcul de  $\sigma_m/a_m$  comme figure de mérite



Estimée idéale :

- pas de dispersion :  $\sigma_m \sim 0$
- pente =  $a_m = 1$
- $\sigma_m/a_m$  tend vers 0

Estimée sans intérêt :

- grande dispersion :  $\sigma_m \gg 1$
- pente =  $a_m \sim 0$
- $\sigma_m/a_m$  tend vers  $\infty$

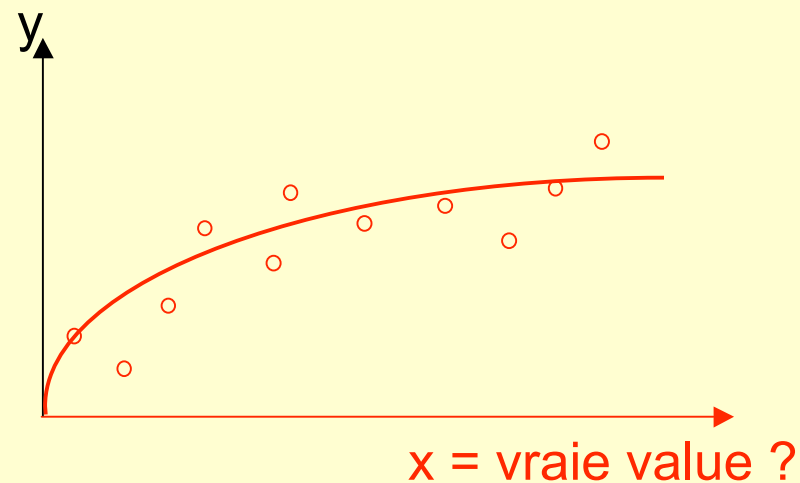


# Généralisation de EWAGS : GEWAGS (1)

Hypothèses :

1) Au moins 2 méthodes  $m$  d'estimation du même paramètre  $x_i$

2)  $y_{im} = a_m x_i^2 + b_m x_i + c_m + \varepsilon_{mi}$  pour  $m = 1, 2, \dots$  et  $\varepsilon_{mi} \sim \mathcal{N}(0, \sigma_m)$



*Buvat et al, J Nucl Med 48: 44P, 2007*

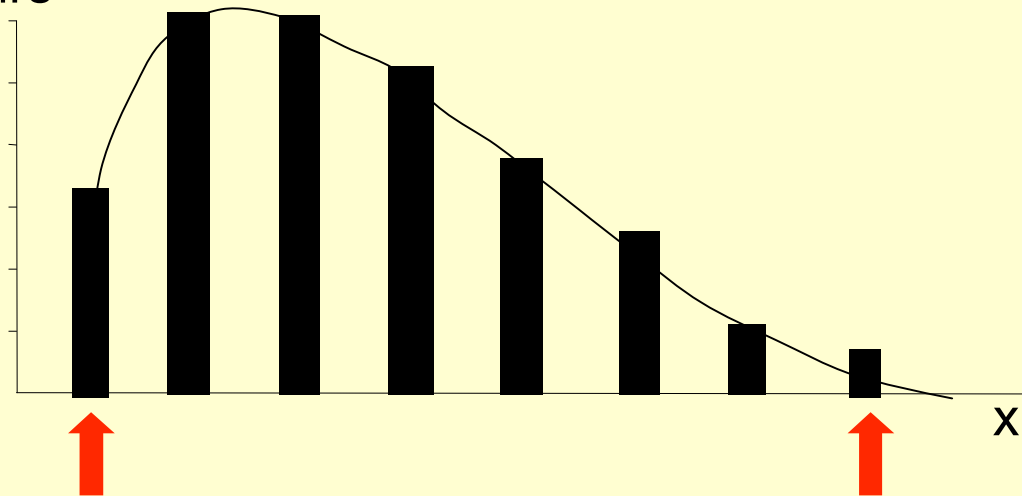
## Généralisation de EWAGS : GEWAGS (2)

Hypothèses :

- 3)  $x_i$  suit une distribution **beta**  $\theta$  définie par 2 paramètres inconnus  $\pi_1$  et  $\pi_2$

Les valeurs **min** et **max** de la distribution sont approximativement connues

nombre de valeurs



Pas de choix à faire entre loi beta et loi gaussienne  
Rend compte des distributions asymétriques

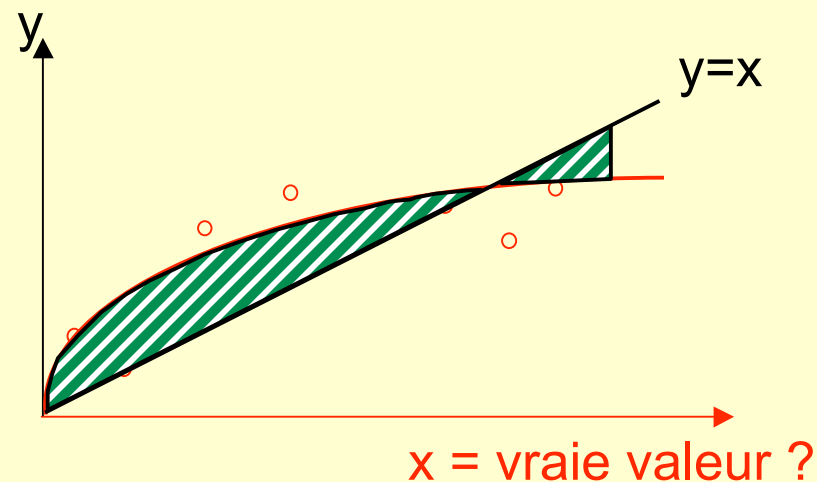
# Généralisation de EWAGS : GEWAGS (3)

