

Université Lumière Lyon2

École Doctorale “Cerveau et Cognition”

Laboratoire d’InfoRmatique en Images et Systèmes d’information

LIRIS - UMR CNRS 5205 - Équipe TURING

UNE MODÉLISATION ÉVOLUTIONNISTE DU LIAGE TEMPOREL

Par David MEUNIER

Thèse de doctorat de **Sciences Cognitives**
Mention Informatique

Dirigée par Professeur Hélène PAUGAM-MOISY

Présenté et soutenue publiquement le **19 Octobre 2007**

Devant un jury composé de :

Docteur Hugues BERRY (Examineur)

Docteur Olivier BERTRAND (Examineur)

Docteur Agnès GUILLOT (Examineur)

Professeur Hélène PAUGAM-MOISY (Directeur de thèse)

Professeur Marc SCHOENAUER (Rapporteur)

Docteur Catherine TALLON-BAUDRY (Rapporteur)

Table des matières

1	Introduction	1
1.1	Problématique du liage des traits	1
1.2	Modèle de liage temporel	2
1.3	Emergence par l'évolution	2
1.4	Théorie des réseaux complexes	3
1.5	Plan de la thèse	3
2	Neurophysiologie	4
2.1	Fonctionnement des neurones biologiques	4
2.1.1	Propriétés des neurones biologiques	4
2.1.2	Codage de l'information par les neurones	5
2.2	Plasticité synaptique	8
2.2.1	Assemblées cellulaires	8
2.2.2	Plasticité synaptique par taux de décharge	9
2.2.3	Plasticité synaptique temporelle	9
2.3	Neuro-imagerie et électrophysiologie	10
2.3.1	Fonction cognitive localisée	11
2.3.2	Fonctionnement "en réseau"	11
3	Hypothèse du liage temporel	13
3.1	Problématique du liage des traits	13
3.2	Liage par convergence	14
3.3	Hypothèse du liage temporel	15
3.3.1	Assemblées temporelles	15
3.3.2	Oscillations	18
3.4	Évidence expérimentale du liage temporel	21
3.4.1	Perception	21
3.4.2	Processus de groupage	22
3.4.3	Attention et influence <i>top-down</i>	24
3.4.4	Mémorisation	25
3.4.5	Intégration multi-modale	26
3.4.6	Pathologies	27
4	Neurosciences computationnelles	29
4.1	Réseaux de neurones artificiels	29
4.1.1	Réseaux de neurones artificiels "classiques"	29

4.1.2	Réseaux de neurones temporels	32
4.2	Modèles de synchronisation neuronale	36
4.2.1	Synchronisation d'oscillateurs	37
4.2.2	Modèle des <i>synfire chains</i>	37
4.2.3	Polychronisation	38
4.2.4	Emergence d'oscillations dans une structure particulière	39
4.2.5	Modèles de liage des traits	42
4.2.6	Modèle <i>a priori</i> versus <i>a posteriori</i>	43
5	Evolution	45
5.1	Evolution biologique	45
5.1.1	Théorie synthétique de l'évolution	45
5.1.2	Effet Baldwin	47
5.2	Algorithme évolutionniste	48
5.2.1	Principes de fonctionnement	48
5.2.2	Explications du fonctionnement	50
5.2.3	Applications	50
5.3	Evolution et cerveau	51
5.3.1	Evolution et apprentissage	51
5.3.2	Simulation de l'évolution du cerveau	53
5.3.3	Applications en robotique	57
5.3.4	Retour à la problématique	60
6	Théorie des réseaux complexes	61
6.1	Introduction	61
6.1.1	Systèmes complexes	61
6.1.2	Réseaux complexes	61
6.1.3	Applications de la théorie des réseaux complexes	62
6.2	Modèles et mesures pour les réseaux complexes	63
6.2.1	Composantes fortement connexes	63
6.2.2	Modèle de réseaux petit-monde et mesures d'efficacité	64
6.2.3	Réseaux invariants d'échelle	68
6.2.4	Modularité	68
6.3	Dynamique dans les modèles de réseaux complexes	73
6.3.1	Dynamique dans les réseaux petit-monde	74
6.3.2	Dynamique dans les réseaux invariants d'échelle	74
7	Réseaux complexes et cerveau	76
7.1	Modularité du cerveau	77
7.1.1	Module Fodorien	77
7.1.2	Module Darwinien	77
7.2	Réseaux complexes et neuro-imagerie	78
7.2.1	Niveau anatomique	78
7.2.2	Niveau dynamique	79
7.3	Réseaux complexes et Neurosciences computationnelles	81
7.4	Réseaux complexes et évolution	83

