

THÈSE DE L'UNIVERSITÉ DE LYON

Délivrée par

UNIVERSITÉ CLAUDE BERNARD LYON 1

DIPLÔME DE DOCTORAT

Ma spécialité

DÉCODAGE DES INTENTIONS ET DES REPRÉSENTATIONS MOTRICES CHEZ L'HOMME : ANALYSE MULTI-ÉCHELLE ET APPLICATION AUX INTERFACES CERVEAU-MACHINE

par

Etienne Combrisson

Thèse soutenue le 09/2016 devant le jury composé de :

M ^{me}	ERIKA RATÉ	Université à la Menthe	(Rapporteur)
M.	JACQUES OUILLE	Université à la Fraise	(Rapporteur)
M.	HENRI ZOTO	Laboratoire laborieux	(Rapporteur)
M.	JEAN FILE	UTC	(Directeur)
	etc.		

*À Isabelle et Didier, mes deux parents,
qui ont tout donné pour que ceci me soit un jour possible.
Merci*

REMERCIEMENTS

Je voudrais tout d'abord exprimer mes plus profonds remerciements à...
AHÂÂÂH!

Je conclurai en remerciant de tout cœur (l'être aimé).

Montréal, le 28 juin 2016.

TABLE DES MATIÈRES

TABLE DES MATIÈRES	vi
LISTE DES FIGURES	viii
NOTATIONS	1
I Introduction générale	3
1 PRÉSENTATION DE LA THÉMATIQUE	7
1.1 ENREGISTREMENT DE L'ACTIVITÉ NEURONALE	7
1.1.1 Enregistrement non-invasif	7
1.1.2 Enregistrement invasif	7
1.2 ÉTAT DE L'ART DES INTERFACES CERVEAU-MACHINE	7
1.2.1 ICM non-invasives	8
1.2.2 ICM invasives	8
1.3 APPRENTISSAGE MACHINE : APPLICATIONS AUX NEUROSCIENCES	8
1.4 ENCODAGE ET DÉCODAGE MOTEUR : BASES PHYSIOLOGIQUES .	8
1.5 INTENTION ET EXÉCUTION	9
1.6 DELAYED TASK : PROTOCOLE EXPÉRIMENTAL	9
2 OBJECTIFS DE LA THÈSE	11
2.1 DÉCODAGE CÉRÉBRALE À PARTIR D'ACTIVITÉ INTRACRÂNIENNE	11
2.2 EXPLORATION ET AMÉLIORATION DES FEATURES	11
2.3 COMPARATIF DES CLASSIFIEURS	11
2.4 EXPLORATION DES RÉGIONS NON-MOTRICES	12
3 MÉTHODOLOGIE	13
3.1 EXTRACTION DES FEATURES	13
3.1.1 Pré-requis	13
3.1.2 Puissance spectrale	14
3.1.3 Phase	15
3.1.4 Phase-amplitude coupling	15
3.2 APPRENTISSAGE SUPERVISÉ	19
3.2.1 Labellisation et apprentissage	20
3.2.2 ClassifiEURS	20
3.2.3 Cross-validation	20
3.2.4 Sur-apprentissage	21
3.3 DU SINGLE ET MULTI-FEATURES	21
3.3.1 Single-feature	21
3.3.2 Multi-features	21

4 DONNÉES EXPÉRIMENTALES	23
4.1 DONNÉES "CENTER-OUT"	23
4.2 AUTRES DONNÉES	23
5 OUVERTURE	25
II Étude 1 : niveau de chance et évaluation statistique des résultats de classification par apprentissage supervisé	27
5.1 PRÉSENTATION DE L'ÉTUDE	31
5.1.1 Contexte	31
5.1.2 Problématique	31
5.1.3 Résultats majeurs	31
5.2 ARTICLE	31
5.3 COMPLÉMENTS D'ÉTUDE	43
III Étude 2 : encodage de l'intention et de l'exécution motrice	45
5.4 RÉSUMÉ DE L'ÉTUDE	49
5.5 ARTICLE	49
CONCLUSION	49
IV Étude 3 : décodage des directions de mouvement pendant et avant l'exécution de mouvement de membres supérieurs	51
5.6 RÉSUMÉ DE L'ÉTUDE	55
5.7 ARTICLE	55
CONCLUSION	55
V Étude 4 : optimisation des paramètres de la bande gamma	57
5.8 RÉSUMÉ DE L'ÉTUDE	61
5.9 ARTICLE	61
CONCLUSION	61
VI Étude 5 : décodage des émotions	63
5.10 RÉSUMÉ DE L'ÉTUDE	67
5.11 ARTICLE	67
CONCLUSION	67
CONCLUSION GÉNÉRALE	69
A ANNEXES	71
A.1 PREUVE DU THÉORÈME TRUC	73
BIBLIOGRAPHIE	75

LISTE DES FIGURES

1.1	Techniques d'enregistrement de l'activité cérébrale	7
1.2	Pipeline général d'un Interface Cerveau-Machine	8
1.3	Contrôle d'un bras robotisé	8
1.4	Comparatif	9
2.1	Mécanismes du couplage phase-amplitude	11
2.2	Localisation des aires sensorimotrices	12
3.1	Densité de probabilité d'une distribution d'amplitudes en fonction de tranches de phases	16
3.2	Comparatif de méthodes d'évaluation de couplage phase-amplitude	18
3.3	(A) Exemple de scalogramme aligné sur la phase du β , (B) Exemple de comodulogramme	18
3.4	Principe du Linear Discriminant Analysis	20
3.5	Principe du Support Vector Machine	20

NOTATIONS

Général

ICM Interface Cerveau-Machine
BCI Brain Computer Interface

Enregistrements

EEG Électroencéphalographie
MEG Magnétoencéphalographie
SUA Single Unit Activity
MUA Multi Unit Activity
SEEG Stéréoélectroencéphalographie
ECoG Électrocorticographie

Features

PAC Phase Amplitude Coupling

Classifeurs

LDA Linear Discriminant Analysis
SVM Support Vector Machine
RF Random Forest
KNN k-Nearest Neighbor
NB Naive Bayes

Première partie

Introduction générale

B_LABLABLABKIBLABLOU INTRO ...

L'objectif de cette thèse a été de ...

Totalité des méthodes explorées durant ma thèse sont présentes dans une toolbox python appelée brainpipe, libre d'accès et de droit.

PRÉSENTATION DE LA THÉMATIQUE

1.1 ENREGISTREMENT DE L'ACTIVITÉ NEURONALE

- Présentation de chacun des types de données
 - Résolution et RSB
 - Avantages // inconvénients
- (Waldert et al., 2009)

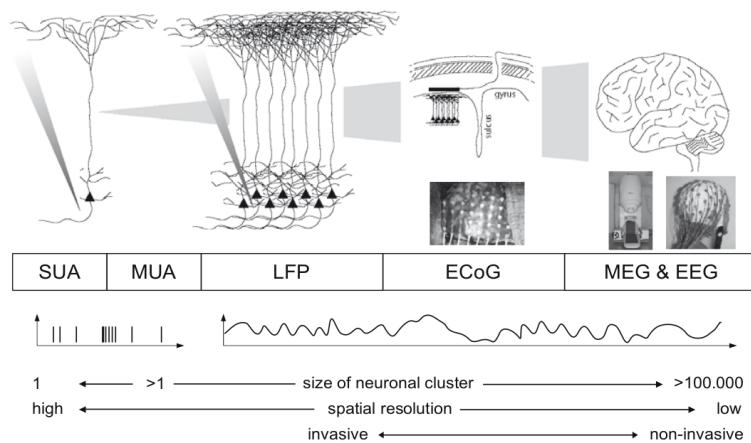


FIGURE 1.1 – Techniques d'enregistrement de l'activité cérébrale

1.1.1 Enregistrement non-invasif

Électroencéphalographie

Magnétoencéphalographie

1.1.2 Enregistrement invasif

Single Unit Activity

Multi Unit Activity

Stéréoélectroencéphalographie

Électrocorticographie

1.2 ÉTAT DE L'ART DES INTERFACES CERVEAU-MACHINE

(Bekaert et al., 2009)

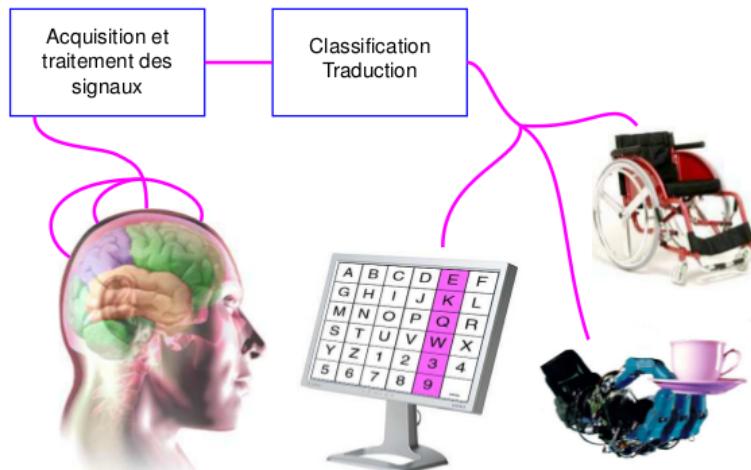


FIGURE 1.2 – Pipeline général d'un Interface Cerveau-Machine

1.2.1 ICM non-invasives

P300 speller

1.2.2 ICM invasives

(Hochberg et al., 2012)

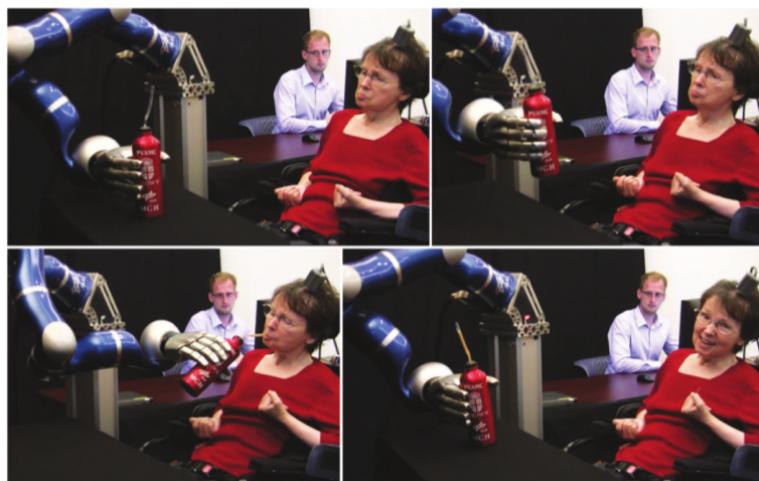


FIGURE 1.3 – Contrôle d'un bras robotisé

1.3 APPRENTISSAGE MACHINE : APPLICATIONS AUX NEUROSCIENCES

1.4 ENCODAGE ET DÉCODAGE MOTEUR : BASES PHYSIOLOGIQUES

La figure 1.4 représente une trajectoire d'un processus markovien de saut, avec les notations associées. (Hanakawa et al., 2008)

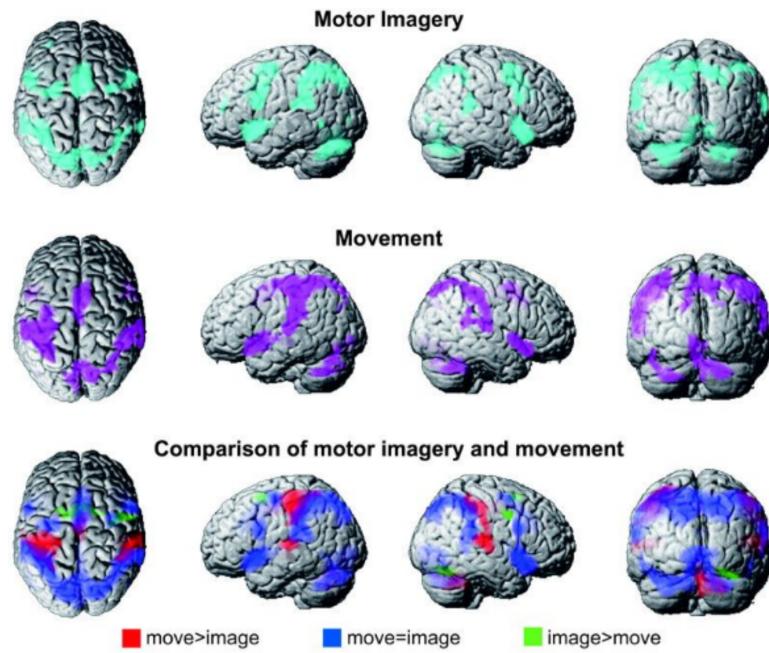


FIGURE 1.4 – Comparatif

1.5 INTENTION ET EXÉCUTION

1.6 DELAYED TASK : PROTOCOLE EXPÉRIMENTAL

OBJECTIFS DE LA THÈSE

2

2.1 DÉCODAGE CÉRÉBRALE À PARTIR D'ACTIVITÉ INTRACRÂ-NIENNE

- Exemple d'un schéma d'implantation + IRM
- Bipolarisation : débruitage et augmentation de la spécificité (article Karim)
- Extraction de features (ici on pourrait mentionner que le deep learning pourrait marcher sur les données brutes)

2.2 EXPLORATION ET AMÉLIORATION DES FEATURES

Rôle physiologique du phase-amplitude coupling

(Hyafil et al., 2015)

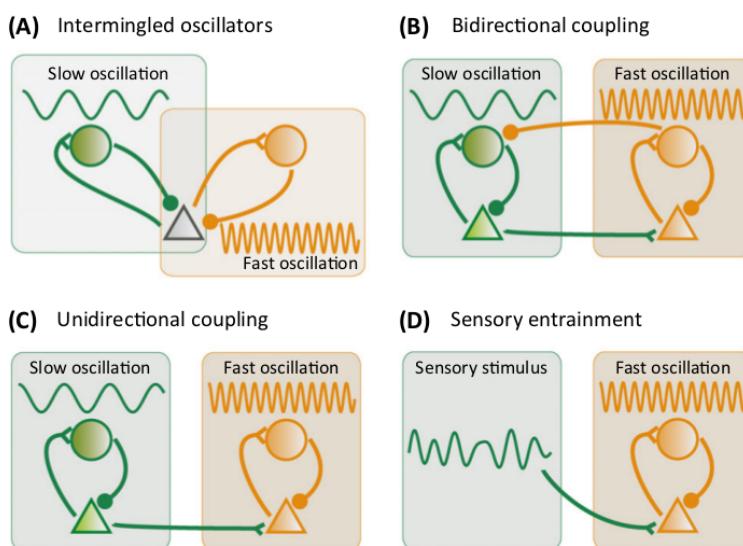


FIGURE 2.1 – Mécanismes du couplage phase-amplitude

2.3 COMPARATIF DES CLASSIFIERS

Expliquer que, chaque classifier possède une méthodologie propre permettant de répondre à des types de données différentes (en fonction des hypothèses de fonctionnement de chacun des classifiers)

2.4 EXPLORATION DES RÉGIONS NON-MOTRICES

(Van Langenhove et al., 2008)

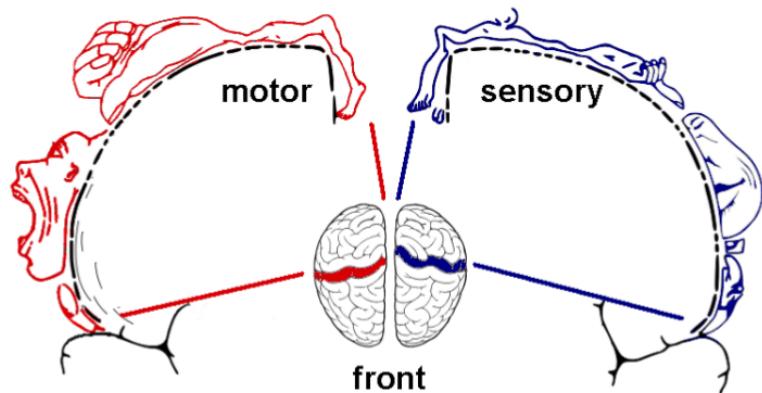


FIGURE 2.2 – Localisation des aires sensorimotrices

MÉTHODOLOGIE

3

Cette partie méthodologique sera divisée en deux grandes sous parties visant à présenter :

1. L'extraction des features : présentation des méthodes utilisées dans le cadre de l'extraction d'attributs issus de l'activité neuronale. De manière générale, nous avons étudiés des attributs spectraux comprenant :
 - Phase et puissance spectrale
 - Attributs de couplage
2. Le machine learning : présentation des principaux algorithmes testées dans le cadre du décodage de l'activité neuronale

3.1 EXTRACTION DES FEATURES

Comme nous l'avons décrit précédemment, l'objectif du décodage de l'activité neuronale est d'arriver à extraire des signaux cérébraux une information suffisamment pertinente pour pouvoir discriminer différents types de classes (exemple : mouvement vers la gauche Vs droite).