Méthode :

1) Maximisation de la log-vraisemblance des paramètres du modèle

$$\mathcal{L}(\{a_m, b_m, c_m, \sigma_m\} | \{y_{im}\})$$

2) Calcul de  $sMSE_m = \sum_{i=1, X} (y_{im} - x_i)^2 / X$  comme figure de mérite



3) Calcul d'intervalles de confiance autour de  $sMSE_m$  au moyen d'une approche bootstrap non paramétrique

# Discussion



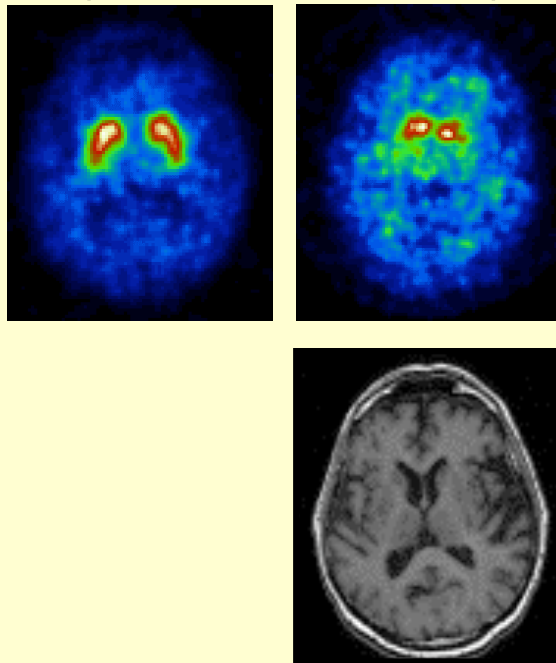
- Plus on dispose de méthodes d'estimation répondant au modèle, meilleur sera l'ajustement, mais 2 méthodes sont suffisantes
- Au moins 25 cas sont nécessaires
- La méthode est robuste même lorsque l'hypothèse sur la distribution des  $x_i$  est approximative

*Hoppin, Kupinski, Kastis, Clarkson, Barrett. IEEE Trans Med Imaging 21: 441-449, 2002*  
*Kupinski, Hoppin, Clarkson, Barrett, Kastis. Acad Radiol 9: 290-297, 2002*

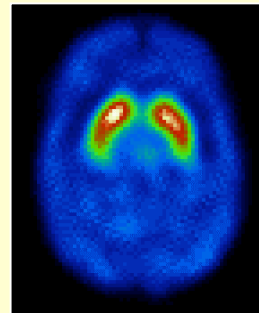
## Exemple d'application 1 - Données

Images 123I-FP-CIT SPECT corrigées de l'atténuation et de la diffusion

23 patients  
(13 MA + 10 DCL)



23 patients simulés\*



Cartographie de  
densité des tissus  
déduites des IRM des  
patients

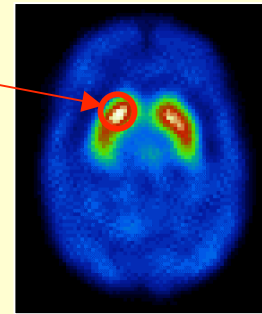
Cartographie d'activité  
déduites des images  
SPECT de patients

\*Soret, Koulibaly, Darcourt, Buvat. *Eur J Nucl Med Mol Imaging* 33:1062-1072, 2006

## Exemple d'application 1 - Méthodes comparées

Mesure du potentiel de liaison (PL) =  
(activité striata - activité fond) / activité fond

ds les 2 noyaux caudé et les 2 putamens



2 méthodes à comparer :

mesure sans correction de volume partiel (no CVP)

