

Extraction du Backbone du Réseau de Transport Aérien : Une Analyse Comparative

Ali Yassin*, Hocine Cherifi**, Hamida Seba***, Olivier Togni*

*Laboratoire d'Informatique de Bourgogne
Univ. Bourgogne - Franche-Comté, Dijon, France

**ICB UMR 6303 CNRS - Univ. Bourgogne
Franche-Comté, Dijon, France

***Univ Lyon, UCBL, CNRS, INSA Lyon
LIRIS, UMR5205, F-69622 Villeurbanne, France
ali_yassin@etu.u-bourgogne.fr

Résumé. Le développement des outils de collecte de données à grande échelle provenant des domaines biologiques, sociaux et technologiques élargit le défi de la visualisation et du traitement des grands graphes de terrain. De nombreuses techniques d'extraction visent à réduire la taille du réseau tout en préservant son essence. Dans cette étude de cas sur le transport aérien, nous réalisons une comparaison expérimentale de sept principales méthodes statistiques. L'analyse des corrélations entre les "backbone" extrait par les différentes méthodes montre que le Marginal Likelihood Filter (MLF), le Locally Adaptive Network Sparsification Filter (LANS) et le Disparity Filter sont biaisés en faveur des arêtes à forte pondération. Nous comparons les structures extraites en utilisant quatre indicateurs : la taille de la plus grande composante, le nombre de nœuds, d'arêtes et le poids total des arêtes. Les résultats montrent que les techniques basées sur un modèle de distribution binomiale (MLF et Noise Corrected Filter) ont tendance à conserver de nombreuses arêtes. En revanche, le Filtre de Disparité, le Filtre de l'Urne de Polya, le Filtre LANS et le Global Statistical Significance Filter (GloSS) sont assez agressifs pour filtrer les arêtes. Le ECM Filter se situe entre ces deux comportements. Ces résultats peuvent guider les utilisateurs dans le choix des techniques appropriées pour leurs applications spécifiques.

1 Introduction

Au cours des dernières décennies, les réseaux sont devenus un outil précieux pour l'analyse des systèmes complexes. Ils modélisent les systèmes complexes en utilisant des nœuds pour représenter les éléments et des arêtes pour représenter leurs interactions. Les tâches analytiques courantes comprennent la détection de communautés Cherifi et al. (2019), l'identification des nœuds influents Chakraborty et al. (2016); Kumar et al. (2018); Rajeh et al. (2020, 2021), et l'investigation de la formation du réseau Orman et al. (2013). Le traitement de réseaux à grande échelle présente des défis importants. Par conséquent, plusieurs méthodes d'extraction

du backbone ont été développées pour réduire la taille du réseau tout en préservant ses caractéristiques essentielles. On peut classer ces méthodes en deux catégories principales : les approches structurelles et statistiques.

Les techniques structurelles impliquent le filtrage des arêtes ou des nœuds en fonction de propriétés topologiques spécifiques telles que la modularité du réseau, les métriques de distance et l'identification des communautés en chevauchement Rajeh et al. (2022); Ghalmane et al. (2021). Les méthodes statistiques telles que le filtre de disparité Serrano et al. (2009) évaluent la signification des arêtes par le biais de tests binaires d'hypothèse et éliminent les arêtes les moins significatives à partir de leur p valeur.

Les méthodes d'extraction du backbone du réseau se sont révélées essentielles pour accélérer l'analyse et améliorer la visualisation des réseaux de transport. Une étude précédente Dai et al. (2018) a principalement comparé des méthodes structurelles pour évaluer les caractéristiques géographiques et topologiques des backbones extraits. Elle a mis en évidence la pertinence de chaque méthode pour diverses applications de recherche sur les transports. Dans une autre étude Yassin et al. (2022), nous avons évalué l'aptitude des méthodes d'extraction statistique du à révéler la structure "Hub and Spoke" du réseau de transport aérien mondial. A ce jour, aucune étude comparative n'a exploré la relation entre les p valeurs des arêtes retenues et leur pondération ainsi que les performance des principales méthodologies statistiques dans le cadre des réseaux de transport pondérés. Notre recherche comble cette lacune.