Tout les attributs testés dans le cadre de cette thèse sont des attributs spectraux, donc issus de bandes de fréquences. La plupart de ces outils partagent donc une partie méthodologique commune à savoir, le filtrage. De plus, la plupart sont extraits en utilisant la transformée d'Hilbert. Pour éviter une redondance à travers les attributs, nous allons tout d'abord introduire quelques pré-requis.

3.1.1 Pré-requis

Filtrage

L'intégralité des filtrages dans cette thèse ont été effectués avec la fonction *eegfilt* (qui a ensuite été reproduite pour le passage à python). De plus, afin d'éviter tout phénomène de déphasage, la fonction *filtfilt* a été systématiquement utilisée afin que le filtre soit appliqué dans les deux sens. Si cette dernière fonctionnalité n'est pas forcément indispensable dans le cadre d'un calcul de puissance, elle est absolument nécessaire pour un calcul de couplage phase-amplitude.

L'ordre du filtre présenté au dessus dépend de la fréquence de filtrage. Il a systématiquement été calculé en utilisant la méthode décrite par Bahramisharif et al. (2013) :

$$FiltOrder = N_{cycle} \times f_s / f_{oi} \quad (3.1)$$

où f_s est la fréquence d'échantillonnage, f_{oi} est la fréquence d'intérêt et N_{cycle} est un nombre de cycles définit par $N_{cycle} = 3$ pour les oscillations lentes et $N_{cycle} = 6$ pour les oscillations rapides.

Transformée d'Hilbert

Transformée permettant de passer un signal temporel $x(t)$ du domaine réel au domaine complexe. Le signal peut ensuite s'écrire $x_H(t) = a(t)e^{j\phi(t)}$ où $a(t)$ est l'amplitude et $\phi(t)$, la phase. Cette transformation est particulièrement exploitée car le module de $x_H(t)$ permet de récupérer l'amplitude et la phase est obtenue en prenant l'angle de $x_H(t)$.

Transformée en ondelettes

(Worrell et al., 2012, Tallon-Baudry et al., 1997)

$$f(a, b) = \int_{-\infty}^{\infty} f(x) \bar{\psi}_{a,b} dx \quad (3.2)$$

$$\psi_{a,b} = \frac{1}{\sqrt{a}} \Psi\left(\frac{x-b}{a}\right)$$

b facteur de translation et a de dilatation.

$$\psi_{a,b} = e^{-\pi x^2} e^{10i\pi x}$$

Checker dans mon implémentation la définition utilisée

A FAIRE

3.1.2 Puissance spectrale

Méthodes explorées

Le calcul de la puissance spectrale a été approché par deux méthodologies :

- La transformée d'Hilbert : souvent exploité dans le cadre du décodage ainsi que pour garder une uniformité entre les attributs de phase et couplage phase-amplitude basés eux aussi sur cette transformée.
- La transformée en ondelettes : principalement utilisée pour la visualisation des cartes temps-fréquence à cause de l'adaptation des ondelettes aux bandes physiologiques.

Normalisation

On utilise la normalisation pour observer l'émergence d'un phénomène par rapport à une période définie comme baseline. A travers la littérature, quatre grands types de normalisation sont rencontrés :

1. Soustraction par la moyenne de la baseline
2. Division par la moyenne de la baseline
3. Soustraction puis division par la moyenne de la baseline
4. Z-score : soustraction de la moyenne puis division par la déviation de la baseline

De manière générale, la normalisation z-score est la plus fréquemment rencontrée à travers la littérature large **PAPIER UTILISANT Z-SCORE**. Le choix du type de normalisation dépend du type de données utilisées. Dans le cadre de nos données, β était clairement la plus adaptée pour tout ce qui était de la visualisation.

En revanche, dans le cadre de la classification, nous obtenions systématiquement de meilleurs résultats sans normalisation.

Évaluation statistique

L'évaluation statistique de la puissance s'est fait comparativement à la baseline. Pour ce faire, nous avons testé deux approches :

1. Permutations : les essais de puissance et de baseline sont aléatoirement mélangés.

A FAIRE

1. Extraction de la puissance (Hibert, wavelet ou PSD)
2. Rôle physiologique des bandes de puissance
3. décodage puissance
4. Normalisation
5. Évaluation statistique + one/two tails
 - (a) Wilcoxon // Kruskal-Wallis
 - (b) Permutations

3.1.3 Phase

- Extraction de la phase (Hilbert)
- Rôle physiologique supposé
- décodage phase
- Evaluation statistique par stat de Rayleigh

3.1.4 Phase-amplitude coupling

Méthodologie du phase-amplitude coupling

Il existe une large variété de méthodes pour calculer le PAC, ce qui complique son exploration. Toutefois, il n'existe pas de consensus sur une méthode plus polyvalente qu'une autre, chacune possédant ses points forts et limitations. Pour aller un peu plus loin, et présenter quelques méthodes, il est nécessaire d'introduire quelques variables. Soit $x(t)$, une série temporelle de données de taille N. Pour cette série temporelle, on souhaite savoir si la phase extraite dans une bande de fréquence $f_\phi = [f_{\phi_1}, f_{\phi_2}]$ est couplée avec l'amplitude contenue dans $f_A = [f_{A_1}, f_{A_2}]$. Pour cela, on va tout d'abord extraire $x_\phi(t)$ et $x_A(t)$ les signaux filtrés dans ces deux bandes. Enfin, la phase $\phi(t)$ est obtenue en prenant l'angle de la transformée d'Hilbert de $x_\phi(t)$ tandis que l'amplitude $a(t)$ est obtenue en prenant le module de la transformée d'Hilbert de $x_A(t)$.

1. Mean Vector Length-Modulation Index

$$MVL = \left| \sum_{j=1}^N a(j) \times e^{j\phi(j)} \right| \quad (3.3)$$

où $a_h(t)$ est ... et $\phi(t)$ est ...

2. Kullback-Leibler divergence :

A l'origine, la divergence de Kullback-Leibler (KLD), qui est issue de la théorie de l'information, permet de mesurer les dissimilarités entre deux distributions de probabilités. Ainsi, pour pouvoir utiliser cette mesure dans le cadre du PAC, Tort et al. (2010) propose une solution élégante qui consiste à générer une distribution de densité probabilités de l'amplitude (DPA) en fonction des valeurs de phase et d'ensuite utiliser le KLD pour comparer cette distribution à la densité de probabilité d'une distribution uniforme (DPU). Plus la DPA s'éloigne de la DPU, plus le couplage entre l'amplitude et la phase est consistant.

Pour construire la DPA, l'astuce consiste à couper le cercle trigonométrique en N tranches (dans l'article il est proposé de couper en 18 tranches de 20°). Puis, si on prend l'exemple de la tranche $[0, 20^\circ]$, on va chercher tout les instants temporels où la phase prend des valeurs comprises entre $[0, 20^\circ]$ ($t, \phi(t) \in [0, 20^\circ]$). On prend ensuite la moyenne de l'amplitude pour ces valeurs de t et on répète cette procédure pour chacune des tranches de phase. On obtient ainsi la densité d'amplitudes en fonction des valeurs de phase. Il ne reste plus qu'à normaliser cette distribution par la somme des amplitudes à travers les tranches et on récupère une distribution de densité de probabilités. La figure 3.1 (Tort et al., 2010) présente un exemple de DPA en fonction de tranches de phase.

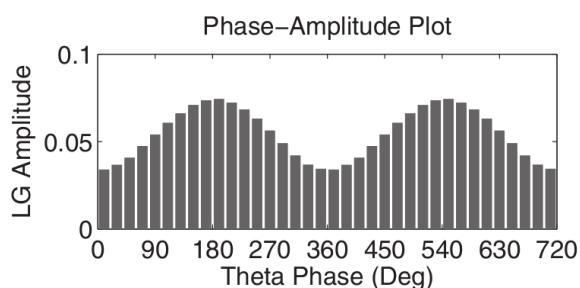


FIGURE 3.1 – Densité de probabilité d'une distribution d'amplitudes en fonction de tranches de phases

Le calcul de la divergence de Kullback-Leibler est ensuite appliqué pour mesurer les dissimilarités entre la DPA et la DPU et c'est cette mesure qui servira d'estimation du couplage phase-amplitude :

$$D_{KL}(P, Q) = \sum_{j=1}^N P(j) \times \log \frac{P(j)}{Q(j)} \quad (3.4)$$

où $P(j)$ est la densité de probabilité de $a(t)$ en fonction de $\phi(t)$ et

$Q(j)$ est la densité de probabilité d'une distribution uniforme.

3. Height Ratio

La méthode du Height Ratio(Lakatos, 2005) est extrêmement proche du Kullback-Leibler divergence. En effet, l'amplitude sera binée de la même façon en fonction des tranches de phase. La mesure du PAC est ensuite donnée par :

$$hr = (f_{max} - f_{min}) / f_{max} \quad (3.5)$$

où f_{max} et f_{min} sont respectivement le maximum et le minimum de la densité de probabilité de l'amplitude en fonction des valeurs de phase.

4. Normalized Direct Phase-Amplitude Coupling

Le Normalized Direct Phase-Amplitude Coupling, qui n'est pas une des méthodes les plus fréquemment rencontrées, présente toutefois une avantage certain. En plus de fournir une estimation fiable du couplage phase-amplitude, Ozkurt (2012) démontre l'existence d'un seuil à partir duquel on peut considérer l'estimation du PAC comme étant statistiquement fiable. La beauté de cette méthode, c'est que ce seuil statistique, qui est une fonction de la valeur p désirée, ne dépend que de la taille de la série temporelle. Ce qui rend son utilisation particulièrement simple.

Pour estimer le PAC, une des hypothèses ayant permis d'aboutir à ce seuil statistique est de devoir normaliser l'amplitude par un z-score dénotée $\tilde{a}(t)$. L'estimation du PAC est quasiment identique au MVL puisque c'est en réalité le carré de celle-ci. Enfin, pour une valeur p désirée, l'article introduit le seuil statistique :

$$x_{lim} = N \times [erf^{-1}(1 - p)]^2 \quad (3.6)$$

où erf^{-1} est la fonction d'erreur inverse. On déduira que l'estimation PAC est significative si et seulement si cette valeur est deux fois supérieure à ce seuil.

Robustesse du phase-amplitude coupling et évaluation statistique

- Calcul de surrogates
- Normalisation du pac par surrogates
- Calcul de la p-value

Comparatif des méthodes

A FAIRE

(Tort et al., 2010)

TABLE 1. Summary of characteristics of the phase-amplitude coupling measures studied

Phase-Amplitude Coupling Measure	Tolerance to Noise	Amplitude Independent	Sensitivity to Multimodality	Sensitivity to Modulation Width
Modulation index	Good	Yes	Good	Good
Heights ratio	Good	Yes	No discrimination	No
Mean vector length	Good	No	Restricted	Reasonable
Amplitude PSD	Low	No	Restricted	Good
Phase-locking value	Low	No*	Restricted	Low
Correlation measure	Low	No*	Restricted	Low
GLM measure	Low	No*	Restricted	Low
Coherence value	Low	No*	Restricted	Low

* Under the presence of noise.

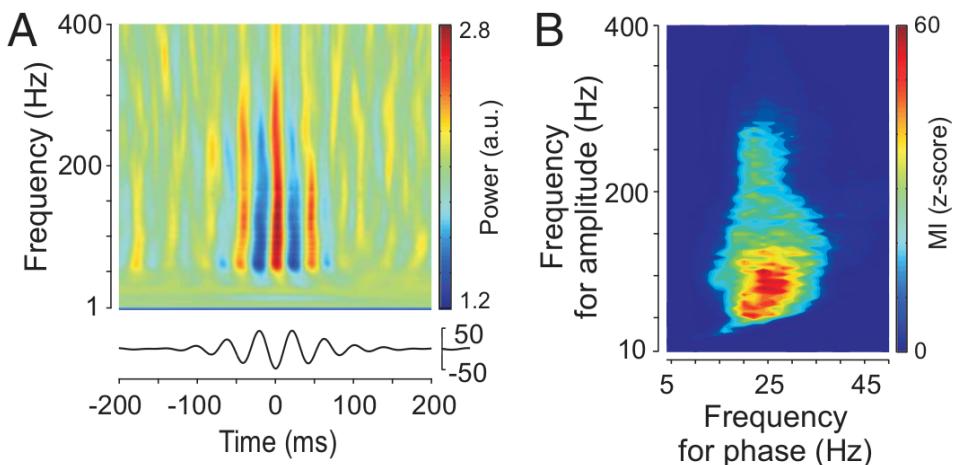
FIGURE 3.2 – Comparatif de méthodes d'évaluation de couplage phase-amplitude

Représentation du phase-amplitude coupling

Comparée à la puissance, l'exploration du PAC peut s'avérer plus complexe dû à sa dimensionnalité plus grande. Il existe donc des outils et des méthodes destinées à simplifier cette exploration et à visualiser ces résultats.