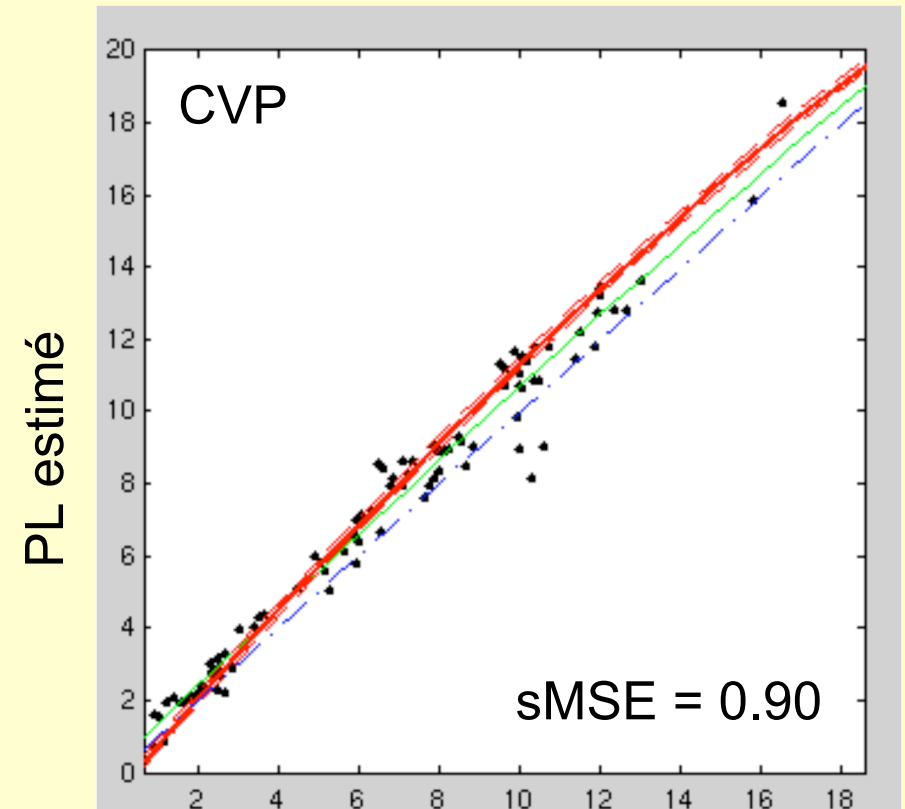
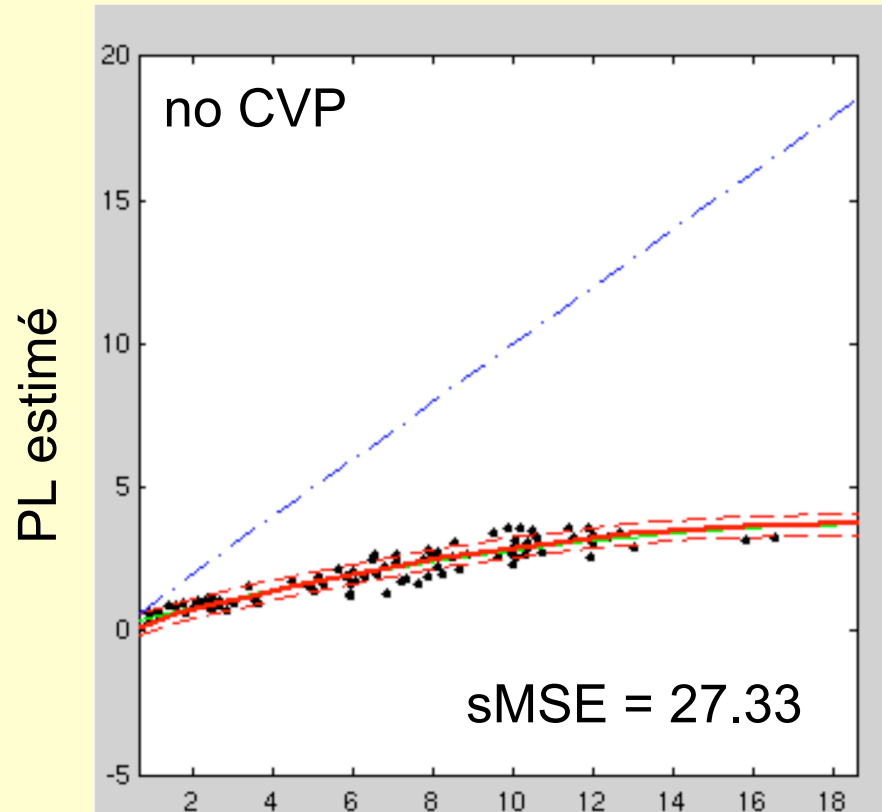
mesure avec correction de volume partiel \* (CVP)

Mise en œuvre de EWAGS et GEWAGS pour  
comparer les 2 méthodes

\*Soret, Koulibaly, Darcourt, Hapdey, Buvat. *J Nucl Med* 44: 1184-1193, 2003

# Exemple d'application 1 - Résultats de GEWAGS (1)

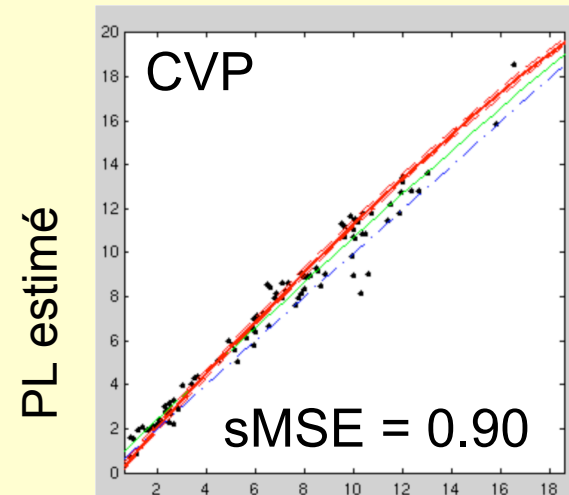
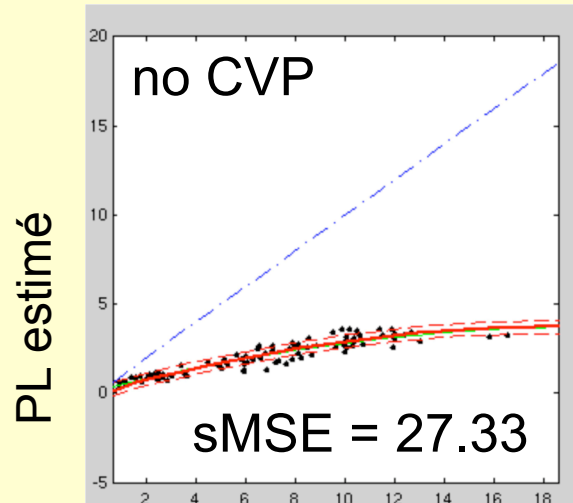
Patients simulés



Soret, Alaoui, Koulibaly, Darcourt, Buvat. Nucl Instrum Meth Phys Res A 571: 173-176, 2007