8	Le modèle EvoSNN	85
8.1	Modèle de neurone	85
8.2	Modèle de synapse	87
8.2.1	STDP	87
8.2.2	Règles additive ou multiplicative	87
8.2.3	Application de la STDP	88
8.3	Modèle d'évolution	88
8.3.1	Construction de la topologie du réseau	88
8.3.2	Environnement virtuel	93
8.3.3	Algorithme évolutionniste	95
9	Résultats sur l'évolution et l'apprentissage	97
9.1	Résultats sur le comportement	97
9.1.1	Résultats sur l'ensemble des simulations	97
9.1.2	Exemple d'une simulation d'évolution	99
9.1.3	Exemple du passage dans l'environnement	100
9.2	Résultats sur l'apprentissage	101
9.2.1	Protocole de test	101
9.2.2	Mise en œuvre du protocole de test	102
10	Résultats sur la topologie	105
10.1	Résultats sur la connectivité	105
10.1.1	Projections "interfaces"	105
10.1.2	Projections internes au réseau	106
10.2	Résultats sur la composante fortement connexe géante	107
10.3	Mesures "petit-monde" et efficacité	108
10.4	Résultats sur la modularité	109
10.4.1	Algorithme NG étendu aux graphes orientés	109
10.4.2	Exemple d'application de l'algorithme NG étendu	110
10.4.3	Résultats avec l'algorithme NG étendu	111
10.5	Interprétation des résultats topologiques	112
11	Résultats sur la dynamique	114
11.1	Protocole de stimulation	114
11.2	Résultats sur les bandes de fréquences	115
11.2.1	Calcul des fréquences à partir des PA	116
11.2.2	Calcul des fréquences à partir des signaux continus	117
11.3	Résultats sur les cross-corrélogrammes	120
11.3.1	Ajustements des cross-corrélogrammes	121
11.3.2	K-moyennes sur les paramètres d'ajustement	123
11.3.3	Résultats des ajustements	124
11.4	Différences entre les stimuli	128
11.5	Interprétation des résultats dynamiques	130

12 Discussion	132
12.1 Evolution	133
12.2 Apprentissage	134
12.2.1 Apprentissage dynamique	134
12.2.2 Connaissances innées	135
12.2.3 Effet Baldwin	135
12.3 Topologie	136
12.3.1 Augmentation du nombre de projections “interfaces”	136
12.3.2 Absence d’émergence de structures modulaires	137
12.3.3 Optimisation du coût de câblage	138
12.4 Dynamique	138
12.4.1 Considérations méthodologiques	139
12.4.2 Activité de fond	140
12.4.3 Injection d’un stimulus	142
12.4.4 Formation d’une assemblée temporelle et liage temporel	145
13 Conclusion	147
14 Perspectives	148
14.1 Améliorations du modèle	148
14.1.1 Inclusion d’un “coût” métabolique	148
14.1.2 Réseau petit-monde	149
14.1.3 Réseau exponentiel	149
14.2 Liens entre la topologie et la dynamique	149
14.3 Application aux études de neuro-imagerie	150
A Théorie des graphes	152
A.1 Plusieurs types de graphes	152
A.2 Représentation matricielle d’un graphe	153
A.3 Mesures sur un graphe	153
A.3.1 Degré	153
A.3.2 Plus court chemin	153
A.4 Réseaux aléatoires et réseaux de voisinage	154
A.5 Parité des circuits	154
B Méthodes d’analyse des signaux continus	156
B.1 Potentiels évoqués	156
B.2 Oscillations induites	157
B.2.1 Diagrammes Temps-Fréquence	157
B.2.2 Oscillations évoquées et induites	157
B.2.3 Synchronisation de phase	158
C Cross-corrélogrammes	159
C.1 Calcul d’un cross-corrélogramme	159
C.2 Cross-corrélogramme corrigé	160
C.3 Normalisation	160

D	Statistiques	161
D.1	ANOVA	161
D.2	χ^2 réduit	161
E	K-moyennes	162

Annexes

Annexe A

Théorie des graphes

La théorie des graphes, issue de la topologie en mathématiques, est majoritairement utilisée en informatique. En effet, aussi bien l'ordonnancement d'un ensemble de processus fonctionnant sur un ordinateur, qu'une arborescence de dossiers et de fichiers, peuvent être décrits par des graphes. On définit un graphe comme un ensemble de nœuds reliés entre eux par des liens, ces liens représentant une forme d'interaction entre les nœuds (relation père-fils entre processus, dossiers ou fichiers contenus dans un autre dossier, etc.).