Exemple concret, si on cherche à connaître les modulations de puissance contenue dans un signal, on peut représenter une carte temps-fréquence. Pour le PAC, idéalement on voudrait visualiser les phases, les amplitudes et le temps mais ces trois dimensions empêche une représentation simple. On peut donc avoir recours à différents types de représentations complémentaires :

- Scalogramme : cette représentation permet de faire émerger l'existence d'un couplage, pour une phase donnée, et d'observer sa durée. Pour cela, on aligne les phase en détectant le pic le plus proche de l'instant temporel étudié. On calcul les cartes temps-fréquence que l'on va ensuite moyenner après les avoir recalées de la même façon que les phases (c'est-à-dire avec la même latence).
- Comodulogramme : pour une tranche temporelle définie, on représente les valeur de PAC pour différentes valeurs de phase et d'amplitude large [Ref comodulogramme ?](#)

FIGURE 3.3 – (A) Exemple de scalogramme aligné sur la phase du β , (B) Exemple de comodulogramme

La figure 3.3 (de Hemtinne et al., 2013) met en évidence que le sca-

rogramme (**A**) est limité d'une part, par la phase sur laquelle on choisit de recaler et d'autre part cette méthode est également limité par l'instant où l'on choisit de recaler. Pour la figure (**B**), le calcul du PAC se faisant à travers la dimension temporelle, on a aucune idée de l'évolution du couplage dans le temps.

Phase-amplitude coupling : résolution temporel ?

Comment peut-on savoir si un ensemble de musiciens jouent ensemble, en rythme ? L'approche traditionnelle consiste à dire que, en fonction de la prestation du groupe, on sera en mesure de dire si ils étaient en rythme ou non. Donc on focalise notre attention sur chaque instant du morceau et on analyse chaque note, chaque décalage. Cela signifie aussi que toute notre attention a été mobilisée par l'analyse du rythme et finalement, on passe à côté de la musique. Notre attention au détail nous a écarté du morceau global. On pourrait dire que l'on a écrasé la dimension temporelle du morceau. Une autre approche consiste à assister à toute les répétitions du fameux groupe. Ce faisant, on est capable de dire si d'une manière générale les musiciens ont tendance à jouer ensemble. Ainsi, le jour d'une représentation, toute notre attention peut rester uniquement sur le concert. On garde donc la dimension temporelle.

C'est par ce changement de positionnement face au problème de résolution temporelle que Voytek et al. (2013) introduit le Event Related Phase-Amplitude Coupling. L'approche traditionnelle du PAC nécessitant de connaître un nombre de cycles afin d'en déduire l'existence ou non du couplage, et donc perdre la dimension temps, l'article propose de calculer le PAC à travers les essais (ou répétitions). Pour un jeu de données de M essais de longueur N , on extrait respectivement les phases et les amplitudes $\phi_M(t)$ et $a_M(t)$ puis, pour chaque point temporel, on calcule la corrélation à travers les essais (corrélation linéaire-circulaire (Berens and others, 2009) qui se fait entre l'amplitude et des sinus/cosinus de la phase). Il en résulte une valeur de corrélation pour chaque instant et donc, de couplage.

Importance du filtrage

- Filtrage dans les deux sens pour éviter les shifts fréquentiels
- Ordre du filtre en fonction du nombre de cycles (Bahramisharif et al., 2013)
- Nombre de cycle : au moins un cycle pour pouvoir estimer le PAC mais comme c'est sensible au bruit, plus on multiplie le nombre de cycles, plus l'estimation est fiable (Tort et al., 2010, Voytek et al., 2013) mais en plus précis : (Bahramisharif et al., 2013)

3.2 APPRENTISSAGE SUPERVISÉ

Présentation du concept Training set et Testing set

3.2.1 Labellisation et apprentissage

3.2.2 Classifieurs

1. Linear Discriminant Analysis

(Fisher, 1936), (Lotte et al., 2007)

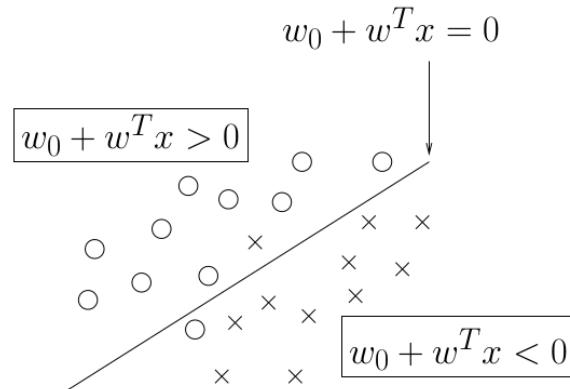


FIGURE 3.4 – Principe du Linear Discriminant Analysis

2. Support Vector Machine

(Cortes and Vapnik, 1995, Vapnik, 2000) (Lotte et al., 2007)

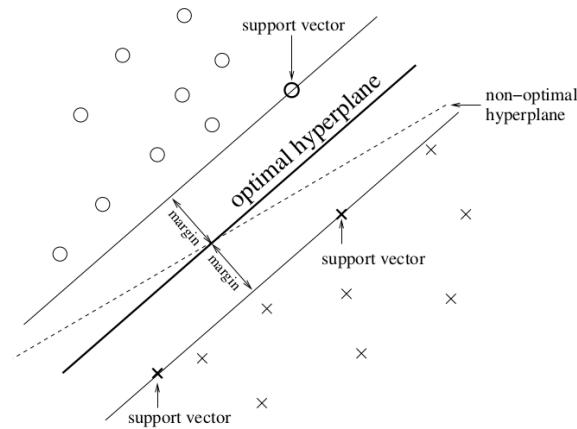


FIGURE 3.5 – Principe du Support Vector Machine

3. k-Nearest Neighbor
4. Naive Bayes
5. Random Forest

3.2.3 Cross-validation

Présentation et utilisation (contexte : séparation Training et Testing set // optimisation des paramètres (classifieurs et multi-features))

1. k-Fold
 - k-Fold, Stratified k-Fold, Shuffle split, et Stratified shuffle split
2. Leave-One-Out

Évaluation de la performance de décodage

1. Decoding Accuracy
2. Receiver operating characteristic

Évaluation statistique de la performance de décodage

1. Loie binomiale
2. Permutation : data driven + différents types de permutations (Ojala and Garriga, 2010)
 - Shuffle y
 - Full-shuffle
 - Intra-class shuffle y

3.2.4 Sur-apprentissage

Optimisation des paramètres de classification

3.3 DU SINGLE ET MULTI-FEATURES

Présentation du concept

3.3.1 Single-feature

3.3.2 Multi-features

1. Sélection statistique
 - (a) Sélection binomiale
 - (b) sélection permutations
2. Sélection séquentielle
 - (a) Forward selection
 - (b) Backward selection
 - (c) exhaustive selection

DONNÉES EXPÉRIMENTALES

4

4.1 DONNÉES "CENTER-OUT"

4.2 AUTRES DONNÉES

OUVERTURE

5

Nos contributions portent sur : ...

Le *premier chapitre* expose la problématique de la thèse.
Le *deuxième chapitre* présente en détail ...

etc.

Cette thèse a fait l'objet de divers travaux écrits : ...

Deuxième partie

Étude 1 : niveau de chance et évaluation statistique des résultats de classification par apprentissage supervisé

SOMMAIRE

5.1	PRÉSENTATION DE L'ÉTUDE	31
5.1.1	Contexte	31
5.1.2	Problématique	31
5.1.3	Résultats majeurs	31
5.2	ARTICLE	31
5.3	COMPLÉMENTS D'ÉTUDE	43
5.4	RÉSUMÉ DE L'ÉTUDE	49
5.5	ARTICLE	49
CONCLUSION		49
5.6	RÉSUMÉ DE L'ÉTUDE	55
5.7	ARTICLE	55
CONCLUSION		55
5.8	RÉSUMÉ DE L'ÉTUDE	61
5.9	ARTICLE	61
CONCLUSION		61
5.10	RÉSUMÉ DE L'ÉTUDE	67
5.11	ARTICLE	67
CONCLUSION		67

Sensibilisation à l'importance du nombre d'essais par exemple

5.1 PRÉSENTATION DE L'ÉTUDE**5.1.1 Contexte****5.1.2 Problématique****5.1.3 Résultats majeurs**

pourquoi cette étude ? Quelles questions ? - Seuil de chance théorique vs pratique ? - Impact sur des méthodes (cross-validation, classifieur) - Validation sur des données réelles (Intra MEG) - dédié aux étudiants - Fournit une toolbox pour reproduire les résultats

5.2 ARTICLE



Contents lists available at ScienceDirect

Journal of Neuroscience Methods

journal homepage: www.elsevier.com/locate/jneumeth



Computational Neuroscience

Exceeding chance level by chance: The caveat of theoretical chance levels in brain signal classification and statistical assessment of decoding accuracy

Etienne Combrisson ^{a,b}, Karim Jerbi ^{a,c,*}

^a DYCOG Lab, Lyon Neuroscience Research Center, INSERM U1028, UMR 5292, University Lyon I, Lyon, France
^b Center of Research and Innovation in Sport, Mental Processes and Motor Performance, University of Lyon I, Lyon, France
^c Psychology Department, University of Montreal, QC, Canada

ARTICLE INFO

Article history:

Received 28 July 2014
Received in revised form 6 January 2015
Accepted 7 January 2015
Available online xxx

Keywords:

k-Fold cross-validation
Small sample size
Classification
Multi-class decoding
Brain-computer-interfaces (BCIs)
Machine learning
Binomial cumulative distribution
Classification significance
Decoding accuracy
MEG
EEG
Intracranial EEG

ABSTRACT

Machine learning techniques are increasingly used in neuroscience to classify brain signals. Decoding performance is reflected by how much the classification results depart from the rate achieved by purely random classification. In a 2-class or 4-class classification problem, the chance levels are thus 50% or 25% respectively. However, such thresholds hold for an infinite number of data samples but not for small data sets. While this limitation is widely recognized in the machine learning field, it is unfortunately sometimes still overlooked or ignored in the emerging field of brain signal classification. Incidentally, this field is often faced with the difficulty of low sample size. In this study we demonstrate how applying signal classification to Gaussian random signals can yield decoding accuracies of up to 70% or higher in two-class decoding with small sample sets. Most importantly, we provide a thorough quantification of the severity and the parameters affecting this limitation using simulations in which we manipulate sample size, class number, cross-validation parameters (*k*-fold, leave-one-out and repetition number) and classifier type (Linear-Discriminant Analysis, Naïve Bayesian and Support Vector Machine). In addition to raising a red flag of caution, we illustrate the use of analytical and empirical solutions (binomial formula and permutation tests) that tackle the problem by providing statistical significance levels (*p*-values) for the decoding accuracy, taking sample size into account. Finally, we illustrate the relevance of our simulations and statistical tests on real brain data by assessing noise-level classifications in Magnetoencephalography (MEG) and intracranial EEG (iEEG) baseline recordings.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Applying machine learning algorithms to brain signals in order to predict intentions or decode cognitive states has become an increasingly popular technique over the last decade. The surge in the use of machine learning methods in neuroscience has been largely fueled by the tremendous increase in brain-computer interface (BCI) and brain signal decoding research either using non-invasive recordings such as Electroencephalography (EEG) or Magnetoencephalography (MEG) (e.g. [Aloise et al., 2012](#); [Besserve et al., 2007](#); [Jerbi et al., 2011](#); [Krusienski and Wolpaw, 2009](#); [Toppi et al., 2014](#); [Waldert et al., 2008](#)) or with intracranial EEG (e.g. [Ball et al., 2009](#); [Derix et al., 2012](#); [Hamamé et al. \(2012\)](#); [Korcyn](#)

[et al. \(2013\)](#); [Lachaux et al., 2007a,b](#); [Leuthardt et al., 2004, 2006](#); [Mehring et al., 2004](#); [Pistohl et al., 2012](#); [Schalk et al., 2008](#); [Jerbi et al., 2007a,2009a,2013](#)). Machine learning and signal classification techniques are powerful and complex tools that have to be used with caution. While most machine learning experts are well aware of the various caveats to watch out for, certain theoretical limitations of these methods can easily elude students and neuroscience researchers new to the field of machine learning and brain-computer interface research.

In supervised learning, samples of a subset of the data and knowledge of their corresponding class (label) are used to train a model to distinguish between two or more classes. The trained classifier is then tested on the remaining data samples (the hold-out samples). This procedure is generally repeated several times by varying the subsets used for training and those used for testing, a standard procedure known as cross-validation. The percent of over-all correct label (or class) prediction across the test samples of the multiple folds is known as the correct classification

* Corresponding author at: Psychology Department, University of Montreal, QC, Canada.

E-mail address: karim.jerbi@umontreal.ca (K. Jerbi).

rate (sometimes called decoding accuracy). Conversely, the mean of misclassified samples over the folds is a measure of classifier prediction error.

The performance of a classifier in neural decoding studies is often assessed by how close its correct classification rate is to the maximum of 100%, or alternatively, how strongly it departs from the *chance-level* rate achieved by a classifier that would randomly associate the samples to the various classes. For instance, in a two-class or four-class classification problem, the probabilistic chance level indicating totally random classification is 50% or 25% respectively. Yet, although such probabilistic chance-levels widely applied in brain signal classification studies, they can be problematic because they are strictly speaking only valid for infinite sample sizes. While it will not come to anyone as a surprise that no study to date was able to acquire infinite data, it is intriguing how rarely brain signal classification studies acknowledge this limitation or take it into account. For a two-class classification problem with small sample size, 60%, 70% or even higher decoding percentages can in theory arise by chance (see simulation results below). As a consequence, for finite samples, a decoding percentage can only be considered reliable if it substantially, or better still, *significantly* departs from the theoretical level in statistical terms. But how can we assess the significance of the departure of a decoder from the outcome of total random classification? For a given sample size and a given number of classes, what would be the statistically significant threshold of correct classification that one needs to exceed in order to consider the decoding *statistically significant*? Although these questions have been widely recognized and addressed in the machine learning field (e.g. Kohavi, 1995; Martin and Hirschberg, 1996a,b), it is unfortunately often overlooked in the emerging field of brain signal classification which, incidentally, is often faced with low sample sizes for which the problem is even more critical.

Not all the previous brain decoding reports suffer from the caveat of using theoretical chance-level as reference. However, numerous studies only apply statistical assessment when testing for significant differences between the performance of multiple classifiers, or when comparing decoding across experimental conditions, but unfortunately neglect to provide a statistical assessment of decoding that accounts for sample size (e.g. Felton et al., 2007; Haynes et al., 2007; Bode and Haynes, 2009; Kellis et al., 2010; Hosseini et al., 2011; Sitaram et al., 2011; Hill et al., 2006; Wang et al., 2010; Bleichner et al., 2014; Babiloni et al., 2000; Ahn et al., 2013; Morash et al., 2008; Neuper et al., 2005; Kayikcioglu and Aydemir, 2010; Momennejad and Haynes, 2012). A number of such studies use theoretical percent chance-levels (e.g. 50% in a 2-class classification) as a reference against which classifier decoding performance is assessed. By doing so, such studies fail to account for the effect of finite sample size. This may have little effect in the case of large sample size or when extremely high decoding results are obtained, however, the bias and erroneous impact of such omissions can be critical for smaller sample sizes or when the decoding accuracies are barely above the theoretical chance levels.