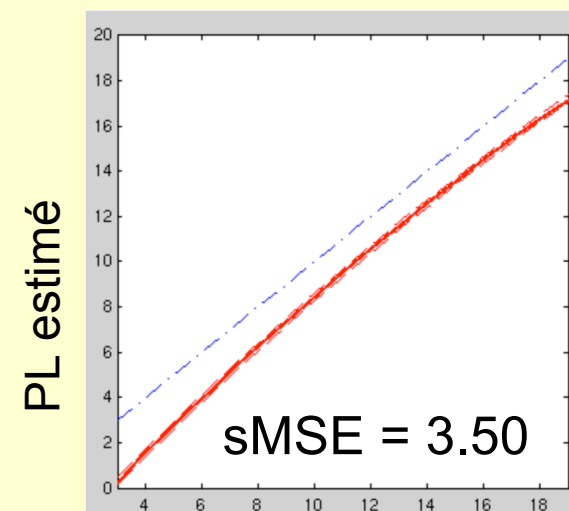
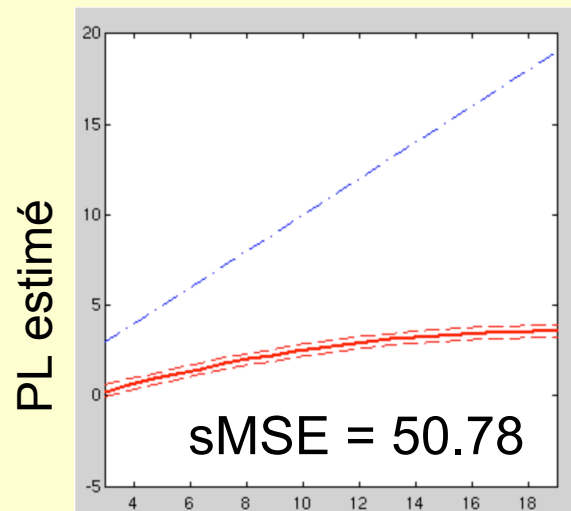
# Exemple d'application 1 - Résultats de GEWAGS (2)

## Patients simulés



$p < 0.01$

## Patients réels

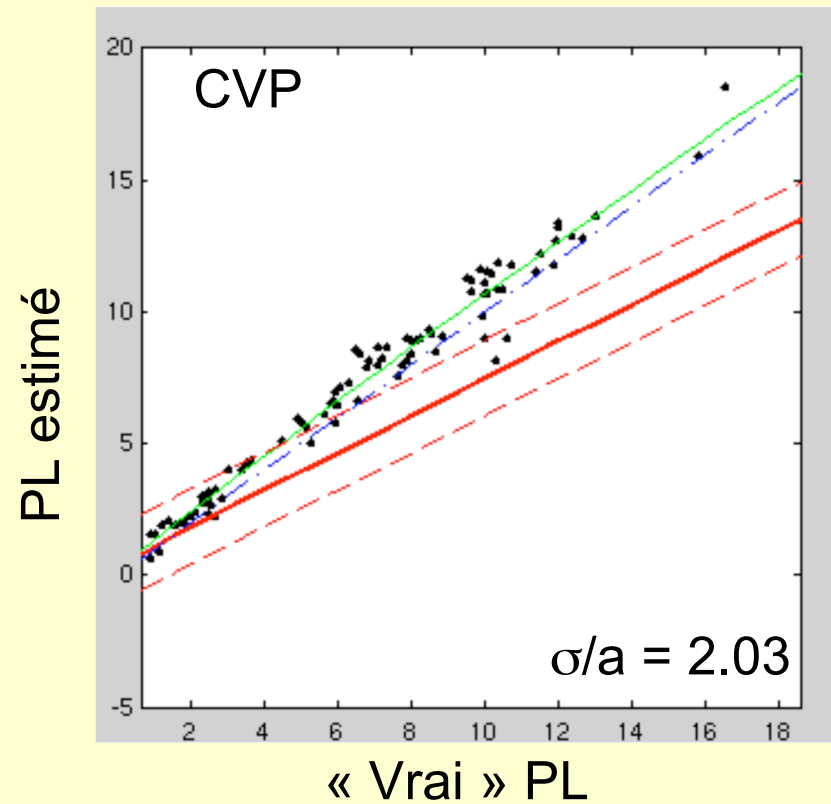
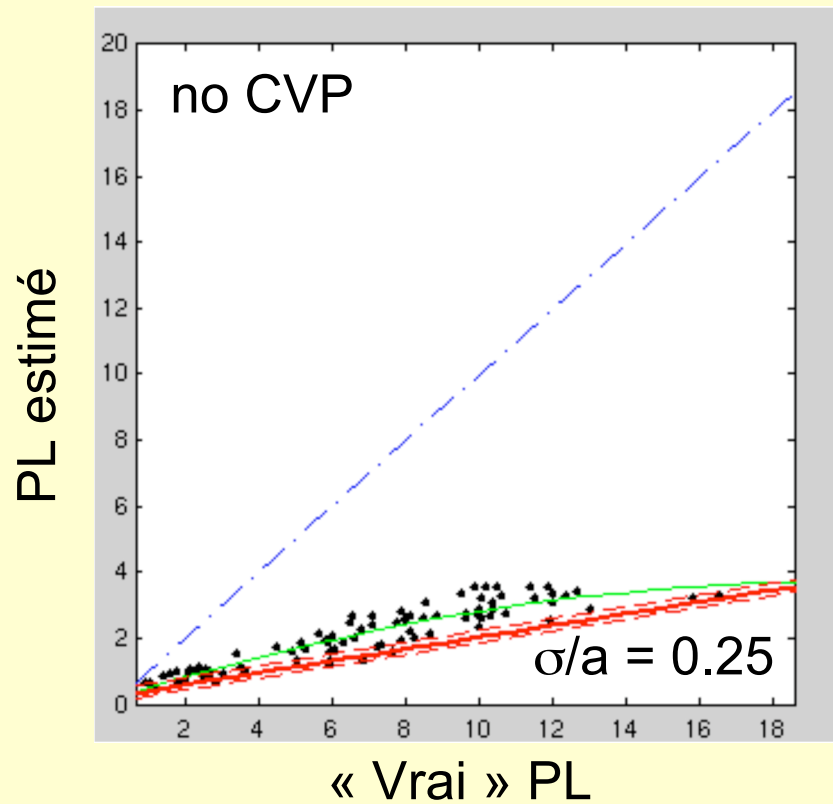