A.1 Plusieurs types de graphes

On définit plusieurs types de graphes, en fonction de la nature des liens qui existent entre les nœuds. Classiquement, un graphe est **non orienté** (Figure A.1, à gauche), c.à.d. qu'on ne tient pas compte, dans sa description, d'une quelconque direction sur les liens, chaque lien représentant une relation réciproque - on parle alors d'**arête**. Une autre famille de graphe est celle des graphes **orientés** (Figure A.1, au centre), où la relation entre les nœuds n'est pas réciproque. La direction du lien du nœud origine vers le nœud cible, est représentée par une flèche - on parle d'**arcs**, dans ce cas. Enfin, une autre famille de graphe est celle des graphes **pondérés** (Figure A.1, à droite), où l'on attribue un poids plus ou moins forts aux liens. Un graphe pondéré peut être orienté, ou non.

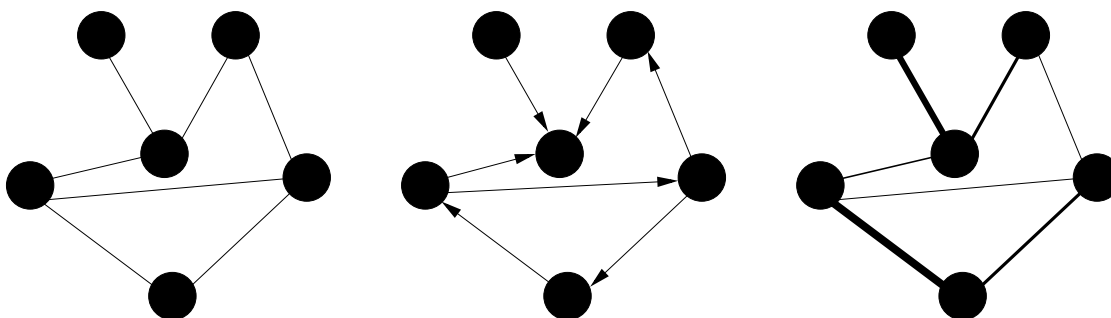


FIG. A.1 – Différents types de graphes : graphe non orienté (à gauche), graphe orienté (au centre) et graphe pondéré (à droite).

A.2 Représentation matricielle d'un graphe

Un graphe peut être représenté sous la forme d'une matrice d'adjacence A , dont les éléments A_{ij} valent 1 s'il existe un lien entre le nœud i et le nœud j , ou 0 sinon. Dans le cas d'un graphe orienté, la matrice d'adjacence n'est pas symétrique. Dans le cas d'un graphe pondéré, les éléments A_{ij} valent 0 s'il n'existe pas de lien entre le nœud i et le nœud j , ou une valeur entière traduisant le force de l'interaction entre les nœuds i et j .

La définition de la topologie d'un RNA (voir section 4.1.1) passe par sa représentation sous forme de graphe pondérée. Chaque neurone est représentée par un nœud, et chaque valeur de la matrice correspond alors à la force du lien synaptique entre deux neurones (voir section 2.1.1).

A.3 Mesures sur un graphe

A.3.1 Degré

Le degré d'un nœud est le nombre de liens dont ce nœud est une des extrémités. On peut définir le degré moyen d'un graphe $\langle k \rangle$ comme la moyenne des degrés de tous les nœuds. Dans les graphes orientés, on différencie le degré entrant (le nombre de liens qui arrivent sur le nœud) et le degré sortant (le nombre de liens qui partent du nœud). On peut également définir le degré entrant moyen d'un graphe $\langle k_{in} \rangle$ et le degré sortant moyen $\langle k_{out} \rangle$. La distribution des degrés dans le réseau permet d'avoir une idée de la répartition des connexions dans le réseau.