Tout d'abord, nous examinons la relation entre les poids des arêtes et les p valeurs de p obtenues par chaque technique d'extraction du backbone. Ensuite, nous extrayons les backbones pour différents niveaux de signification et les comparons en utilisant quatre indicateurs : la taille des composantes, le nombre de nœuds, d'arêtes et le poids total. Les principales contributions de l'article sont résumées comme suit :

- Nous réalisons une analyse comparative de sept techniques statistiques d'extraction du backbone sur le réseau mondial de transport aérien.
- Nous explorons la relation entre les poids des arêtes et la significativité statistique des techniques d'extraction de backbone
- Nous évaluons les techniques pour différents niveaux de signification en utilisant quatre indicateurs : la taille des composantes, le nombre de nœuds, le nombre d'arêtes et le poids total.

2 Les Techniques d'extraction du Backbone

Cette section présente brièvement les techniques statistiques d'extraction du backbone évaluées. Nous invitons les lecteurs à se référer aux articles originaux pour une description détaillée de ces méthodes.

Disparity Filter Serrano et al. (2009) : Cette technique très populaire suppose que les poids normalisés des arêtes d'un nœud suivent une distribution uniforme. Les comparaisons des poids normalisés des arêtes observées avec ce modèle nul permettent de filtrer les arêtes à un niveau de signification α souhaité.

Polya Urn Filter Marcaccioli et Livan (2019) : Dans l'urne de Pólya, l'observation d'un événement rend plus probable la répétition du même événement. De manière similaire, les poids des arêtes peuvent résulter des interactions entre les nœuds. Les auteurs créent un modèle nul

pour chaque arête en se basant sur les nœuds qu'elle relie, et calculent la probabilité qu'un nœud distribue sa force (somme des poids) à travers un processus de Pólya contrôlé par un paramètre de renforcement, α .

Marginal Likelihood Filter Dianati (2016) : Contrairement aux méthodes précédentes, le Marginal Likelihood Filter prend en compte à la fois les nœuds connectés par une arête et traite les arêtes pondérées en nombres entiers comme plusieurs arêtes unitaires. Le modèle nul repose sur l'idée que chaque arête unitaire sélectionne deux nœuds de manière aléatoire, suivant une distribution binomiale. Il calcule la probabilité de tirer au moins " w " arêtes unitaires à partir de la force du réseau, avec la probabilité liée aux forces des deux nœuds.

Noise Corrected Filter Coscia et Neffke (2017) : Similaire au Marginal Likelihood Filter, il suppose que les poids des arêtes sont tirés d'une distribution binomiale. Cependant, il estime la probabilité d'observer un poids reliant deux nœuds en utilisant un cadre bayésien. Ce cadre nous permet de générer des variances a posteriori pour toutes les arêtes. Cette variance a posteriori nous permet de créer un intervalle de confiance pour chaque poids d'arête. En fin de compte, nous retirons une arête si son poids est inférieur à δ écarts-types plus faible que l'espérance, où δ est le seul paramètre de l'algorithme. Il fournit également une approximation directe par le biais d'une distribution binomiale similaire au Marginal Likelihood Filter.

Global Statistical Significance Filter (GloSS) Radicchi et al. (2011) : Le filtre GloSS suppose qu'on ne peut pas évaluer les arêtes de manière indépendante de la topologie globale du réseau. Par conséquent, il définit un modèle nul global pour évaluer la signification d'une arête. Cependant, il ne fait aucune hypothèse concernant la distribution des poids. À la place, il utilise la distribution empirique. Le modèle est un réseau ayant la même topologie que le réseau d'origine, et les poids des arêtes sont tirés au hasard à partir de la distribution empirique des poids. En d'autres termes, il estime la probabilité d'observer un poids d'arête entre deux nœuds donnés en considérant les degrés et les forces observés des nœuds comme contraintes.

Locally Adaptive Network Sparsification Filter (LANS) Foti et al. (2011) : Il ne fait aucune hypothèse sur la distribution des poids sous-jacents. Au lieu de cela, il utilise la fonction de densité cumulative empirique pour évaluer la signification statistique. Ainsi, du point de vue des nœuds incidents de l'arête, il calcule la probabilité de choisir une arête au hasard avec un poids égal au poids observé.