$p < 0.01$



# Exemple d'application 1 - Limites de EWAGS

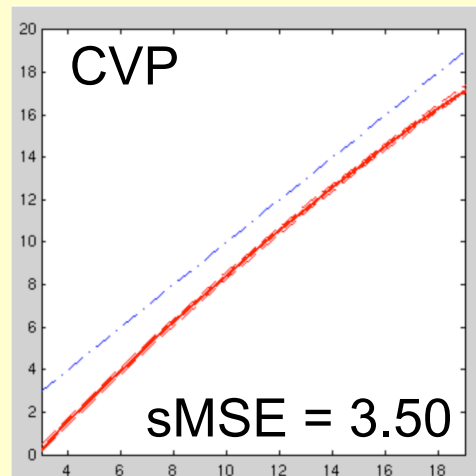
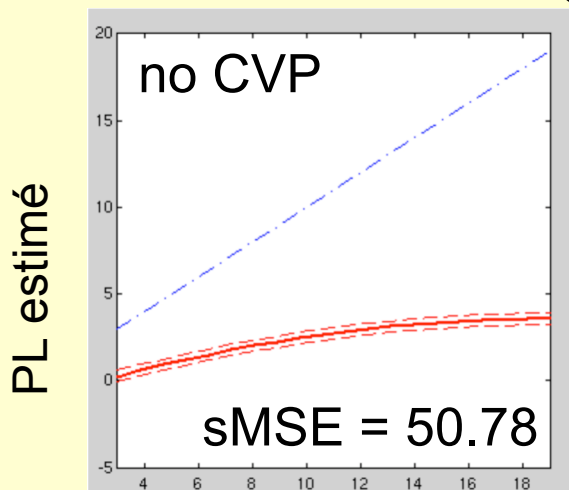
Patients simulés



# Exemple d'application 1 - Limites de EWAGS

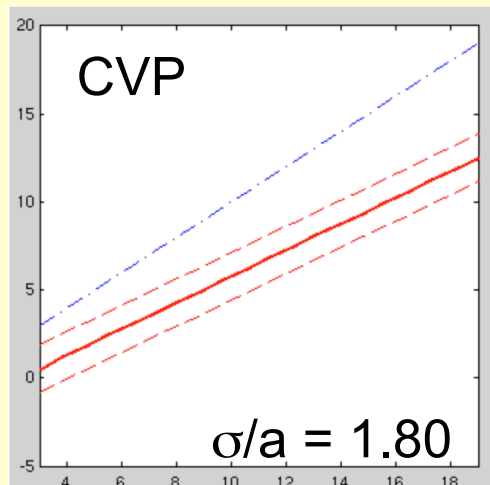
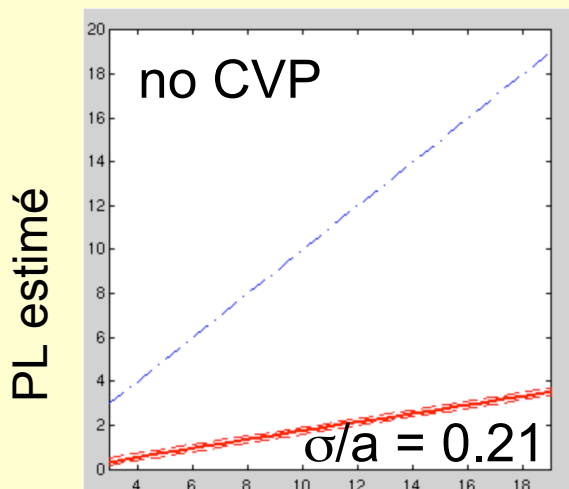
Patients réels

GEWAGS



$p < 0.01$

EWAGS



# Robustesse de GEWAGS

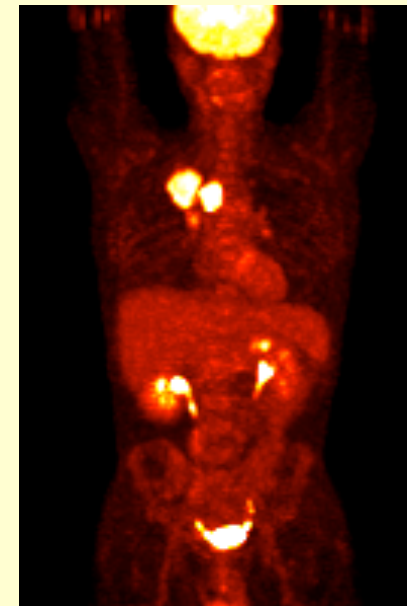
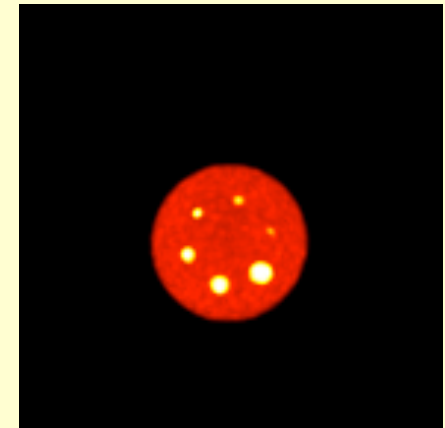
	GEWAGS sMSE		EWAGS $\sigma/a$	
	no CVP	CVP	no CVP	CVP
MA simulés	66.6	2.6*	0.22	2.80
DCL simulés	11.4	2.2*	0.17	0.08
MA+DCL simulés	27.3	0.9*	0.25	2.03
MA réels	35.4	2.9*	0.16	0.06
DCL réels	8.46	1.6*	0.17	0.07
MA+DCL réels	50.8	3.5*	0.21	1.80