A.3.2 Plus court chemin

L'intérêt de la représentation d'un système sous forme de graphe est que l'on peut parcourir ce graphe. A partir d'un nœud, on se déplace vers un des nœuds auquel il est lié, puis on répète la procédure à partir du nouveau nœud atteint, etc.

Dans un graphe non orienté, il existe un **chemin** entre deux nœuds s'il existe un tel parcours, permettant d'aller d'un nœud i à un nœud j . La longueur de ce chemin se mesure au nombre de liens intermédiaires traversés pour aller de i à j . S'il existe un chemin entre le nœud i et j , on dit que le nœud j est **atteignable** depuis le nœud i . On définit un **cycle** comme un chemin dont les deux extrémités sont identiques.

Dans un graphe orienté, on parle de **chemin orienté**. L'existence d'un chemin orienté du nœud i au nœud j ne garantit pas l'existence d'un chemin orienté de j vers i , tandis que cette réciproque est toujours vraie dans un graphe non orienté. On définit un **circuit** comme un chemin orienté dont les nœuds origine et cible sont identiques.

Un graphe est dit **connexe** s'il existe un chemin entre chaque paire de nœuds. Pour un graphe orienté, un graphe est **faiblement connexe** s'il existe un chemin entre chaque paire de nœuds du graphe, **fortement connexe** s'il existe un chemin orienté entre chaque paire de nœuds du graphe.

Dans un graphe connexe, il existe un ou plusieurs chemins entre chaque paire de nœuds. On appelle plus court chemin de longueur d_{ij} entre deux nœuds i et j , le chemin dont la longueur est la plus courte. Dans un graphe non orienté, d_{ij} peut être différente de d_{ji} . Le **diamètre** D d'un graphe \mathbf{G} est la moyenne de la longueur des plus courts chemins reliant deux nœuds distincts :

$$D(\mathbf{G}) = \frac{1}{N(N-1)} \sum_{i \neq j \in \mathbf{G}} d_{ij} \quad (\text{A.1})$$

A.4 Réseaux aléatoires et réseaux de voisinage

Les premiers modèles de réseaux aléatoires ont été étudiés par Erdős et Rényi (1960), puis repris par Bollobás (1985). Dans ces modèles, on définit une probabilité (appelée **densité** du graphe) d de lier deux nœuds entre eux par une arête. Avec N le nombre de noeuds et M le nombre d'arêtes, la densité vaut $d = 2M/N(N-1)$. Elle vaut 1 s'il existe un lien entre chaque paire de noeuds du réseau. On dit que le réseau est clairsemé (*sparse*) si M est petit devant le nombre maximal de connexions possible, ou encore si $d \ll 1$.

On peut définir Un **seuil de percolation** d_{perc} sur la densité d . Si la valeur de d est inférieure à d_{perc} , le graphe est composé de plusieurs sous-graphes déconnectés les uns des autres. En revanche, si la valeur de d est supérieure à d_{perc} , le graphe est essentiellement composé une composante unique (c.à.d. un sous-graphe connexe) contenant presque tous les nœuds.

Un autre type de réseaux fréquemment étudié en physique statistique est le réseau de voisinage, ou treillis régulier (*regular lattice*). Dans ce modèle, les nœuds sont placés sur une ligne (ou un anneau, pour éviter les effets de bord). Chaque nœud est lié à ses voisins sur la ligne. Par exemple, pour un réseau de degré 2, chaque nœud est lié à ses deux voisins les plus proches ; pour un réseau de degré 4, chaque nœud est lié à ses deux voisins les plus proches, et aux deux voisins de ces voisins, etc.

L'extension de ce modèle en dimension 2 est appelée une **carte**. Les liens sont là aussi définis entre les plus proches voisins sur la carte, cette fois en tenant compte des quatre directions (haut, bas, gauche, droite).

A.5 Parité des circuits

L'introduction de signes pour les arcs d'un graphe orienté a pour conséquence de faire naître des **circuits positifs ou négatifs** dans le réseau. Un circuit est un chemin dont les nœuds origine et cible sont identiques (voir annexe A). Dans un graphe signé (+1 pour un lien excitateur, -1 pour un lien inhibiteur), la parité d'un circuit est déterminée par le produit des parités des arêtes qui le composent. Un circuit pair est appelé circuit positif, un circuit impair est appelé circuit négatif.