Note however, that the rigorous assessment of significant classification thresholds is not equally ignored across the various types of neuronal decoding studies; it seems that the omissions (or unfortunate tendency to rely on the theoretical chance levels) are more common in more recent sub-branches of the neuronal decoding field. This is the case for signal classification and BCI studies based on non-invasive (fMRI, EEG and MEG) brain recordings in humans, and possibly electrocorticographic macro-electrode recordings in patients, where the methods (including classifiers, features and statistics) are less well-established than in the field of neuronal spike decoding in primates for instance.

In this brief article, we address caveats related to interpreting brain classification performances with small sample sizes. The paper is written with the broad neuroscience readership in mind

and is oriented, in particular, to students and researchers new to neural signal classification. First of all, we describe how applying signal classification to randomly generated signals can yield decoding accuracies (correct classification rates) that strongly depart from theoretical chance levels, with values up to 70% and higher with small sample sizes (instead of the expected theoretical 50% for 2-class decoding). Most importantly, we illustrate and quantify the phenomenon by using simulations in which we manipulate sample size, class number, cross-validation parameters and classifier type. In addition to raising a red flag of caution, we recommend practical alternatives to overcome the problem. We describe a straight-forward method to derive a statistically significant threshold that accounts for sample size and provides confidence intervals for the classification accuracy achieved by cross-validations. A reference table is also provided to allow readers to quickly look-up the percent correct classification thresholds that need to be exceeded in order to assert statistical significance of the findings for a range of possible sample sizes, classes and significance levels.

2. Materials and methods

2.1. Data simulation and classification

2.1.1. Generating normally distributed random data

In order to simulate a situation with classification results that approach the theoretical chance level, we generated 100 data sets of zero-mean Gaussian white noise. The normally distributed variables in each data set were generated in MATLAB (Mathworks Inc., MA, USA) via a pseudo-random number generator. Each one of the 100 data sets was randomly split into c subsets data (here we used $c = 2$ - or 4-classes) and we then evaluated the classification performance obtained by applying different classification algorithms to these simulated datasets. Because the variables in each 'simulated class' were drawn from the exact same Gaussian random distribution data set, applying supervised machine learning algorithms should fail to distinguish between classes and should theoretically yield chance-level classification rates (50% for $c = 2$ and 25% for $c = 4$). To examine the effect of sample size on how close the empirical classifications are to the theoretically expected chance level we varied the total number of samples n from 24 to 500. In other words, in the 2-class simulation for instance, the number of samples in each class varied from 12 to 250. Note that the code we implemented for the generation of random data for classification purposes is provided online (see Appendix A).

2.1.2. Classification algorithms

We implemented three types of machine learning algorithms: linear discriminant analysis (LDA), naïve Bayes (NB) classifier and a support vector machine (SVM), the latter with two different kernels: a linear kernel and a radial basis function (RBF) kernel. These three methods, which are frequently used for neural signal classification in the context of brain-computer interface research are briefly described in the following.

Linear discriminant analysis: LDA (Fisher, 1936) is a straight-forward and fast algorithm which assumes that the independent variables in each class are normally distributed with identical covariance (homoscedasticity assumption). For a two dimension problem, the LDA tries to find a hyperplane that maximizes the mean distance between the two classes while minimizing the inter-class variance. A multiclass problem can be tackled as a multiple two-class problem by discriminating each class from the rest using multiple hyperplanes.

Naive Bayesian classifier: The NB model (e.g. Fukunaga, 1990) is a probabilistic classifier that assigns features to the class to which they have the highest probability of belonging. NB assumes that the

features in each class are normally distributed and independent. The name arises from the fact that it is based on applying Bayes' theorem with strong (naive) independence assumptions.

Support vector machine: SVM ([Boser et al., 1992](#); [Burges, 1998](#); [Cortes and Vapnik, 1995](#); [Vapnik, 1995](#)) classifiers originate from statistical learning theory. An SVM searches for a hyperplane that maximize margins between the hyperplane and the closest features in the training set. For non-linearly separable classes, SVM uses a kernel function to project features in a higher dimensional space in order to reduce the nonlinear problem to a linear one, which is then separable by a hyperplane. The (Gaussian) Radial Basis Function (RBF) kernel is a popular choice. In this study, both linear and RBF kernels were used for SVM classification.

Details of the theoretical background of various classifiers can be found in standard statistics and machine learning textbooks and various reviews (e.g. [Lotte et al., 2007](#); [Wieland and Pittore, 2014](#)). Here, we used MATLAB implementation for the LDA and NB and the libsvm library for multi-class SVM.

2.1.3. Repeated and stratified k -fold cross-validation

To compute the decoding accuracy achieved by each one of the classifiers on the random data, we used standard stratified k -fold cross-validation. For a given data set size, all available N samples are partitioned into k folds, where $(k - 1)$ folds are used for training the classifier model (training set) and the remaining fold is used for validation (test set). This procedure is then repeated k times so that each fold is used once as test set. The stratified option ensures that each fold has approximately the same proportion of samples from each class as in the original dataset as a whole. The case $k = N$ (e.g. 200 folds in a data set of 200 samples) is called leave-one-out (LOO) cross-validation because one element is used to test the performance of a classifier trained on the rest of the data. Because k -fold cross-validation involves a random partition, the variance of the classifier can in theory be reduced by repeating the full cross validation procedure q times. Therefore, in addition to testing different classifier types, this study explores the effect of the following parameters: n (sample size, 20–500), k (number of cross-validation folds: 5, 10 and leave-one-out) and q (number of repetitions: 1, 5 and 20).

2.2. Statistical significance of classification using a binomial cumulative distribution

For a given number of classes c , the percent theoretical chance level of classification is given by $100/c$. For example, for a 4-class problem, the chance level is $100/4 = 25\%$. This threshold is based on the assumption of infinite sample size. In practice, the empirical chance level depends on the number of samples available. One way to address this limitation is to test for the statistical significance of the decoding accuracy. This can be done by assuming that the classification errors obey a binomial cumulative distribution, where for a total of n samples and c classes, the probability to predict the correct class at least z times by chance is given by:

$$P(z) = \sum_{i=z}^n \binom{n}{i} \times \left(\frac{1}{c}\right)^i \times \left(\frac{c-1}{c}\right)^{n-i}$$

Although neural signal classification studies predominantly evaluate decoding performance by how well the results depart from the theoretical chance level, several BCI studies have in addition, used the binomial cumulative distribution to derive statistical significance thresholds (e.g. [Ang et al., 2010](#); [Demandt et al., 2012](#); [Pistohl et al., 2012](#); [Waldert et al., 2007, 2008, 2012](#)). In this study, we use the MATLAB (Mathworks Inc., MA, USA) function *binoinv* to compute the statistically significant threshold $St(\alpha) = binoinv(1 - \alpha, n, 1/c) \times 100/n$, where α is the significance level given by $\alpha = z/n$

(i.e. the ratio of tolerated false positives z – i.e. number of observations correctly classified by chance with respect to all observations n). For instance, for a sample size of $n = 40$ and a 2-class classification problem ($c = 2$), computing the threshold for statistical significance of the decoding at $\alpha = 0.001$ using the above formulation yields 70.0%. In other words, at $n = 40$, any decoding percentage below 70% is not statistically significant ($p < 0.001$), whereas if one relied on the theoretical threshold for two classes (i.e. 50%) a decoding accuracy of 67% might have been considered relevant. [Table 1](#) provides the minimal thresholds as a function of selected sample sizes, class number and significance levels. Note that code for the calculation of these analytical significance levels is provided online (see [Appendix A](#)).

2.3. Statistical significance of classification using permutation tests

The statistical significance of decoding can also be assessed by non-parametric statistical methods, namely using permutation tests ([Good, 2000](#); [Nichols & Holmes, 2002](#)). By randomly permuting the observations across classes and calculating classification accuracy at each permutation, it is possible to establish an empirical null distribution of classification accuracies on random observations. The tails of this distribution can then be used to determine significance boundaries for a given rate of tolerated false positives (i.e. correct classifications that occur by chance). For instance, if the original (without randomization) classification accuracy is higher than the 95 percentile of empirical performance distribution established by randomly permuting the data, then one can assert that the original classification is significant with $p < 0.05$. The advantage of this empirical approach is that it does not require particular assumption about statistical properties of the samples.

An intuitive illustration of this procedure would be as follows: one performs for example 99 random permutations of the labels (classes) in the data and computes the classification accuracy for each permutation. This provides an empirical distribution of 99 classification accuracy values. Now if the classification performance obtained with the original (unpermuted) data is higher than the maximum of the empirical distribution, one can conclude that it is significant with $\alpha = 0.01$.

Permutations test provide a useful empirical approach to deriving statistical significance of classifier performance (e.g. [Golland and Fischl, 2003](#); [Ojala and Garriga, 2010](#); [Meyers and Kreiman, 2011](#)). To demonstrate the utility to derive significance boundaries as a function of sample size and thus compare it to the use of the binomial formula. To this end, we used simulated random data with associated labels (as described in Section 2.1) and computed the classification performance (using LDA) for 10,000 permutations (randomly exchanging labels of the original observations). From this we derived the accuracy thresholds that correspond to the 99%, 99.9% and 99.99% percentile of the distribution (i.e. $p < 0.01$, $p < 0.001$, and $p < 0.0001$ respectively). This was done for each sample size value n (20–500), which allowed us to depict the evolution of the empirical significance boundaries as a function of sample size. Note that code for the calculation of permutation-based empirical significance levels is provided online (see [Appendix A](#)).

2.4. Classification of baseline data from real brain signals

Because real data does not necessarily have the same properties as those implemented in our random data simulations (zero-mean Gaussian white noise), we also calculated the correct classification rate (as a function of sample size) that is achieved when classifying real brain data that do not contain any true discrepancies. This was carried out for pre-stimulus or baseline recordings in MEG (4 subjects) and with intracranial EEG recordings (4 patients). The

Table 1

Look-up table for statistically significant classification performance. Minimal correct classification rate (%) to assert statistical significance (at a given p -value) as a function of sample size n and number of classes c . Threshold values are based on the binomial cumulative distribution function and are rounded to the first digit.

n	c	2-Classes				4-Classes				8-Classes						
		$p < 0.05$		$p < 0.01$		$p < 10^{-3}$		$p < 10^{-4}$		$p < 0.05$		$p < 0.01$		$p < 10^{-3}$		
		$p < 0.05$	$p < 0.01$	$p < 10^{-3}$	$p < 10^{-4}$	$p < 0.05$	$p < 0.01$	$p < 10^{-3}$	$p < 10^{-4}$	$p < 0.05$	$p < 0.01$	$p < 10^{-3}$	$p < 10^{-4}$	$p < 0.05$	$p < 0.01$	
20	70.0%	75.0%	85.0%	90.0%	40.0%	50.0%	55.0%	65.0%	25.0%	30.0%	40.0%	45.0%				
40	62.5%	67.5%	75.0%	77.5%	37.5%	42.5%	47.5%	52.5%	22.5%	25.0%	30.0%	35.0%				
60	60.0%	65.0%	70.0%	73.3%	35.0%	38.3%	43.3%	46.7%	20.0%	23.3%	26.7%	30.0%				
80	58.7%	62.5%	67.5%	70.0%	32.5%	36.2%	41.2%	43.7%	18.7%	21.2%	25.0%	27.5%				
100	58.0%	62.0%	65.0%	68.0%	32.0%	35.0%	39.0%	42.0%	18.0%	21.0%	24.0%	26.0%				
200	56.0%	58.0%	61.0%	63.0%	30.0%	32.5%	35.0%	37.0%	16.5%	18.0%	20.0%	22.0%				
300	54.7%	56.7%	59.0%	60.7%	29.0%	31.0%	33.0%	34.7%	15.7%	17.0%	18.7%	20.0%				
400	54.0%	55.7%	57.7%	59.2%	28.5%	30.0%	31.7%	33.2%	15.2%	16.5%	17.7%	19.0%				
500	53.6%	55.2%	57.0%	58.2%	28.2%	29.6%	31.2%	32.4%	15.0%	16.0%	17.2%	18.2%				

rationale here is that baseline (pre-stimulus) data is not expected to show any genuine discriminative brain patterns related to post-stimulus events, and as such, it is comparable to random background noise. Therefore, signal classification on these baseline periods should fail, and the accuracies that classifiers achieve can be taken as an empirical representation of chance-level decoding.

2.4.1. Illustrative data from MEG rest activity

We used illustrative data from 4 subjects scanned with a whole-head MEG system (151 sensors; VSM MedTech, BC, Canada) acquired at 1250 Hz sampling rate and with a band pass filter of 0–200 Hz. The participants provided written informed consent, and the experimental procedures were approved by the Institutional Review Board and by the National French Science Ethical Committee. The MEG data segments used for the purpose of the current analysis were extracted from the pre-stimulus baseline of a visuomotor MEG experiment (Jerbi et al., 2007b), and each trial was assigned one of 2 (or of 4) arbitrary labels for the 2-class (or 4-class) classification. Oscillatory alpha (8–12 Hz) power was computed using Hilbert transform and subsequently used as feature in an LDA-based classification procedure. We used 10-fold cross-validation and the whole procedure was repeated for increasing values of trial numbers (sample size n) ranging from 20 to 200 (in steps of 8).

2.4.2. Illustrative data from intracranial EEG baseline activity

We used illustrative data from 4 epilepsy patients stereotactically implanted with intracranial depth electrodes (0.8 mm diameter, 10–15 contact leads, DIXI Medical Instruments, Besançon, France). The intracerebral EEG (iEEG) recordings were conducted using a video-SEEG monitoring system (Micromed, Treviso, Italy), which allowed for the simultaneous recording from 128 depth-EEG electrode sites (More details of the routine SEEG acquisitions in Jerbi et al., 2009b). The data were bandpass filtered online from 0.1 to 200 Hz and sampled at 1024 Hz. The recordings were performed at the epilepsy department of the Grenoble University Hospital (headed by Dr. Philippe Kahane). All participants provided written informed consent, and the experimental procedures were approved by the Institutional Review Board and by the National French Science Ethical Committee.

The data segments used here were extracted from the pre-stimulus (baseline) of a standard motor task and each trial was associated with one of 2 (or of 4) labels for the 2-class (or 4-class) classification. The labels assigned to each pre-stimulus baseline trial were in fact the genuine post-stimulus events for the same trials (but no true discrimination can be expected prior to stimulus onset as the post-stim event could not be known or inferred during the pre-stimulus period). Broadband gamma (60–250 Hz) power was computed using Hilbert transform and subsequently used as feature in an LDA-based classification procedure. As for the MEG

data, we used 10-fold cross-validation and the whole procedure was repeated for increasing values of trial numbers (sample size n) ranging from 20 to 200 (in steps of 8).