\*  $p < 0.01$

## Exemple d'application 2 - Données

### Images TEP au FDG

- Fantôme :
  - Acquis sur le TEP/TDM Gemini GXL (HEGP)
  - 6 réplicats d'une acquisition de 6 sphères de volume 0,5 à 16mL
  - 1 SUV par sphère, variant de 4,3 à 6,07 (mesuré à l'activimètre)
- Patients :
  - Acquis sur le TEP/TDM Discovery LS (Institut J. Bordet, Bruxelles)
  - 14 examens PET/CT
  - 32 tumeurs pulmonaires + métastases



*\*Tylski, Dusart, Garcia, Vanderlinden, Buvat SNM 2009 (soumis)*

## Exemple d'application 2 - Méthodes comparées

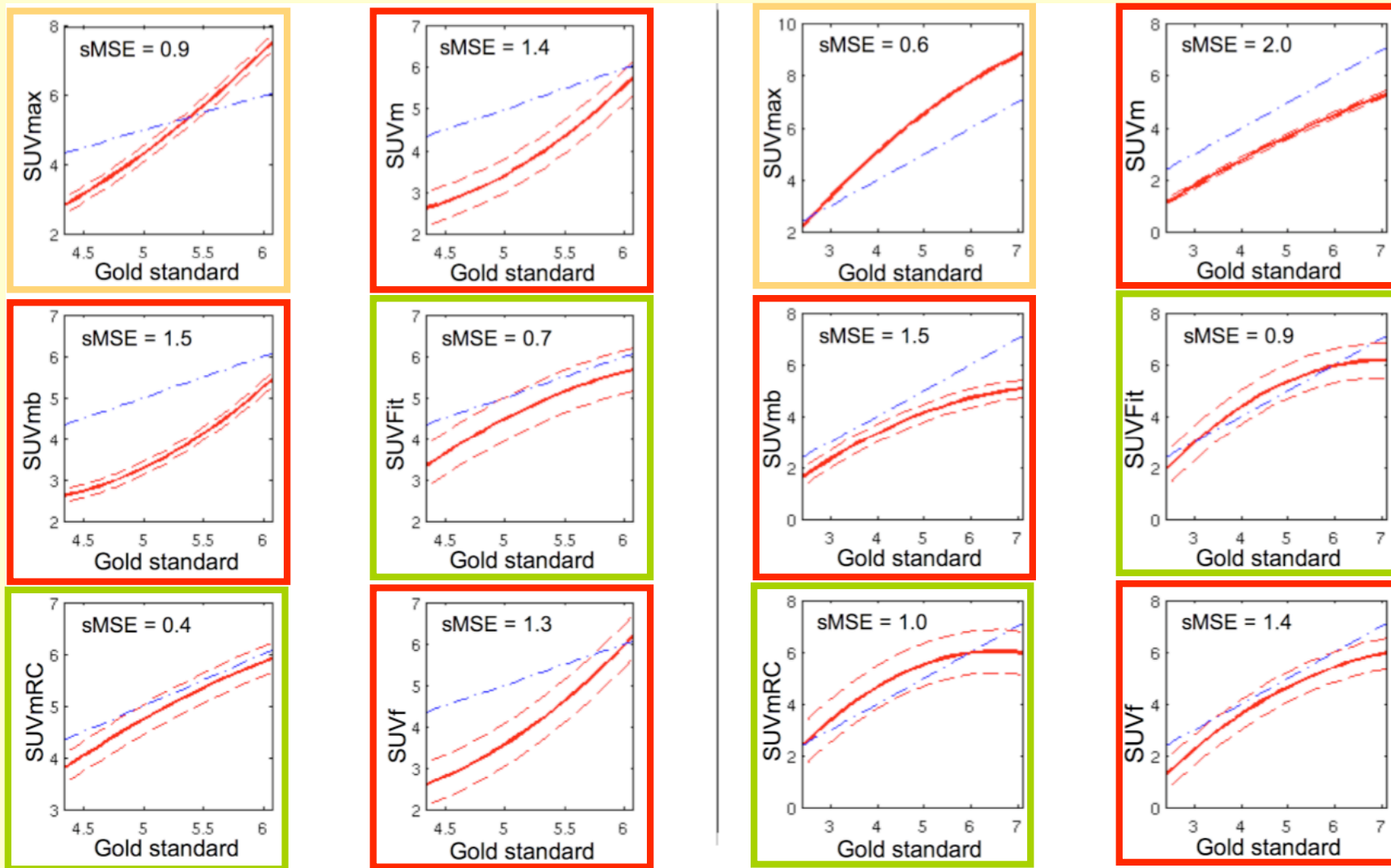
6 méthodes d'estimation de SUV :