Thomas (1981) a émis la conjecture selon laquelle l'existence d'un circuit positif dans le graphe d'interaction d'un système génère un comportement chaotique (activité auto-

entretenu), tandis que l'existence d'un circuit négatif engendre un comportement stable (activité régulée). Plusieurs démonstrations de cette conjecture ont été proposées, dans des contextes différents (Gouzé, 1998; Soulé, 2003; Remy et Ruet, 2006).

Annexe B

Méthodes d'analyse des signaux continus

B.1 Potentiels évoqués

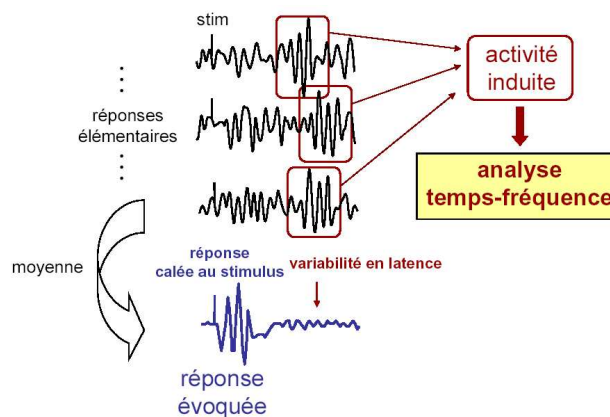


FIG. B.1 – Méthodes de calcul des potentiels évoqués. Repris de Bertrand (2006).

La méthode des potentiels évoqués permet de voir comment le potentiel mesuré par une électrode PCL, EEG (ou MEG, où le potentiel est déduit du champ magnétique mesuré) varie en fonction du temps, après la présentation d'un stimulus. Pour cela, on moyenne les signaux obtenus pour tous les essais, puis pour l'ensemble des sujets (voir figure B.1). On parle de potentiel évoqué, ou de réponse évoquée.

Les réponses que l'on observe sont répertoriées par rapport à la mesure d'une onde, positive ou négative, apparaissant avec une certaine latence après la présentation du stimulus, latence spécifique d'un traitement cognitif. On parle ainsi de P300 (onde positive apparaissant avec une latence de 300ms après l'apparition du stimulus), ou de N400 (onde négative apparaissant avec une latence de 400ms après l'apparition du stimulus). Ces ondes sont également associées à certaines structures cérébrales.

Il est possible d'étudier le spectre du signal correspondant à la moyenne de tous les essais par une analyse temps-fréquence (voir section B.2.1). On parle alors d'oscillations

évoquées (voir figure B.2, à gauche).

B.2 Oscillations induites

B.2.1 Diagrammes Temps-Fréquence

L'analyse spectrale de signaux continus est largement utilisée dans des domaines tels que l'acoustique, ou l'électronique. Une des méthodes classiques pour étudier le spectre d'un signal est la transformée de Fourier, qui permet de mesurer la contribution de chaque fréquence sur l'ensemble du signal. Cependant, cette analyse ne permet pas de mesurer des signaux dont la contribution spectrale est variable dans le temps. Une des solutions serait de réaliser une transformée de Fourier pour une partie finie du signal. Toutefois, ce traitement est dépendant de la fenêtre temporelle considérée, et le choix de cette fenêtre a une grande influence sur la représentation qu'elle engendre. Ainsi, l'analyse avec une fenêtre temporelle de courte durée ne permet pas de mesurer la contribution des fréquences correspondant à une période supérieure à la moitié de la durée considérée, selon le théorème de Shannon. Dans ce cas, la transformation la mieux adaptée est la convolution d'ondelettes, qui permet d'obtenir une représentation temps-fréquence présentant moins d'effets de bord.

Les ondelettes de Morlet $w(t, f_0)$ sont caractérisées par une forme gaussienne dans le domaine temporel (d'écart-type σ_t) et dans le domaine fréquentiel (d'écart-type σ_f), autour d'une fréquence centrale f_0 :

$$w(t, f_0) = Ae^{(-t^2/\sigma_t^2)}e^{(2i\pi f_0 t)}, \quad (\text{B.1})$$

avec $\sigma_f = 1/(2\pi\sigma_t)$. Les ondelettes sont normalisées, de manière à avoir une énergie totale de 1, impliquant que le facteur de normalisation soit $A = (\sigma_t\pi)^{-1/2}$. Une famille d'ondelettes est caractérisée par un rapport constant f_0/σ_f qui, en pratique, doit être supérieur à 5 (Tallon-Baudry et al., 1997a). Les ondelettes utilisées dans la partie 11.2.2 ont un rapport f_0/σ_f de 7, avec f_0 allant de 1 à 100 Hz par pas de 1 Hz.