3. Results

3.1. Empirical evaluation of chance level decoding as a function of sample size

Fig. 1 shows the decoding accuracies obtained by conducting 10-fold cross validation on 100 randomly generated data sets. The decoding is depicted as a function of increasing sample size (from 24 to 500) and for the case of 2-class (left column) and 4-class (right column) classification. Although the theoretical chance levels for these configurations are 50% and 25% respectively, the results show how much the empirical decoding accuracies obtained with random data deviate from these probabilistic values.

The small sample size problem: as expected, the variance of the decoding accuracy across the 100 simulated random data sets is high, and the more so for small sample sizes. As illustrated in Fig. 1, while the decoding does converge toward the theoretical chance level as the sample size increases, the values achieved with small sample size ($n < 100$) can be disturbingly high. For instance, the highlighted examples (solid black line) in panels (a) to (f) illustrate how decoding accuracies as high as 70% for 2-class classification (or 50% for 4-classes) can be observed even when conducting classification on subsets of randomly generated data with randomly associated labels.

The small sample issue is persistent and qualitatively similar across all classifiers used. The first three rows of Fig. 1 show the results obtained with LDA, NB and SVM (with an RBF kernel). Panels (g) and (h) of Fig. 1 show that cross-validation results in all three classifiers have comparable deviation across the 100 simulated data sets. The variance of cross-validation over the 100 random data sets is high for small sample sizes (<200 observations) and drops off with increasing sample size.

3.2. Tweaking cross-validation parameters does not solve the small sample problem

It might be tempting to think that changing the cross-validation parameters might be a way to get around the small sample problem illustrated here. To address this we evaluated the impact of varying (a) the number of cross-validation folds, and (b) the number of repetitions of the cross-validation, on the reported deviation of the cross-validation results (cf. Fig. 1g and h) across the 100 data sets and all sample sizes. The results in Fig. 2(a–c) show that applying 5- and 10-fold cross-validation to the random data yielded substantially the same results, and that leave-one-out (LOO) cross-validation actually provided worse results (i.e.

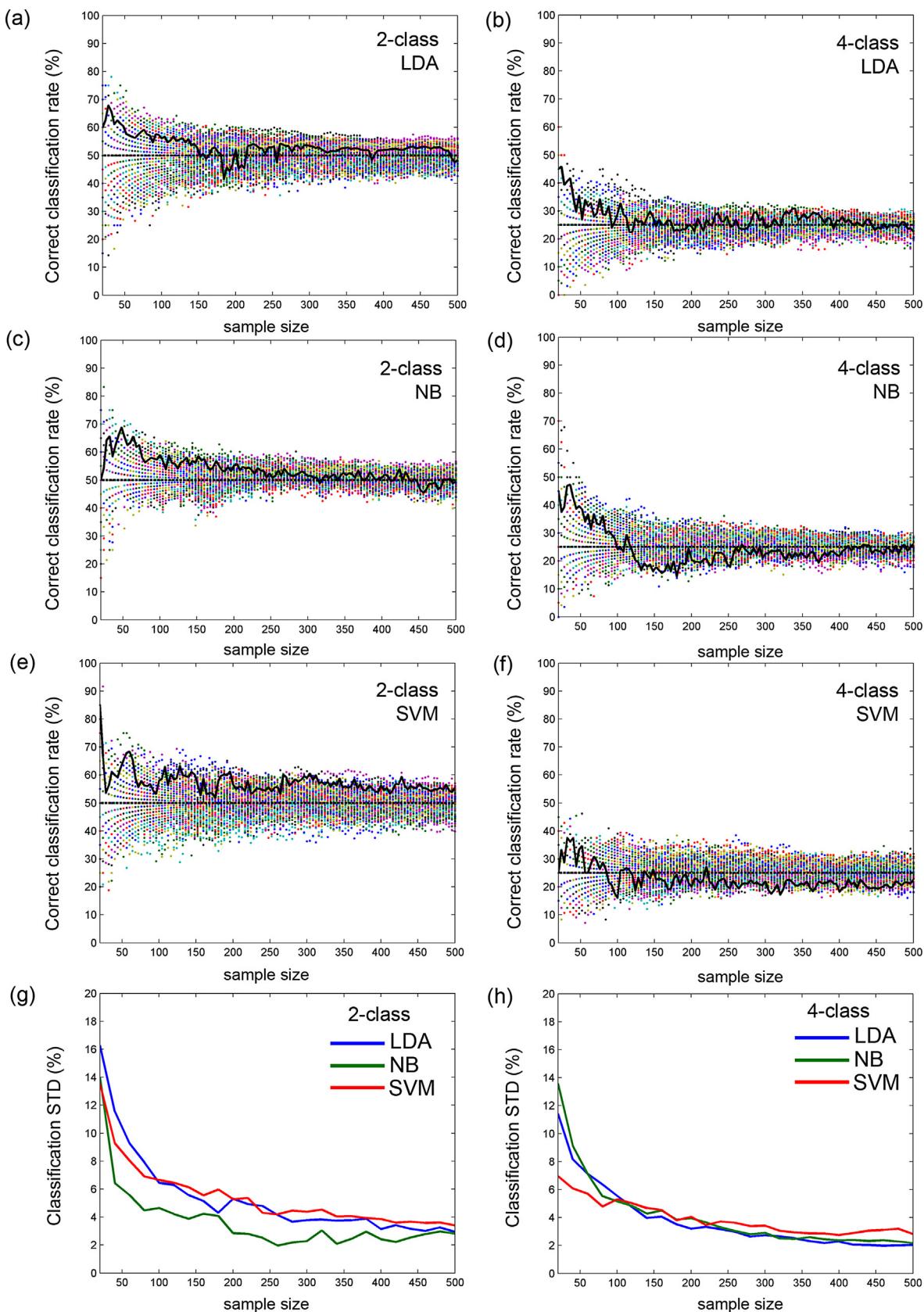


Fig. 1. Classifier decoding rates as a function of sample size when applied to random data sets using 10-fold cross-validation. (a) Two-class LDA classification rate (%) as a function of sample size (empirical results increasingly deviate from the 50% chance-level as the sample size gets smaller). The backline shows the evolution of cross-validation results for one specific data set out of the 100 depicted in multiple colors. (b) Same as panel (a) but using 4-class classification, i.e. at each sample size n , the data is split into 4 virtual classes instead of two. (c and d) Same as (a and b) but for a Naïve Bayesian classifier. (e and f) Same as (a and b) but for an SVM classifier using an RBF kernel. (g) Evolution of cross-validation standard deviation across the 100 data sets for each of the three classifiers for 2-class decoding. (h) Same as panel (g) but for 4-class decoding.

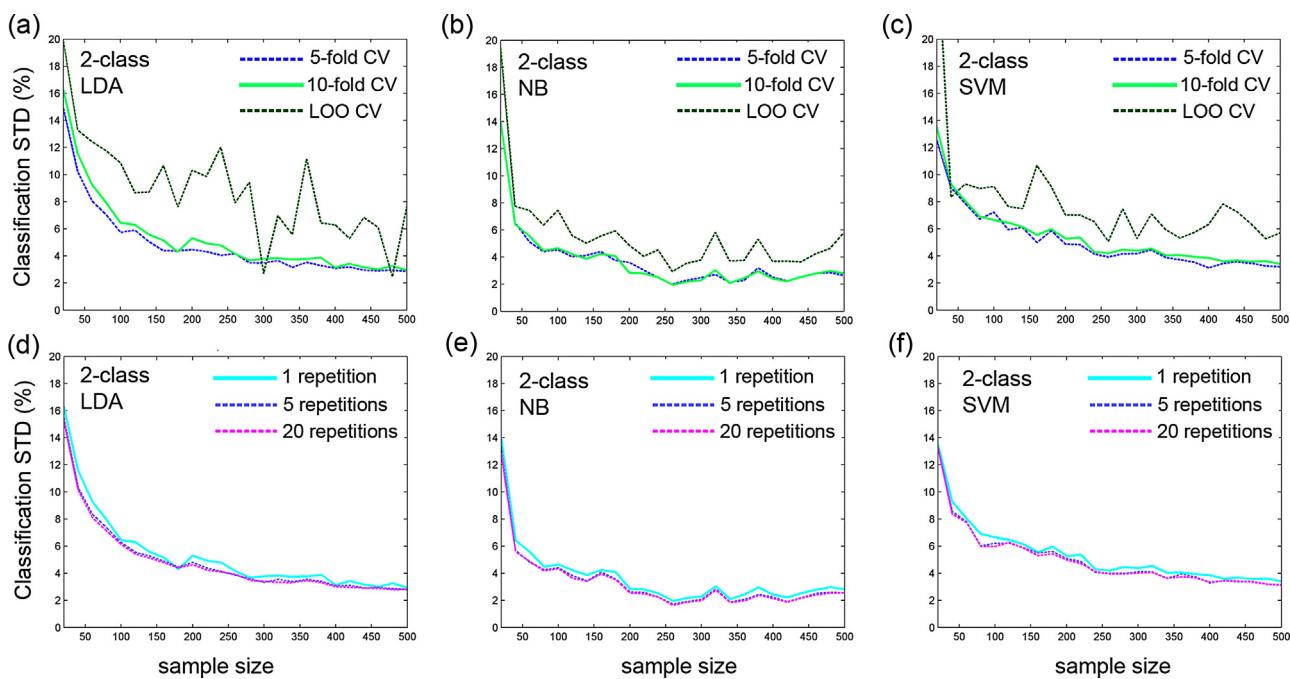


Fig. 2. Effect of cross-validation parameters on the variability of 2-class decoding performance computed across 100 sets of random data. (a–c) Effect of the number of folds (k): drop in cross-validation variance as sample size n increases, shown for $k=5$, $k=10$ (default), and $k=n$ (i.e. leave-one-out) and for all three classifiers LDA (panel a), NB (panel b) and SVM (panel c). (d–f) Effect of cross-validation repetition number: drop in cross-validation variance as sample size n increases, shown for repetition values $q=1$ (default), $q=5$, and $q=20$ and for all three classifiers LDA (panel d), NB (panel e) and SVM (panel f). Note that the strong deviation from 50% chance-level for small sample sizes is persistent across all panels, and appears to be worst for LOO cross-validation with LDA.

higher variance). Moreover, repeating the cross-validation procedure (whether 5 or 20 times) achieved a negligible reduction of variance (Fig. 2d and e). Overall, these observations indicate that neither changing the number of folds nor to the number of overall repetitions has an impact on the variance of decoding accuracy (i.e. cross-validation results) across the 100 sets of Gaussian white noise.

3.3. Estimating statistical significance of decoding accuracy: binomial formula and permutation tests

Panels (a) and (b) in Fig. 3 show the evolution of the minimal statistically significant decoding rate as a function of sample size (respectively for 2- and 4-classes) using the binomial cumulative distribution (described in Section 2.2). The plots depicted for three distinct significance levels (10^{-2} , 10^{-3} and 10^{-4}) all show that the minimal correct decoding rate that is required in order to assert significance, decreases as the number of samples increases. Given small sample sizes (e.g. below 100 observations), to be statistically significant, the decoding accuracy must be substantially higher than the probabilistic chance level. For example, for 40 observations, a 2-class decoding is statistically significant (at $p < 0.001$) only if it exceeds the threshold of 75%. Note that for sample sizes as high as 500 observations, statistical significance still requires correct decoding higher than 55% (at $p < 0.01$), i.e. at least 5% above the theoretical chance level. A more comprehensive overview of the statistical decoding thresholds (wider ranges of p -values and of class number) computed for selected sample sizes (20–500), is provided in Table 1.

Panels (c) and (d) in Fig. 3 depict not only the evolution of the decoding boundaries for 2-class and 4-class decoding, using the binomial formula but also using the permutation test approach (see Section 2.3). Interestingly, the boundaries (for each level of admitted false positives) using both methods are reasonably close. The boundaries obtained with permutations show a slight tendency to

be more restrictive than the binomial formula. While this is a little more apparent for small values of n , the difference between the two methods rapidly vanishes as n increases.

3.4. MEG and iEEG baseline data reveal erroneously high decoding results

Fig. 4 depicts the results of the empirical estimation of *de facto* chance-level decoding in illustrative MEG and iEEG data segments taken during pre-stimulus baseline periods (where no decoding is theoretically expected). Similarly to our findings using random data simulations (Fig. 1 a–f), the baseline MEG and iEEG data trials also led to decoding rates that strongly departed from the theoretical chance levels of 50% for 2-class classification and 25% for 4-class classification. Also in line with the results of the simulated data, the effects observed here were again highest for small sample sizes and dropped off slowly with increasing n . Note that the results in Fig. 4 show consistent performances across the 4 subjects at each value of n (with MEG and with iEEG). Finally, the superimposed gray curves (which depict the significance boundary given by the binomial formula as a function of sample size) nicely follow the trend of the % correct classification rate, and also illustrate cases of tolerated false positives for a given alpha.