Enhanced Configuration Model Filter (ECM) Gemmetto et al. (2017) : Cette technique améliore le modèle nul du filtre Marginal Likelihood. En utilisant le Modèle de Configuration Amélioré pour la reconstruction du réseau, son modèle nul repose sur l'ensemble canonique de l'entropie maximale des réseaux pondérés ayant la même distribution de degrés et de forces que le réseau réel.

3 Présentation des données et méthodes

Cette section présente l'ensemble de données en cours d'examen et décrit la méthodologie de l'analyse comparative.

Les Données : Dans nos expériences, nous utilisons le Réseau Mondial de Transport Aérien Alves et al. (2020), où les aéroports sont représentés sous forme de nœuds, et des vols directs les relient par des arêtes. Le poids de chaque arête correspond au nombre de vols sur

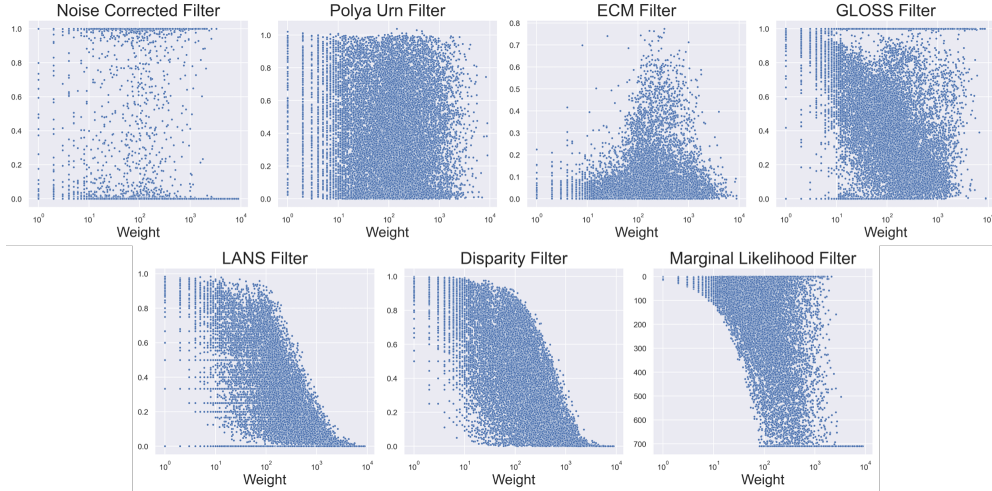
cette route. Ce réseau comprend 2734 nœuds et 16665 arêtes, avec en moyenne 12 connexions par nœud. La densité du réseau est d'environ 0,004, et il a un diamètre de 12.

Les Méthodes : Nous avons extrait les backbones à l'aide du package netbone Yassin et al. (2023). Dans un premier temps, nous examinons la relation entre les p valeurs obtenues à partir des techniques d'extraction du backbone et les poids des arêtes. L'objectif est d'évaluer si les techniques ne sont pas biaisées envers les poids forts. En effet, accorder la même importance aux arêtes indépendamment de leur poids permet de préserver les hiérarchies à toutes les échelles de poids et de fournir une meilleure représentation du réseau. Nous utilisons la corrélation de Pearson pour quantifier cette relation. Les valeurs de la corrélation de Pearson se situent dans la plage $[-1, +1]$. Plus la valeur absolue du coefficient de corrélation est élevée, plus la relation est forte. Les valeurs extrêmes indiquent une relation linéaire parfaite. Une valeur de 0 implique l'absence de corrélation entre les variables.

Dans un second temps, nous extrayons le backbone pour différents niveaux de signification. Pour évaluer les performances des techniques d'extraction nous comparons le nombre de nœuds, le nombre d'arêtes et le poids total des backbones extrait pour un même niveau de signification. Idéalement les backbones extraits préservent le maximum d'information possible dans une composante géante.