L'énergie $E(t, f_0)$ est la norme quadratique de la convolution de l'ondelette $w(t, f_0)$ avec le signal continu s à l'instant t :

$$E(t, f_0) = |w(t, f_0) * s(t)|^2 \quad (\text{B.2})$$

B.2.2 Oscillations évoquées et induites

La méthode des potentiels évoqués (annexe B.1) permet de ne garder que la réponse qui se reproduit toujours avec la même latence vis-à-vis de l'apparition du stimulus. Toutes les réponses apparaissant à une latence variable d'un essai sur l'autre sont supprimées lors du moyennage des signaux (voir figure B.1).

Tallon-Baudry et al. (1997b) développe une méthode permettant de mesurer des phénomènes dont la latence vis-à-vis de l'instant d'apparition du stimulus est variable d'un essai sur l'autre. Le calcul de la convolution d'ondelettes à partir d'un signal étant une opération non linéaire, la somme des diagrammes temps-fréquence de chaque essai ne donne pas le même résultat que le diagramme temps-fréquence du potentiel évoqué

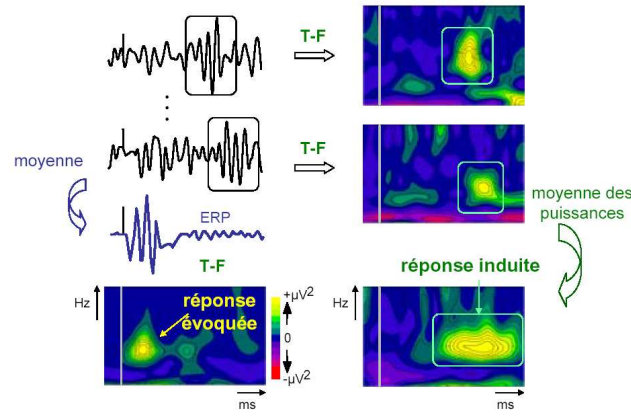


FIG. B.2 – Différences entre les méthodes de calcul d’oscillations évoquées et induites. Repris de Bertrand (2006).

(i.e. la somme des signaux). En plus de la contribution du potentiel évoqué, une seconde contribution apparaît sur le diagramme obtenu par la méthode de Tallon-Baudry et al. (1997b). On parle alors d’oscillation induites, ou de réponse induite, par opposition à la réponse évoquée. Cette réponse induite, dont le caractère variable est supposé refléter la variabilité de l’état attentionnel du sujet à chaque essai, serait plus pertinente pour mesurer le traitement cognitif réalisé par le sujet, alors que la réponse évoquée serait le reflet de l’impact physique du stimulus dans le système nerveux du sujet.

B.2.3 Synchronisation de phase

La recherche de synchronisations de phase (*Phase-Locking Factor*) dans les signaux EEG ou MEG permet de mesurer une correspondance de phase entre deux signaux continus en fonction d’une fréquence, ou d’une bande de fréquences Lachaux et al. (1999). Pour chaque essai, on convolue le signal avec une exponentielle complexe à une fréquence donnée, et on en déduit une phase complexe pour chaque fréquence. On peut alors calculer une différence de phase entre deux signaux, correspondant à deux électrodes, pour une fréquence donnée.

La technique se base sur la somme complexe de vecteurs issus de différents essais. Le vecteur somme résultant est appelé “facteur de calage en phase”. Ce facteur permet de déterminer si deux signaux présentent une différence de phase constante au cours des différents essais. S’il existe un lien de phase reproductible entre les deux signaux, le vecteur somme résultant est non nul. Il est nul si le lien de phase varie arbitrairement d’un essai sur l’autre. L’intérêt de cette mesure est de pouvoir mesurer des cohérences entre signaux, indépendamment de leurs amplitudes.

Annexe C

Cross-corrélogrammes

Une des représentations les plus utilisées pour déterminer l'influence d'un train de PA sur un autre est le cross-corrélogramme (CC), ou diagramme de cross-corrélation (pour des exemples de CC, voir la section 11.3.1). Cette représentation permet de déterminer avec quelle latence les PA émis par un neurone (le neurone de référence) sont en général suivis, ou précédés par les PA d'un autre neurone (le neurone cible).