4. Discussion

The current study has two primary take-home messages. The first is emphasizing the importance of watching out for a potential caveat that may arise when using departure from the theoretical chance-level as evidence for meaningful decoding. By launching various classifiers on normally distributed random data (Gaussian white noise), we demonstrate and quantify to which extent small samples lead to decoding accuracies that overshoot the chance-level merely by chance. This observation follows from the fact that

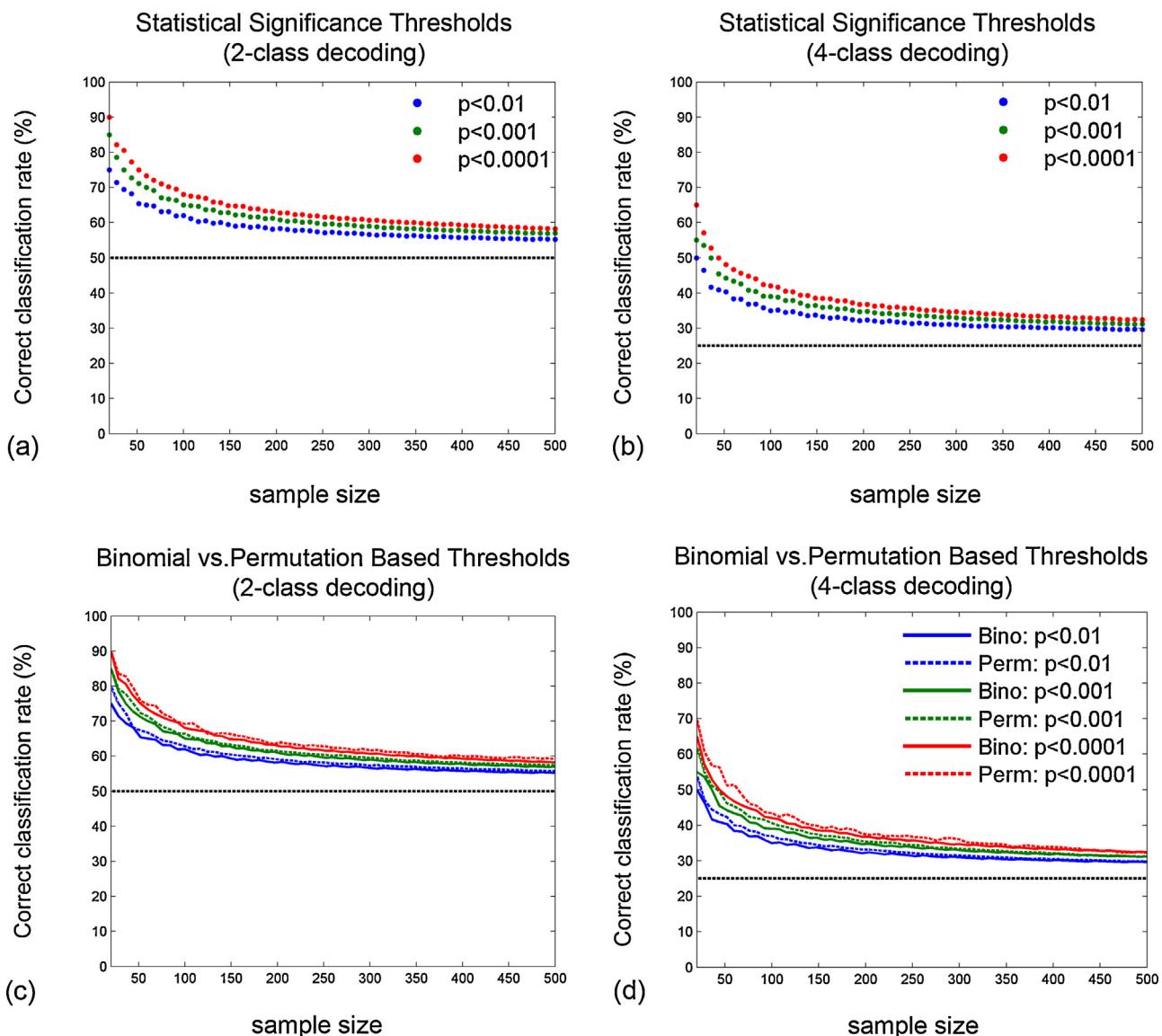


Fig. 3. Estimation of the statistical significance thresholds for 2- and 4-class classification as a function of sample size (assuming prediction errors are binomially distributed). Panels (a) and (b) show the evolution of the minimal statistically significant decoding rate as a function of sample size (respectively for 2 and 4 classes) using the binomial cumulative distribution (see Section 2.2). The plots were derived for significance levels 10^{-2} , 10^{-3} and 10^{-4} . As an example: panel (a) indicates that given a total of 100 data samples, a 2-class decoding result can only be considered statistically significant (at $p < 0.001$) if it exceeds 65%. This minimal value drops to 59% for 300 samples, but rises up to 75% if only 40 data points are available (See Table 1). Panels (c) and (d) show the same statistically significant decoding rate as a function of sample size (respectively for 2 and 4 classes) but now using both the binomial cumulative distribution (continuous lines) and the data-driven permutation-based approach (dashed lines) applied to the simulated random data (see Section 2.3 for details).

small samples are a bad approximation of true randomness and that as a result, the level $100/c$ (where c is the number of classes) is a purely theoretical chance-level that only holds for infinite sample sizes and that is particularly violated for small sample sizes. This basic fact is often overlooked in the neuronal decoding literature, where it is sometimes tempting to interpret for instance a 65% decoding accuracy in a 2-class classification as reflecting true neuronal decoding, without taking sample size into account. We have shown here that such levels of classification can be achieved with small samples of randomly generated data. This issue is not problematic for huge data samples, however, in data obtained from brain signal recordings in humans (such EEG or MEG), sample size can often be small. The effect of small samples on the reliability of probabilistic thresholds is therefore of particular importance in neural decoding and brain-computer interface studies. This effect is possibly even more critical when attempting to decode neuronal

signals acquired using intracranial recordings (electrocorticography or stereoacoustic-EEG) and in clinical BCI applications where even less data samples might be available.

Furthermore, our exploration of the effects of classifier type (LDA, NB and SVM), cross-validation partition (number of folds) and cross-validation repetition number (up to 20), indicates that none of these parameters has a noticeable impact on the variance of the classification when applied to random data. The small sample size problem cannot be circumvented by tweaking these parameters and even for larger sample sizes of white noise any reduction in classification variance remains negligible. Note that the explored parameters and classifier comparisons performed here only address the variance and bias of the techniques when applied to normally distributed random data, reviews and comparisons of classifiers can be found elsewhere (e.g. [Lotte et al., 2007](#)).

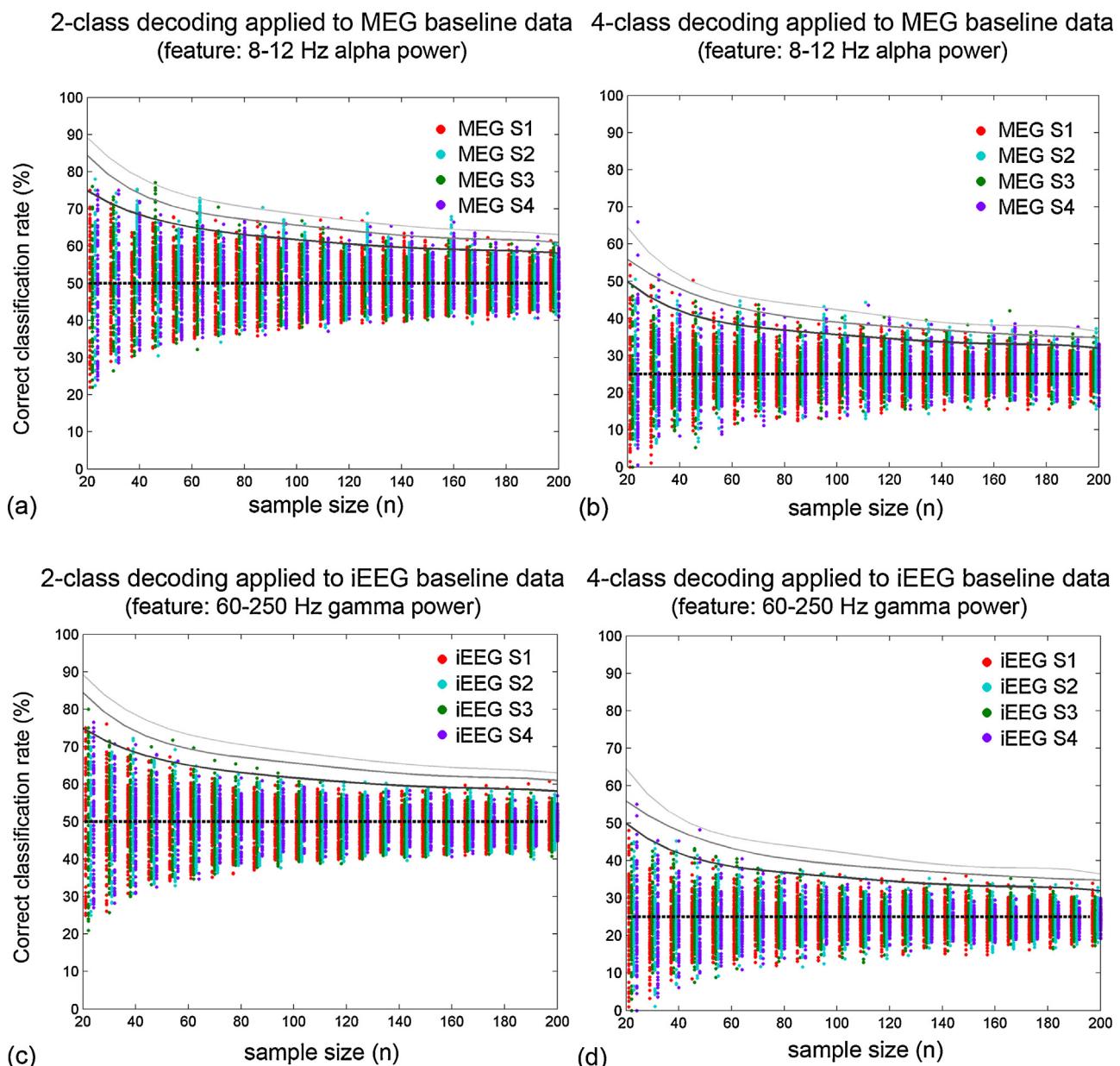


Fig. 4. Experimental assessment of chance-level classification accuracy in baseline (pre-stimulus) MEG and intracranial EEG data. (a) Two-class LDA classification rate (%) of MEG baseline data (alpha power features) as a function of sample size (illustrative data in 4 participants MEG S1–S4). The gray lines show the evolution of statistical significance boundaries computed with the binomial formula. Points lying above the gray lines thus represent false positives (type I errors) (b) Same as panel (a) but using 4-class classification, (c and d) Same as (a and b) but using baseline data (gamma power features) from intracranial EEG recordings (illustrative data in 4 epilepsy patients iEEG S1–S4).

Ten-fold cross-validation, which we used here as default, has been shown to be a reasonable choice providing low variance (Kohavi, 1995; Martin and Hirschberg, 1996a). Nevertheless, we also explored 5-fold and LOO cross-validation, alongside repetition number (Fig. 2). We found that none of these parameters could help reduce the cross-validation variance for low sample sizes. What is more, leave-one-out cross-validation showed even higher variability (in particular when using LDA), which is in agreement with previous reports suggesting that, despite its low bias, its high variance leads to unreliable estimates (Efron, 1983). Note that estimating the variance of cross-validation results across its k folds is generally problematic. Naive estimators that do not take into account error correlations due to the overlap between training and test sets (across the cross-validation folds) can severely underestimate variance (Bengio and Grandvalet, 2004). The cross-validation

variances reported here were computed across the 100 independent data sets of Gaussian white noise.

The second take-home message from our study is an important reminder that one way to overcome this limitation is to seek statistically significant thresholds on decoding accuracy, rather than relying solely on the theoretical chance-level to claim successful decoding. This has been demonstrated here using a sample-size dependent threshold computation derived from the binomial cumulative distribution function. The underlying assumption that the number of errors is binomially distributed is commonly used in statistical learning (Kohavi, 1995; Breiman et al., 1984) but the statistical bounds it provides are unfortunately rarely exploited in brain signal decoding studies (e.g. Quiroga and Panzeri, 2009; Müller-Putz et al., 2008; Ang et al., 2010; Arvaneh et al., 2013; Demandt et al., 2012; Galan et al., 2014; Lampe et al., 2014; Pisto

et al., 2012; Waldert et al., 2008, 2007, 2012). Kohavi (1995) provides a proof that k -fold cross-validation is binomial if the classifier induction method is stable under cross-validation. Note also that the validity of the assumption that prediction errors are binomially distributed has also been demonstrated for the specific case of 10-fold cross-validation with small samples (Martin and Hirschberg, 1996b). The latter study also emphasizes that the textbook formula based on the normal approximation to the binomial is not a good approximation to the confidence interval of an error rate estimate for small samples.

In addition to the binomial formula, we have also demonstrated the use of permutation tests as an alternative method to derive statistical significance boundaries for classifier performance as a function of sample size (Fig. 3c and d). Permutation tests provide a reliable and data-driven approach to the problem and has been proposed and used in numerous previous studies (e.g. Golland and Fischl, 2003; Ojala and Garriga, 2010; Meyers and Kreiman, 2011). Our analysis shows how, via multiple random shuffling of the data (or class labels), permutation tests can provide an estimate of sample-size dependent chance-level decoding accuracy. These empirical chance levels need to be exceeded in order to assert significance of a classification for a given rate of tolerated false positives. When applied to random noise signals, we found that the significance boundaries derived using permutations are reasonably close to those obtained using the binomial formula. Deciding which of the two approaches is more convenient when applied to real brain signals will likely depend on the data at hand. Permutation tests do not make any assumptions about the distribution of the data and provide a data-driven approach; however they also come with the burden of high computational cost, which dramatically increases with sample size, and with the level of statistical significance required.

Meyers and Kreiman (2011) note that deriving significance thresholds via the binomial formula as discussed here and elsewhere (e.g. Quiroga and Panzeri, 2009) comes with theoretical limitations that one should keep in mind, in particular, when combined with cross-validation; its application to mean performance over all folds violates the assumption of data point independence and leads to p -values that are too small. From a practical perspective, the impact of this theoretical limitation is likely to depend on the data at hand and on the selected cross-validation parameters. Simulations show that cross-validation parameters (number of cross-validation folds and repetitions) have an impact on the cumulative distribution function of classification accuracies (e.g. Noirhomme et al., 2014). As a result, cross-validation parameters, alongside classifier type and feature space, collectively lead to deviations from a binomial cumulative distribution. These deviations can be significant for small sample sizes (e.g. $N < 100$), which would advocate against using the binomial formula for statistical assessments under such circumstances (Noirhomme et al., 2014). In contrast, permutation tests being inherently data-driven, do take cross-validation parameters into account. As far as the Gaussian white noise simulated in the current study is concerned, permutation tests and the binomial formula appear to provide reasonably similar significance boundaries. Comparing the output of the binomial formula and (the more time consuming) permutation test, on at least a portion of the data, could be a pragmatic way to decide on whether the former provides a suitable and fast approximation of the latter.

Moreover, our analysis of decoding accuracy using real brain signals (with random labeling) is in line with our simulation results. This is a reassuring finding, as the latter were based on zero-mean Gaussian white noise while the former were based on power features (alpha and gamma-bands) derived from real brain data. The baseline-period MEG and iEEG data suggest that the binomial formula provides a reasonable estimation of chance-level

decoding in these data sets. As a general rule, whenever possible, it is highly recommended to use baseline data as a recording in which no task-dependent encoding occurs and thus within chance-level decoding is expected. Comparisons with pre-stimulus (baseline) decoding performances should be used as an additional sanity check whenever such data is available (Meyers and Kreiman, 2011).

An alternative framework that can be applied to measure and compare classifier performance, is the use of receiver operating characteristic (ROC) analysis and in particular the area under the ROC curve (AUC) (Ling et al., 2003; Huang and Ling, 2005; Bradley, 1997). It has also been shown that calculating the probability density function (pdf) for each point on a ROC curve for any given sample size can be used to produce confidence intervals for ROC curves that are valid for small sample sizes (Tilbury et al., 2000). Adaptations of this method might be particularly suited to assessing classifier performance in BCI research (Hamadicharef, 2010). Other solutions that have been proposed to tackle the small sample size problem include frameworks that combine cross-validation with bootstrapping (e.g. Fu et al., 2005) and the use of class-dependent PCA in conjunction with linear discriminant feature extraction (Das and Nenadic, 2009). It is noteworthy that a few authors have even suggested that classification studies should be based primarily on effect size estimation with confidence intervals, rather than on significance tests and p -values (Berrar and Lozano, 2013).