- SUV maximum :  $SUV_{max}$
- SUV moyen dans une région fixe (cylindre de 15 mm de diamètre et de 15mm de hauteur) :  $SUV_f$
- SUV moyen dans une région définie par un seuil d'intensité :  $SUV_m$
- SUV moyen dans une région définie par un seuil d'intensité avec prise en compte de l'activité environnante :  $SUV_{mb}$
- SUV moyen corrigé de l'effet de volume partiel par une méthode d'ajustement  $SUV_{meanfit}$
- SUV moyen corrigé de l'effet de volume partiel par un coefficient de recouvrement  $SUV_{mRCt}$



Mise en œuvre de GEWAGS pour  
comparer les 6 méthodes

# Exemple d'application 2 - Résultats



Données fantômes

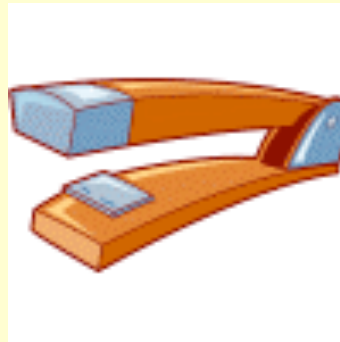
Données cliniques

# STAPLE

---

## Simultaneous Truth and Performance Level Estimation

Idée un peu similaire, mais dédiée à la problématique de segmentation de structure(s)

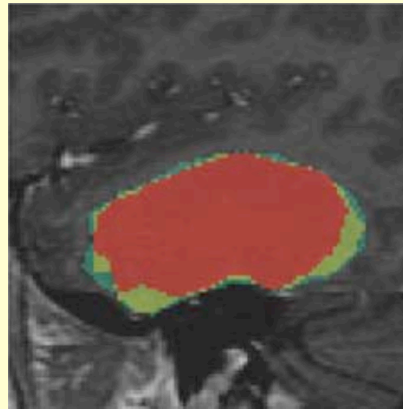


*Warfield, Zou, Wells, IEEE Trans Med Imaging 23: 903-921, 2004*

# STAPLE : hypothèses

Hypothèses :

- 1) Au moins 2 résultats de segmentation indépendants disponibles ( $m$ ) :  
 $m$  décisions de segmentation dans chacun des  $N$  voxels :  $D(N, m)$



*IRM prostate*

- 2) Segmentation vraie = variable binaire dans chaque voxel :  $T(N)$
- 3) Chaque méthode de segmentation peut être complètement caractérisée par sa sensibilité  $p(m)$  et spécificité  $q(m)$



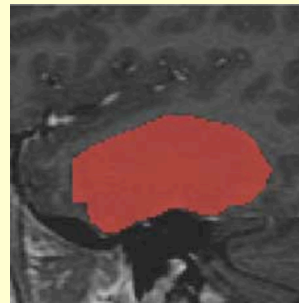
# STAPLE : méthode

Méthodes :

1) Initialisation des  $p_m$  et  $q_m$  à la même valeur qqsoit m OU initialisation de la segmentation vraie

2) Maximisation itérative (EM) de la log-vraisemblance :

$$\mathcal{L}(\{D, T\} | \{p_m, q_m\}_{\text{pour tt } m})$$



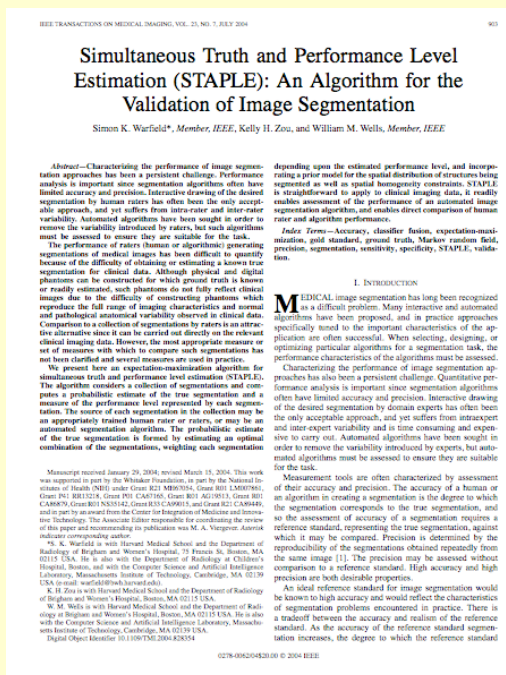
T

3) Caractérisation des performances de chaque méthode par  $p_m$  et  $q_m$

Expert	$\hat{p}$	$\hat{q}$	PV
1	0.8951	0.9999	0.998
2	0.9993	0.9857	0.977
3	0.9986	0.9982	0.954

# STAPLE : courte discussion

- Possibilité d'évaluer la segmentation de plusieurs structures
- Possibilité d'introduire des a priori (approche MAP)
- Possibilité d'introduire des contraintes spatiales
- Utilisé en IRM