L'existence d'un pic central sur le CC indique généralement une dépendance directe entre les trains de PA des deux neurones, le décalage de ce pic vers les valeurs positives ou négatives indiquant respectivement que les PA du neurone de référence sont toujours en avance, ou en retard sur les PA du neurone cible. Une autre forme d'information est l'existence d'harmoniques périodiques, en plus du pic central, qui apparaissent sur le CC si les PA des deux neurones sont émis avec des latences relatives identiques pendant une certaine période de temps. L'amplitude de ces harmoniques diminue en proportion inverse à la longueur de la période de temps pendant laquelle les neurones émettent avec des latences relatives identiques.

C.1 Calcul d'un cross-corrélogramme

Le cross-corrélogramme entre les trains de PA de deux neurones nécessite la définition d'une fenêtre temporelle (sur la figure C.1, [-100ms, 100ms]). Cette fenêtre est centrée sur l'instant d'émission du premier PA du neurone de référence. Dans la fenêtre temporelle, pour chaque PA du neurone cible, on ajoute +1 dans la corbeille (*bin*) temporelle correspondant à la différence entre les deux instants d'émission de PA du neurone de référence et du neurone cible (voir figure C.1, en haut). La fenêtre temporelle est ensuite déplacée au PA suivant du neurone de référence (voir figure C.1, en bas), et on recommence la procédure.

On obtient alors le **cross-corrélogramme grossier**, qui se présente sous la forme d'un histogramme, indiquant le nombre de PA du neurone cible ayant une différence temporelle donnée avec un PA du neurone de référence. La précision des corbeilles utilisée dans la section 11.3 est de 1ms, correspondant à la précision des instants d'émissions de PA. Une précision des corbeilles de 1ms est également couramment utilisée en neurosciences.

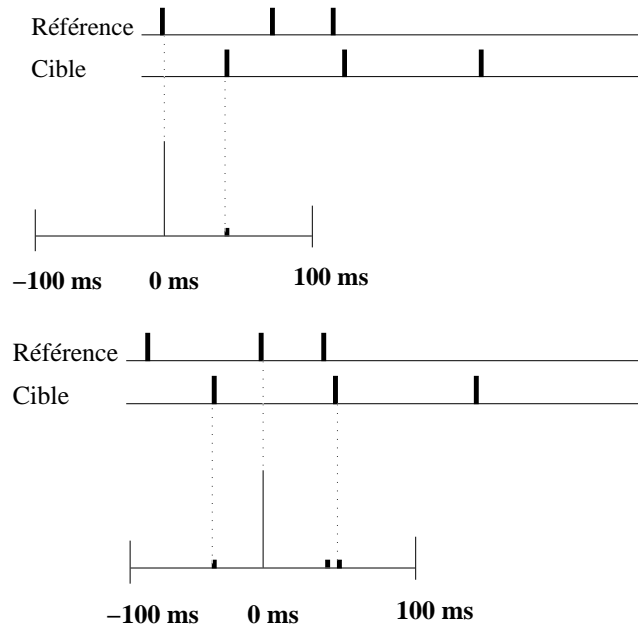


FIG. C.1 – Méthode classique de calcul d'un cross-corrélogramme (deux premières étapes).

C.2 Cross-corrélogramme corrigé

Pour ne prendre en compte que les relations temporelles entre les neurones, et non pas les variations des fréquences de décharge (par exemple une augmentation simultanée du taux de décharge des deux neurones), il faut ensuite corriger le diagramme en fonction de la prédiction du décalage d'un essai sur l'autre (*trial-based shift predictor*). Il s'agit de calculer le cross-corrélogramme entre l'essai *A* et un essai *B*, ou l'essai *B* correspond à un essai pris aléatoirement dans l'ensemble des essais, hormis l'essai *A*. Le cross-corrélogramme obtenu est ensuite soustrait au cross-corrélogramme grossier, pour obtenir le **cross-corrélogramme corrigé**. Contrairement au cross-corrélogramme grossier, le cross-corrélogramme corrigé est un histogramme qui peut prendre des valeurs négatives.