In summary, the notion of statistical significance for decoding rates (or prediction error) and the small sample size problem have been tackled in the field of statistical learning for a long time (e.g. Raudys and Jain, 1991). However, these notions have not been sufficiently acknowledged in the relatively recent surge in application of machine learning methods in neuroscience. In the worse cases, this can unfortunately lead to erroneous interpretation of decoding results. Beyond its importance for brain-computer interface research specifically, signal classification is also increasingly used in neuroscience with the broader aim of elucidating the functional role of specific neuronal features (i.e. unraveling neuronal encoding by investigating single-trial neuronal decoding). Incidentally, this is where researchers are likely to be tempted to consider low (but above chance-level) decoding accuracies (e.g. 68% in a two-class classification) as being relevant. The use of confidence intervals and robust estimation of statistical significance is of particular importance in such studies, and even more so in cases with low trial numbers (e.g. below 150 observations). Machine learning and cross-validation accuracy in multi-class decoding may therefore not be thought of as a less-strict approach that can circumvent traditional rigorous statistical comparisons of data from multiple experimental conditions. Finally, whether signal classification is used in a BCI context *stricto sensu* or within a framework to conduct basic neuroscience analysis, we highly recommend systematically reporting the decoding accuracy as well as its statistical significance. We hope that the simulation results, statistical approaches and practical recommendations discussed here will be helpful in illustrating the problem and providing ways of tackling it.

Acknowledgements

Etienne Combrisson is currently supported by a Ph.D. Scholarship awarded by the Ecole Doctorale Inter-Disciplinaire Sciences-Santé (EDISS), Lyon, France. This work was partly performed within the framework of the LABEX CORTEX (ANR-11-LABX-0042) of Université de Lyon, within the program ANR-11-IDEX-0007. This research was undertaken, in part, thanks to funding from the Canada Research Chairs program. The authors are grateful for the collaboration with the research and clinical staff of the Magnetoencephalography (MEG) center at the Pitié-Salpêtrière Hospital

in Paris and the University Hospital in Grenoble (Dr. Philippe Kahane).

Appendix A.

A. Software availability: The MATLAB scripts and functions that were developed and used in this study have been made available online for the community. The provided code can be used to generate, label and classify random data. It also provides routines to compute and plot, as a function of sample size, (a) analytical chance levels via the binomial formula as well as (b) empirical chance levels via permutation tests. We hope that this set of tools will help students and researchers replicate and extend our analyses. The code can be downloaded from Mathwork's File Exchange platform at the following URL: <http://www.mathworks.fr/matlabcentral/fileexchange/48274-random-data-classification>

References

- Ahn M, Ahn S, Hong JH, Cho H, Kim K, Kim BS, et al. Gamma band activity associated with BCI performance: simultaneous MEG/EEG study. *Front Hum Neurosci* 2013;7. Available from: <http://journal.frontiersin.org/journal/10.3389/fnhum.2013.00848/full>.
- Aloise F, Schettini F, Aricò P, Salinari S, Babiloni F, Cincotti F. A comparison of classification techniques for a gaze-independent P300-based brain-computer interface. *J Neural Eng* 2012;9(4):045012.
- Ang KK, Guan C, Sui Geok Chua K, Ang BT, Kuah C, Wang C, et al. Clinical study of neurorehabilitation in stroke using EEG-based motor imagery brain-computer interface with robotic feedback. *IEEE Trans Rehabil Eng* 2010;8(June (2)):5549–52 [cited 2014 Jul 4]. Available from: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5626782.
- Arvaneh M, Guan C, Ang KK, Quek C. EEG data space adaptation to reduce intersession nonstationarity in brain-computer interface. *Neural Comput* 2013;25(May (8)):2146–71.
- Babiloni F, Cincotti F, Lazzarini L, Millán J, Mourino J, Varsta M, et al. Linear classification of low-resolution EEG patterns produced by imagined hand movements. *IEEE Trans Rehabil Eng* 2000;8(June (2)):186–8.
- Ball T, Schulze-Bonhage A, Aertsen A, Mehring C. Differential representation of arm movement direction in relation to cortical anatomy and function. *J Neural Eng* 2009;6(1):016006.
- Bengio Y, Grandvalet Y. No unbiased estimator of the variance of k-fold cross-validation. *J Mach Learn Res* 2004;5:1089–105.
- Berrada D, Lozano JA. Significance tests or confidence intervals: which are preferable for the comparison of classifiers? *J Exp Theor Artif Intell* 2013;25(June (2)):189–206.
- Besserve M, Jerbi K, Laurent F, Baillet S, Martinier J, Garnier L, et al. Classification methods for ongoing EEG and MEG signals. *Biol Res* 2007;40(4):415–37.
- Bleichner MG, Jansma JM, Sellmeijer J, Raemaekers M, Ramsey NF. Give me a sign: decoding complex coordinated hand movements using high-field fMRI. *Brain Topogr* 2014;27(March (2)):248–57.
- Bode S, Haynes J-D. Decoding sequential stages of task preparation in the human brain. *Neuroimage* 2009;45(April (2)):606–13.
- Boser BE, Guyon IM, Vapnik VN. A training algorithm for optimal margin classifiers. *ACM* 1992;144–52 [cité 25.07.14].
- Bradley AP. The use of the area under the ROC curve in the evaluation of machine learning algorithms. *Pattern Recognit* 1997;30(July (7)):1145–59.
- Breiman L, Friedman JH, Olshen R, Stone CJ. Classification and regression trees; 1984.
- Burges CJ. A tutorial on support vector machines for pattern recognition. *Data Min Knowl Discov* 1998;2(2):121–67.
- Cortes C, Vapnik V. Support-vector networks. *Mach Learn* 1995;20(September (3)):273–97.
- Das K, Nenadic Z. An efficient discriminant-based solution for small sample size problem. *Pattern Recognit* 2009;42(May (5)):857–66.
- Demandt E, Mehring C, Vogt K, Schulze-Bonhage A, Aertsen A, Ball T. Reaching movement onset- and end-related characteristics of eeg spectral power modulations. *Front Neurosci* 2012;6. <http://www.frontiersin.org/journal/10.3389/fnhum.2012.00065/full>.
- Derix J, Ilijina O, Schulze-Bonhage A, Aertsen A, Ball T. "Doctor" or "darling"? Decoding the communication partner from ECg of the anterior temporal lobe during non-experimental, real-life social interaction. *Front Hum Neurosci* 2012;6. <http://www.frontiersin.org/journal/10.3389/fnhum.2012.00251/full>.
- Efron B. Estimating the error rate of a prediction rule: improvement on cross-validation. *J Am Stat Assoc* 1983;78(June (382)):316–31.
- Felton EA, Wilson JA, Williams JC, Garell PC. Electrocorticographically controlled brain-computer interfaces using motor and sensory imagery in patients with temporary subdural electrode implants. Report of four cases. *J Neurosurg* 2007;106(March (3)):495–500.
- Fisher RA. The use of multiple measurements in taxonomic problems. *Ann Eugen* 1936;7(September (2)):179–88.
- Fu WJ, Carroll RJ, Wang S. Estimating misclassification error with small samples via bootstrap cross-validation. *Bioinformatics* 2005;21(May (9)):1797–86.
- Fukunaga K. *Introduction to statistical pattern recognition*. 2nd ed. Boston: Academic Press; 1990.
- Galan F, Baker MR, Alter K, Baker SN. Missing kinaesthesia challenges precise naturalistic cortical prosthetic control, May. Report no: 004861; 2014, <http://biorxiv.org/lookup/doi/10.1101/004861>.
- Golland P, Fischl B. Permutation tests for classification: towards statistical significance in image-based studies. In: Proc. IPMI: international conference on information processing and medical imaging, LNCS, vol. 2732; 2003. p. 330–41.
- Good PI. *Permutation tests: a practical guide to resampling methods for testing hypotheses*. 2nd ed. New York: Springer; 2000.
- Hamadicharef B. AUC confidence bounds for performance evaluation of Brain-Computer Interface. In: IEEE 3rd International (Volume:5) Conference on Biomedical Engineering and Informatics (BMEI); 2010. p. 1988–91. Available from: http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5639671.
- Hamamé CM, Vidal JR, Ossandón T, Jerbi K, Dalal SS, Minotti L, et al. Reading the mind's eye: online detection of visuo-spatial working memory and visual imagery in the inferior temporal lobe. *NeuroImage* 2012;59(January (1)):872–9.
- Haynes J-D, Sakai K, Rees G, Gilbert S, Frith C, Passingham RE. Reading hidden intentions in the human brain. *Curr Biol* 2007;17(February (4)):323–8.
- Hill NJ, Lal TN, Schröder M, Hinterberger T, Wilhelm B, Nijboer F, et al. Classifying EEG and ECoG signals without subject training for fast BCI implementation: comparison of nonparalyzed and completely paralyzed subjects. *IEEE Trans Neural Syst Rehabil Eng* 2006;14(June (2)):183–6.
- Hosseini SMH, Mano Y, Rostami M, Takahashi M, Suguri M, Kawashima R. Decoding what one likes or dislikes from single-trial fNIRS measurements. *NeuroReport* 2011;22(April (6)):269–73.
- Huang J, Ling CX. Using AUC and accuracy in evaluating learning algorithms. *IEEE Trans Knowl Data Eng* 2005;17(3):299–310.
- Jerbi K, Bertrand O, Schoendorff B, Hoffmann D, Minotti L, Kahane P, et al. Online detection of gamma oscillations in ongoing intracerebral recordings: From functional mapping to brain computer interfaces. In: Noninvasive Funct Source Imaging Brain Heart Int Conf Funct Biomed Imaging 2007 NFSI-ICFBI 2007. It Meet 6th Int Symp On. IEEE; 2007a. p. 330–3. http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4387767.
- Jerbi K, Lachaux JP, Karim N, Pantazis D, Leahy RM, Garnero L, et al. Coherent neural representation of hand speed in humans revealed by MEG imaging. *Proc Natl Acad Sci* 2007b;104(18):7676–81.
- Jerbi K, Freyermuth S, Minotti L, Kahane P, Berthoz A, Lachaux J. Watching brain TV and playing brain ball International review of neurobiology. In: Brain machine interfaces for space applications: enhancing astronaut capabilities. San Diego: Elsevier Academic Press; 2009a. p. 159–68 [chapter 12]. <http://linkinghub.elsevier.com/retrieve/pii/S0074774209860121>.
- Jerbi K, Ossandón T, Hamamé CM, Senova S, Dalal SS, Jung J, et al. Task-related gamma-band dynamics from an intracerebral perspective: review and implications for surface EEG and MEG. *Hum Brain Mapp* 2009b;30(June (6)):1758–71.
- Jerbi K, Vidal JR, Mattout J, Maby E, Lecaillard F, Ossandón T, et al. Inferring hand movement kinematics from MEG, EEG and intracranial EEG: from brain-machine interfaces to motor rehabilitation. *IRBM* 2011;32(February (1)):8–18.
- Jerbi K, Combrisson E, Dalal SS, Vidal JR, Hamamé CM, Bertrand O, et al. Decoding cognitive states and motor intentions from intracranial EEG: how promising is high-frequency brain activity for brain-machine interfaces? *Epilepsy Behav* 2013;28(2):283–302.
- Kayikcioglu T, Aydemir O. A polynomial fitting and k-NN based approach for improving classification of motor imagery BCI data. *Pattern Recognit Lett* 2010;31(August (11)):1207–15.
- Kellis S, Miller K, Thomson K, Brown R, House P, Greger B. Decoding spoken words using local field potentials recorded from the cortical surface. *J Neural Eng* 2010;7(October (5)):056007.
- Kohavi R. A study of cross-validation and bootstrap for accuracy estimation and model selection; 1995. p. 1137–45. <http://frostiebek.free.fr/docs/Machine%20Learning/validation-1.pdf>.
- Korcyn AD, Schachter SC, Brodie MJ, Dalal SS, Engel J, Guekht A, et al. Epilepsy, cognition, and neuropsychiatry (Epilepsy, Brain, and Mind, part 2). *Epilepsy Behav* 2013;28(2):283–302.
- Krusienski DJ, Wolpaw JR. Brain-computer interface research at the wadsworth center developments in noninvasive communication and control. *Int Rev Neurobiol* 2009;86:147–57.
- Lachaux JP, Jerbi K, Bertrand O, Minotti L, Hoffmann D, Schoendorff B, et al. A Blueprint for Real-Time Functional Mapping via Human Intracranial Recordings. *PLoS ONE* 2007a;2(October (10)):e1094.
- Lachaux JP, Jerbi K, Bertrand O, Minotti L, Hoffmann D, Schoendorff B, et al. BrainTV: a novel approach for online mapping of human brain functions. *Biol Res* 2007b;40(January (4)):401–13.
- Lampe T, Fiederer LDJ, Voelker M, Knorr A, Riedmiller M, Ball T. A brain-computer interface for high-level remote control of an autonomous, reinforcement-learning-based robotic system for reaching and grasping. In: Proceedings of the 19th international conference on intelligent user interfaces. New York, NY, USA: ACM; 2014. p. 83–8. <http://dx.doi.org/10.1145/2557500.2557533>.
- Leuthardt EC, Schalk G, Wolpaw JR, Ojemann JG, Moran DW. A brain-computer interface using electrocorticographic signals in humans. *J Neural Eng* 2004;1(June (2)):63.
- Leuthardt EC, Miller KJ, Schalk G, Rao RPN, Ojemann JG. Electrocorticography-based brain computer interface—the Seattle experience. *IEEE Trans Neural Syst Rehabil Eng* 2006;14(June (2)):194–8.