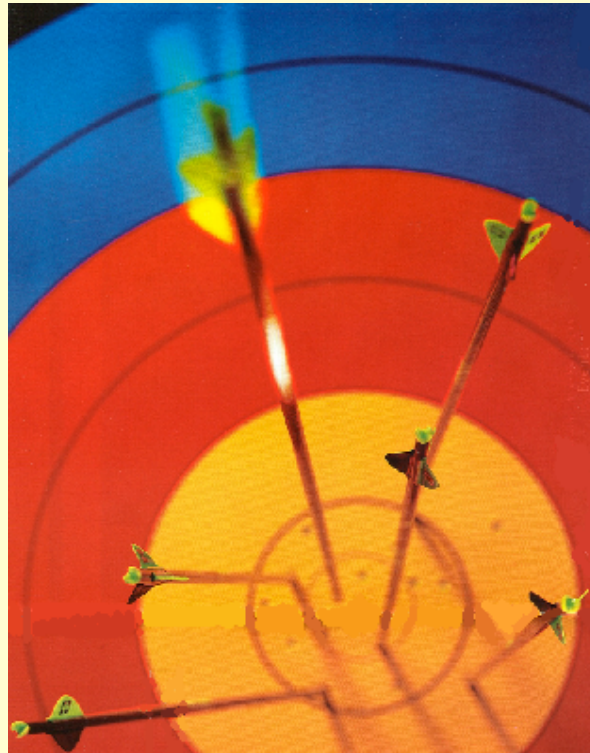


Warfield, Zou, Wells, IEEE Trans Med Imaging 23: 903-921, 2004

# Conclusion

---

- Une évaluation rigoureuse est possible même en l'absence de gold standard



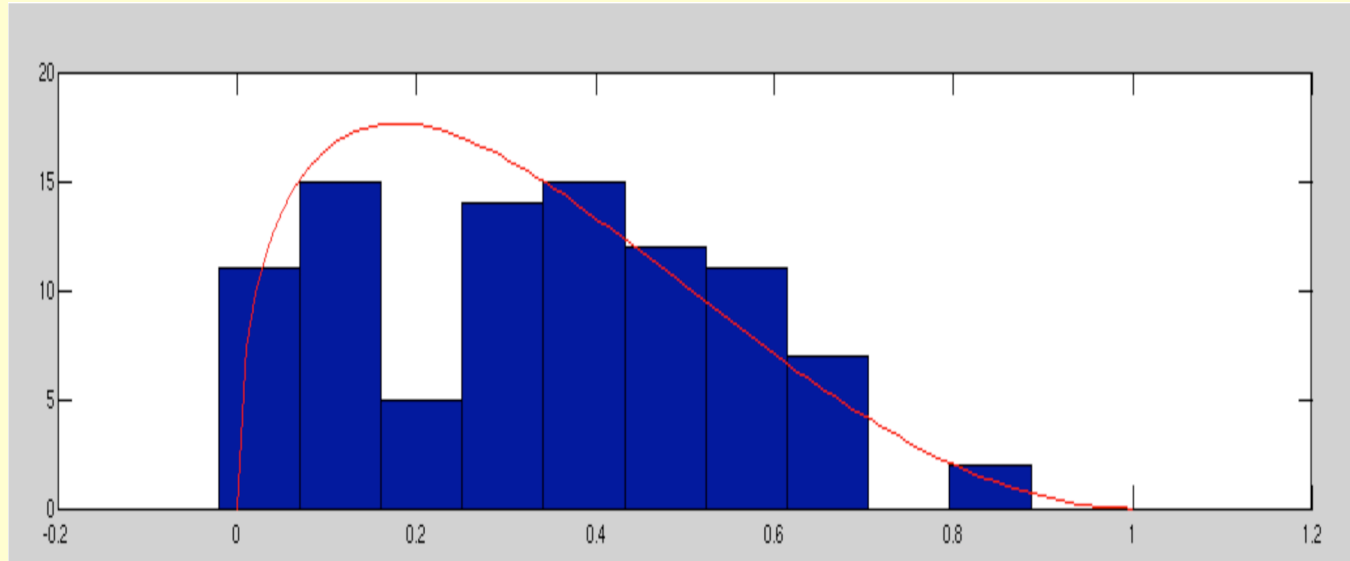
---

Diaporama disponible sur  
<http://www.guillemet.org/irene/conferences>

# Exemple de distribution beta

Patients simulés de l'étude MA / DCL

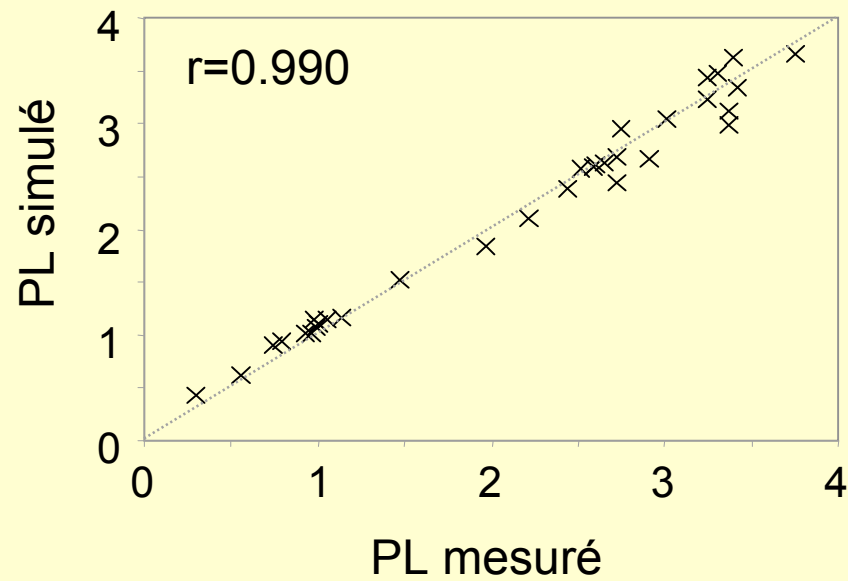
nombre de  
valeurs



PL normalisé

# Cohérence des données réelles et simulées (étude MA - DCL)

No CVP



CVP