4 Corrélation entre les poids et la p valeur des Techniques d'Extraction du Backbone

Avant d'examiner les backbones extraits et leurs propriétés, nous examinons les p valeurs calculées par chaque technique d'extraction du backbone. Les techniques d'extraction du backbone abordent la limitation de l'approche de filtrage naïve, qui fixe un seuil sur les poids. Ainsi, une bonne technique d'extraction du backbone devrait accorder une importance égale aux arêtes de différentes échelles de poids. La Figure 1 montre un nuage de points pour les poids et les valeurs de p correspondant aux techniques d'extraction du backbone. Dans la rangée supérieure, on peut voir que le Noise Corrected Filter, le Polya Urn Filter, le ECM Filter et le GloSS Filter accordent de l'importance aux poids faibles et aux poids forts. En effet, ils attribuent de petites et grandes valeurs de p pour toutes les échelles de poids. En revanche, la rangée du bas montre une distribution plus étroite des valeurs de p aux extrêmes. En effet, poids faibles et les poids forts aux extrémités sont associés à des valeurs de p élevées et faibles, respectivement. Ce qui signifie que ces techniques ont plus tendance à préserver des poids fort. Nous calculons la corrélation de Pearson entre les poids et les valeurs de p de chaque technique d'extraction du backbone pour valider ces constatations. On peut distinguer la aussi deux catégories : 1) Nous ne constatons aucune corrélation entre les poids et les p valeurs du Noise Corrected Filter (0.01), du Polya Urn Filter (0.01), du ECM Filter (0.08) et du GloSS Filter (0.13). On peut donc en déduire que ces techniques ne sont pas biaisées envers les arêtes de poids fort. Elle préservent plus ou moins une large échelle de poids. 2) Nous remarquons une faible corrélation ($0,3 < \rho < 0,5$) entre les poids et les p valeurs du Marginal Likelihood Filter (0.51), du Disparity Filter (0.48) et du LANS Filter (0.37). Ces techniques privilégie plus ou moins les arêtes de poids fort.

FIG. 1 – Les poids par rapport aux valeurs p des bords du squelette pour chaque méthode.

5 Comparaison des Techniques d'Extraction du Backbone pour Différents Niveaux de Signification

Idéalement, une technique de filtrage robuste devrait maximiser la préservation de l'information tout en filtrant autant de connexions que possible, tout en veillant à ce que le réseau reste intact. Pour évaluer les méthodes d'extraction du backbone, nous considérons quatre indicateurs clés : la taille de la composante, le nombre de nœuds, le nombre d'arêtes et le poids total. Ces propriétés sont calculées en variant les niveaux de signification pour l'extraction du backbone.

La Figure 2 présente les résultats de cette expérience. Nous observons un régime de filtrage robuste, correspondant à $\alpha < 10^{-2}$. Dans le panneau en haut à gauche, nous pouvons distinguer trois comportements typiques concernant la proportion d'arêtes préservées. Tout d'abord, il y a des backbones qui maintiennent un pourcentage d'arêtes constant sur l'ensemble de la plage de niveaux de signification, même dans le régime de filtrage fort. Des exemples de cette catégorie incluent le Marginal Likelihood Filter et le Noise Corrected Filter, qui préservent environ 85% des arêtes. La deuxième catégorie comprend le ECM Filter, où le pourcentage d'arêtes augmente progressivement avec les niveaux de signification dans le régime de filtrage fort, et plus rapidement par la suite. La troisième catégorie se caractérise par un pourcentage d'arêtes constant dans le régime de filtrage fort, suivi d'une augmentation exponentielle. Les filtres comme le Poly Urn Filter, le Disparity Filter, le GloSS Filter et le LANS Filter entrent dans cette catégorie. Cependant, pendant le régime de filtrage fort, le LANS Filter maintient un pourcentage d'arêtes relativement constant (environ 18%), tandis que les autres en conservent très peu.

En passant au panneau en haut à droite, nous examinons l'évolution des poids. Les courbes montrent un schéma similaire, et nous pouvons les catégoriser comme dans le cas précédent,

à l'exception du Disparity Filter, qui rejoint le ECM Filter dans la deuxième catégorie. Dans le régime fort ($\alpha < 10^{-2}$), le Disparity Filter préserve une fraction de poids similaire au ECM Filter, même s'il conserve un petit nombre d'arêtes. Pendant ce temps, le ECM Filter conserve un plus grand nombre d'arêtes, atteignant 40% pour $\alpha = 10^{-2}$, ce qui indique que le Disparity Filter priorise les poids élevés, tandis que le ECM Filter conserve des poids de diverses échelles. Pour $\alpha > 10^{-2}$, le Disparity Filter préserve un pourcentage similaire de poids au LANS Filter, confirmant son accent sur les poids élevés.