- Ling CX, Huang J, Zhang H. AUC: a statistically consistent and more discriminative measure than accuracy; 2003. p. 519–24. <http://arion.csd.uwo.ca/faculty/ling/papers/ijcai03.pdf>.
- Lotte F, Congedo M, Lécuyer A, Lamarche F, Arnaldi B. A review of classification algorithms for EEG-based brain-computer interfaces. *J Neural Eng [Internet]* 2007 [cited 2012 Oct 3];4. Available from: <http://hal.archives-ouvertes.fr/docs/00/13/49/50/PDF/article.pdf>.
- Martin JK, Hirschberg DS. Small sample statistics for classification error rates I: error rate measurements. Irvine: Information and Computer Science, University of California; 1996a. <http://www.ics.uci.edu/~dan/pubs/TR96-21.pdf>.
- Martin JK, Hirschberg DS. Small sample statistics for classification error rates II: confidence intervals and significance tests [Internet]. Information and Computer Science. Irvine: University of California; 1996b. Disponible sur: <http://www.ics.uci.edu/~dan/pubs/TR96-22.pdf>.
- Mehring C, Nawrot MP, de Oliveira SC, Vaadia E, Schulze-Bonhage A, Aertsen A, et al. Comparing information about arm movement direction in single channels of local and epicortical field potentials from monkey and human motor cortex. *J Physiol – Paris* 2004;98(July (4–6)):498–506.
- Meyers EM, Kreiman G. Tutorial on pattern classification in cell recordings. In: Kriegeskorte N, Kreiman G, editors. Understanding visual population codes. Boston: MIT Press; 2011.
- Momennejad I, Haynes J-D. Human anterior prefrontal cortex encodes the “what” and “when” of future intentions. *Neuroimage* 2012;61(May (1)):139–48.
- Morash V, Bai O, Furlani S, Lin P, Hallett M. Classifying EEG signals preceding right-hand, left hand, tongue, and right foot movements and motor imaginations. *Clin Neurophysiol* 2008;119(November (11)):2570–8.
- Müller-Putz GR, Scherer R, Brunner C, Leeb R, Pfurtscheller G. Better than random? A closer look on BCI results. *Int J Bioelectromagn* 2008;10(1):52–5.
- Neuper C, Scherer R, Reiner M, Pfurtscheller G. Imagery of motor actions: differential effects of kinesthetic and visual-motor mode of imagery in single-trial EEG. *Brain Res Cogn Brain Res* 2005;25(December (3)):668–77.
- Nichols TE, Holmes AP. Nonparametric permutation tests for functional neuroimaging: a primer with examples. *Hum Brain Mapp* 2002;15(1):1–25.
- Noirhomme Q, Lesenfants D, Gomez F, Soddu A, Schrouff J, Garraux G, et al. Biased binomial assessment of cross-validated estimation of classification accuracies illustrated in diagnosis predictions. *NeuroImage: Clin* 2014;4:687–94.
- Ojala M, Garriga GC. Permutation tests for studying classifier performance. *J Mach Learn Res* 2010;11(June):1833–63.
- Pistohl T, Schulze-Bonhage A, Aertsen A, Mehring C, Ball T. Decoding natural grasp types from human ECoG. *NeuroImage* 2012;59(January (1)):248–60.
- Quiroga RQ, Panzeri S. Extracting information from neuronal populations: information theory and decoding approaches. *Nat Rev Neurosci* 2009;10:173.
- Raudys SJ, Jain AK. Small sample size effects in statistical pattern recognition: recommendations for practitioners. *IEEE Trans Pattern Anal Mach Intell* 1991;13(March (3)):252–64.
- Schalk G, Miller KJ, Anderson NR, Wilson JA, Smyth MD, Ojemann JG, et al. Two-dimensional movement control using electrocorticographic signals in humans. *J Neural Eng* 2008;5(March (1)):75–84.
- Sitaram R, Lee S, Ruiz S, Rana M, Veit R, Birbaumer N. Real-time support vector classification and feedback of multiple emotional brain states. *NeuroImage* 2011;56(May (2)):753–65.
- Tilbury JB, Van Etetvelt WJ, Garibaldi JM, Curnsw JSH, Ifeachor EC. Receiver operating characteristic analysis for intelligent medical systems—a new approach for finding confidence intervals. *IEEE Trans Biomed Eng* 2000;47(7):952–63.
- Toppi J, Risetti M, Quigadamo LR, Petti M, Bianchi L, Salinari S, et al. Investigating the effects of a sensorimotor rhythm-based BCI training on the cortical activity elicited by mental imagery. *J Neural Eng* 2014;11(June (3)):035010.
- Vapnik V. The nature of statistical learning theory. New York: Springer Science & Business Media; 1995.
- Waldert S, Braun C, Preissl H, Birbaumer N, Aertsen A, Mehring C. Decoding performance for hand movements: EEG vs. MEG. *IEEE* 2007:5346–8.
- Waldert S, Preissl H, Demandt E, Braun C, Birbaumer N, Aertsen A, et al. Hand movement direction decoded from MEG and EEG. *J Neurosci* 2008;28(January (4)):1000–8.
- Waldert S, Tüshaus L, Kaller CP, Aertsen A, Mehring C. fNIRS exhibits weak tuning to hand movement direction. *PLoS ONE* 2012;7(November (11)):e49266.
- Wang W, Sudre GP, Xu Y, Kass RE, Collinger JL, Degenhart AD, et al. Decoding and cortical source localization for intended movement direction with MEG. *J Neurophysiol* 2010;104(November (5)):2451–61.
- Wieland M, Pittore M. Performance Evaluation of Machine Learning Algorithms for Urban Pattern Recognition from Multi-spectral Satellite Images. *Remote Sens* 2014;6(March (4)):2912–39.

5.3 COMPLÉMENTS D'ÉTUDE

compléments sur les différents types de permutation (Ojala and Garriga, 2010)

Troisième partie

Étude 2 : encodage de l'intention et de l'exécution motrice

SOMMAIRE

5.1	PRÉSENTATION DE L'ÉTUDE	31
5.1.1	Contexte	31
5.1.2	Problématique	31
5.1.3	Résultats majeurs	31
5.2	ARTICLE	31
5.3	COMPLÉMENTS D'ÉTUDE	43
5.4	RÉSUMÉ DE L'ÉTUDE	49
5.5	ARTICLE	49
	CONCLUSION	49
5.6	RÉSUMÉ DE L'ÉTUDE	55
5.7	ARTICLE	55
	CONCLUSION	55
5.8	RÉSUMÉ DE L'ÉTUDE	61
5.9	ARTICLE	61
	CONCLUSION	61
5.10	RÉSUMÉ DE L'ÉTUDE	67
5.11	ARTICLE	67
	CONCLUSION	67

Ce chapitre introductif gnagnagna.
Pas obligatoire!

5.4 RÉSUMÉ DE L'ÉTUDE

5.5 ARTICLE

CONCLUSION DU CHAPITRE

Ceci est la conclusion. Personnellement, je n'aime pas que la conclusion soit numéroté, mais je veux qu'elle apparaisse dans la table des matière, d'où la commande addcontentsline.

Quatrième partie

**Étude 3 : décodage des
directions de mouvement
pendant et avant l'exécution de
mouvement de membres
supérieurs**

SOMMAIRE

5.1	PRÉSENTATION DE L'ÉTUDE	31
5.1.1	Contexte	31
5.1.2	Problématique	31
5.1.3	Résultats majeurs	31
5.2	ARTICLE	31
5.3	COMPLÉMENTS D'ÉTUDE	43
5.4	RÉSUMÉ DE L'ÉTUDE	49
5.5	ARTICLE	49
	CONCLUSION	49
5.6	RÉSUMÉ DE L'ÉTUDE	55
5.7	ARTICLE	55
	CONCLUSION	55
5.8	RÉSUMÉ DE L'ÉTUDE	61
5.9	ARTICLE	61
	CONCLUSION	61
5.10	RÉSUMÉ DE L'ÉTUDE	67
5.11	ARTICLE	67
	CONCLUSION	67

Ce chapitre introductif gnagnagna.
Pas obligatoire!

5.6 RÉSUMÉ DE L'ÉTUDE

5.7 ARTICLE

CONCLUSION DU CHAPITRE

Ceci est la conclusion. Personnellement, je n'aime pas que la conclusion soit numéroté, mais je veux qu'elle apparaisse dans la table des matière, d'où la commande addcontentsline.

Cinquième partie

Étude 4 : optimisation des paramètres de la bande gamma

SOMMAIRE

5.1	PRÉSENTATION DE L'ÉTUDE	31
5.1.1	Contexte	31
5.1.2	Problématique	31
5.1.3	Résultats majeurs	31
5.2	ARTICLE	31
5.3	COMPLÉMENTS D'ÉTUDE	43
5.4	RÉSUMÉ DE L'ÉTUDE	49
5.5	ARTICLE	49
CONCLUSION		49
5.6	RÉSUMÉ DE L'ÉTUDE	55
5.7	ARTICLE	55
CONCLUSION		55
5.8	RÉSUMÉ DE L'ÉTUDE	61
5.9	ARTICLE	61
CONCLUSION		61
5.10	RÉSUMÉ DE L'ÉTUDE	67
5.11	ARTICLE	67
CONCLUSION		67

Ce chapitre introductif gnagnagna.
Pas obligatoire!

5.8 RÉSUMÉ DE L'ÉTUDE

5.9 ARTICLE

CONCLUSION DU CHAPITRE

Ceci est la conclusion. Personnellement, je n'aime pas que la conclusion soit numéroté, mais je veux qu'elle apparaisse dans la table des matière, d'où la commande addcontentsline.

Sixième partie

Étude 5 : décodage des émotions

SOMMAIRE

5.1	PRÉSENTATION DE L'ÉTUDE	31
5.1.1	Contexte	31
5.1.2	Problématique	31
5.1.3	Résultats majeurs	31
5.2	ARTICLE	31
5.3	COMPLÉMENTS D'ÉTUDE	43
5.4	RÉSUMÉ DE L'ÉTUDE	49
5.5	ARTICLE	49
CONCLUSION		49
5.6	RÉSUMÉ DE L'ÉTUDE	55
5.7	ARTICLE	55
CONCLUSION		55
5.8	RÉSUMÉ DE L'ÉTUDE	61
5.9	ARTICLE	61
CONCLUSION		61
5.10	RÉSUMÉ DE L'ÉTUDE	67
5.11	ARTICLE	67
CONCLUSION		67

Ce chapitre introductif gnagnagna.
Pas obligatoire!

5.10 RÉSUMÉ DE L'ÉTUDE

5.11 ARTICLE

CONCLUSION DU CHAPITRE

Ceci est la conclusion. Personnellement, je n'aime pas que la conclusion soit numéroté, mais je veux qu'elle apparaisse dans la table des matière, d'où la commande addcontentsline.

CONCLUSION GÉNÉRALE

Enfin : la conclusion générale!!!

Au cours de ce mémoire, nous avons développé un modèle ...

1. **Modélisation**

2. **Inférence statistique**

PERSPECTIVES

Dans la continuité directe de notre travail de thèse, nous pouvons ...

ANNEXES

A

SOMMAIRE

A.1	PREUVE DU THÉORÈME TRUC	73
-----	-----------------------------------	----

A.1 PREUVE DU THÉORÈME TRUC

Ce théorème est un résultat classique donné, par exemple, par...

BIBLIOGRAPHIE

- Bahramisharif, A., van Gerven, M. A. J., Aarnoutse, E. J., Mercier, M. R., Schwartz, T. H., Foxe, J. J., Ramsey, N. F., and Jensen, O. (2013). Propagating Neocortical Gamma Bursts Are Coordinated by Traveling Alpha Waves. *Journal of Neuroscience*, 33(48) :18849–18854.
- Bekaert, M. H., Botte-Lecocq, C., Cabestaing, F., Rakotomamonjy, A., et al. (2009). Les interfaces Cerveau-Machine pour la palliation du handicap moteur sévère. *Sciences et Technologies pour le Handicap*, 3(1) :95–121.
- Berens, P. and others (2009). CircStat a MATLAB toolbox for circular statistics. *J Stat Softw*, 31(10) :1–21.
- Cortes, C. and Vapnik, V. (1995). Support-vector networks. *Machine learning*, 20(3) :273–297.
- de Hemptinne, C., Ryapolova-Webb, E. S., Air, E. L., Garcia, P. A., Miller, K. J., Ojemann, J. G., Ostrem, J. L., Galifianakis, N. B., and Starr, P. A. (2013). Exaggerated phase–amplitude coupling in the primary motor cortex in Parkinson disease. *Proceedings of the National Academy of Sciences*.
- Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2) :179–188.
- Hanakawa, T., Dimyan, M. A., and Hallett, M. (2008). Motor Planning, Imagery, and Execution in the Distributed Motor Network : A Time-Course Study with Functional MRI. *Cerebral Cortex*, 18(12) :2775–2788.
- Hochberg, L. R., Bacher, D., Jarosiewicz, B., Masse, N. Y., Simeral, J. D., Vogel, J., Haddadin, S., Liu, J., Cash, S. S., van der Smagt, P., and Donoghue, J. P. (2012). Reach and grasp by people with tetraplegia using a neurally controlled robotic arm. *Nature*, 485(7398) :372–375.
- Hyafil, A., Giraud, A.-L., Fontolan, L., and Gutkin, B. (2015). Neural Cross-Frequency Coupling : Connecting Architectures, Mechanisms, and Functions. *Trends in Neurosciences*, 38(11) :725–740.
- Lakatos, P. (2005). An Oscillatory Hierarchy Controlling Neuronal Excitability and Stimulus Processing in the Auditory Cortex. *Journal of Neurophysiology*, 94(3) :1904–1911.
- Lotte, F., Congedo, M., Lécuyer, A., Lamarche, F., Arnaldi, B., et al. (2007). A review of classification algorithms for EEG-based brain–computer interfaces. *Journal of neural engineering*, 4.

- Ojala, M. and Garriga, G. C. (2010). Permutation tests for studying classifier performance. *The Journal of Machine Learning Research*, 11 :1833–1863.
- Ozkurt, T. E. (2012). Statistically Reliable and Fast Direct Estimation of Phase-Amplitude Cross-Frequency Coupling. *Biomedical Engineering, IEEE Transactions on*, 59(7) :1943–1950.
- Tallon-Baudry, C., Bertrand, O., Delpuech, C., and Pernier, J. (1997). Oscillatory γ -band (30–70 Hz) activity induced by a visual search task in humans. *The Journal of neuroscience*, 17(2) :722–734.
- Tort, A. B. L., Komorowski, R., Eichenbaum, H., and Kopell, N. (2010). Measuring Phase-Amplitude Coupling Between Neuronal Oscillations of Different Frequencies. *Journal of Neurophysiology*, 104(2) :1195–1210.
- Van Langhenhove, A., Bekaert, M. H., Cabestaing, F., N'Guyen, J. P., et al. (2008). Interfaces cerveau-ordinateur et rééducation fonctionnelle : étude de cas chez un patient hémiparésique. *Sciences et Technologies pour le Handicap*, 2(1) :41–54.
- Vapnik, V. (2000). *The Nature of Statistical Learning Theory*. Springer Science & Business Media.
- Voytek, B., D'Esposito, M., Crone, N., and Knight, R. T. (2013). A method for event-related phase/amplitude coupling. *NeuroImage*, 64 :416–424.
- Waldert, S., Pistohl, T., Braun, C., Ball, T., Aertsen, A., and Mehring, C. (2009). A review on directional information in neural signals for brain-machine interfaces. *Journal of Physiology-Paris*, 103(3-5) :244–254.
- Worrell, G., Jerbi, K., Kobayashi, K., Lina, J., Zelmann, R., and Le Van Quyen, M. (2012). Recording and analysis techniques for high-frequency oscillations. *Progress in neurobiology*, 98(3) :265–278.

Titre Décodage des intentions et des représentations motrices chez l'homme : analyse multi-échelle et application aux interfaces cerveau-machine

Résumé Le résumé en français (\approx 1000 caractères)

Mots-clés Les mots-clés en français

Title Le titre en anglais

Abstract Le résumé en anglais (\approx 1000 caractères)

Keywords Les mots-clés en anglais