C.3 Normalisation

Afin de tenir seulement compte de la forme de la courbe, le cross-corrélogramme corrigé est ensuite normalisé pour obtenir une courbe centrée-réduite. Ce calcul s'appelle également un calcul des Z-scores. Si la moyenne des points sur l'ensemble du CC vaut M , et l'écart-type vaut σ , alors la valeur $Z(x_i)$ du point x_i vaut :

$$Z(x_i) = \frac{x_i - M}{\sigma} \quad (\text{C.1})$$

On considère classiquement qu'un point diffère de la distribution de manière significative lorsque sa valeur de Z-score, en valeur absolue, dépasse la valeur 2.

Annexe D

Statistiques

D.1 ANOVA

L'ANOVA (pour *ANalysis Of VAriance*) est un test statistique généralisant le test-t à plusieurs populations. Il permet de comparer plusieurs moyennes liées à différentes distributions, en faisant l'hypothèse que ces distributions sont gaussiennes. Il permet ainsi de déterminer si ces distributions sont significativement différentes. La preuve de cette significativité apparaît par la valeur de p , indiquant la probabilité que les deux distributions soient similaires. On considère classiquement qu'un p inférieur à 0,05, indiquant que les distributions se chevauchent à 5% au plus, permet de conclure à l'existence d'une différence significative entre les deux distributions.

On parle d'ANOVA à un facteur lorsque l'on ne considère que deux distributions. L'ANOVA multivariée permet de prendre en compte plusieurs facteurs simultanément, et ainsi de déterminer s'il existe une interaction entre les différents facteurs.

D.2 χ^2 réduit

Considérons une fonction f (fonctions sinusoïde ou ondelette dans notre cas), qui doit être ajustée à un ensemble de N points $\{M_i\}_{i=1\dots N}$, M_i de coordonnées (x_i, y_i) . Idéalement, si la fonction obtenue s'ajuste parfaitement avec les données, on a : $\forall i \in \{1 \dots N\}, f(x_i) = y_i$ pour les N points. La valeur du χ^2 est calculée par l'équation suivante :

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - f(x_i))^2}{\sigma_i^2} \quad (\text{D.1})$$

où σ_i^2 est la variance liée aux mesures d'erreur entre y_i et $f(x_i)$ pour les N points.

Le χ^2 réduit correspond à χ^2/ν , avec ν le degré de liberté : $\nu = N - p - 1$, où p est le nombre de paramètres à ajuster (dans le chapitre 11.3, $p = 3$ pour les ajustement avec une sinusoïde, $p = 4$ pour des ajustements avec une ondelette). $\chi^2/\nu \sim s^2/\sigma^2$, s^2 la variance estimée et σ^2 la variance de la distribution réelle. Pour un ajustement optimal, $s^2 \sim \sigma^2$, donc $\chi^2/\nu \sim 1$ (Laub et Kuhl, 2005). D'autre part, on peut considérer que l'ajustement est correct (au plus de 10% d'erreur) lorsque $\chi^2 < 1,5$.

Annexe E

K-moyennes

L'algorithme des K-moyennes (*K-means*), appelé également nuées dynamiques ou méthode des centres mobiles, est une technique de recherche d'agrégats (*clusters*), rassemblant des points ayant des propriétés similaires dans leur espace de représentation (Duda et Hart, 1973). Le nombre K d'agrégats recherchés doit être spécifié avant l'exécution de l'algorithme. A l'initialisation de l'algorithme, chaque point est attribué arbitrairement à un agrégat. Les barycentres de chaque agrégat sont calculés. Ensuite, de manière itérative, un point est choisi aléatoirement, et les distances de ce point à tous les barycentres sont calculées. Ainsi, le barycentre le plus proche de ce point n'est pas nécessairement celui de l'agrégat auquel il est actuellement attribué. Le point est alors alloué à l'agrégat dont il est le plus proche du barycentre. L'algorithme se termine lorsqu'il n'y a plus de modifications dans les allocations des points aux agrégats (Bottou et Bengio, 1995).

Pour n points x_j , et K clusters, l'algorithme cherche à minimiser la fonction de coût F , définie par

$$F = \sum_{i=1}^K \sum_{x_j \in S_i} \|x_j - c_i\| \quad (\text{E.1})$$

où c_i est le barycentre du cluster S_i , $\|\cdot\|$ est une mesure de distance dans l'espace considéré (ici la distance euclidienne).