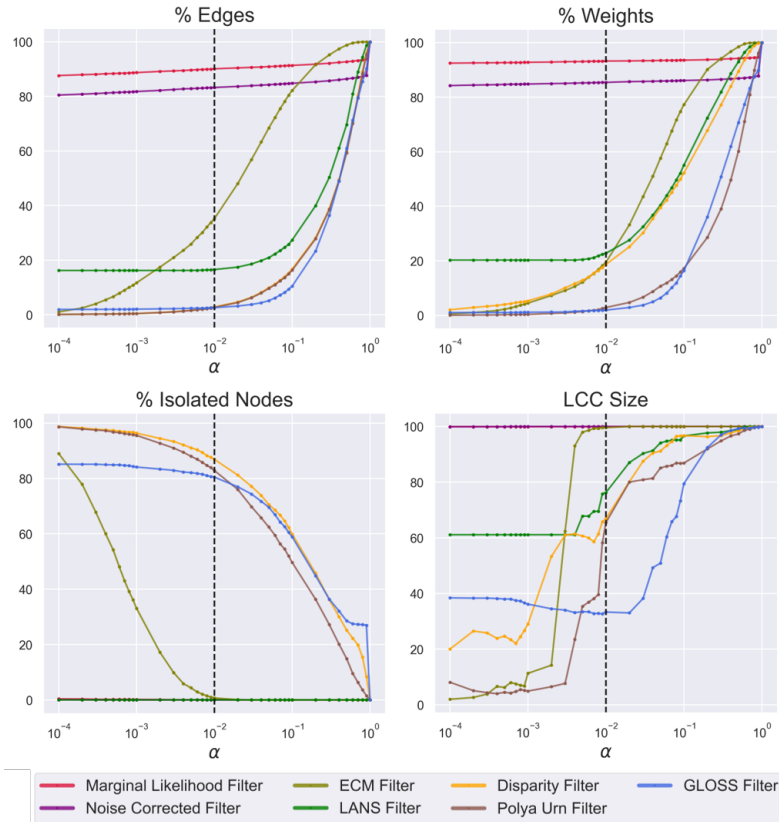
Dans le panneau inférieur gauche, nous observons l'évolution de la fraction de nœuds isolés. Le Marginal Likelihood Filter et le Noise Corrected Filter n'isolent aucun nœud car ils filtrent à peine des arêtes. Le LANS Filter préserve également tous les nœuds même en filtrant environ 80% des arêtes. En revanche, à l'exception du ECM Filter, les autres méthodes isolent un nombre substantiel de nœuds dans le régime fort. Le pourcentage de nœuds isolés diminue jusqu'à ce qu'il ne reste plus de nœuds isolés lorsque nous atteignons $\alpha = 10^{-2}$. Après ce point, le pourcentage de nœuds isolés diminue à mesure que le niveau de signification augmente, ce qui illustre comment les techniques ont du mal à conserver tous les nœuds tout en filtrant les arêtes.

Dans le dernier panneau, nous notons que, contrairement au Marginal Likelihood Filter et au Noise Corrected Filter, les autres filtres ne maintiennent pas une seule composante géante. Le LANS Filter et le GloSS Filter conservent une composante géante de taille fixe dans le régime fort. La taille de cette composante augmente progressivement jusqu'à ce qu'elle forme une seule composante seulement lorsque toutes les arêtes sont ajoutées. En revanche, nous observons l'émergence d'une composante géante dans le Disparity Filter, le Polya Urn Filter et le ECM Filter à mesure que nous approchons des limites du régime fort. Après ce point, le ECM Filter cesse d'isoler les nœuds, tandis que les autres ne s'arrêtent que lorsque toutes les arêtes sont conservées.

6 Conclusion

Réduire la taille d'un réseau tout en préservant ses propriétés est un enjeu essentiel pour de ses nombreuses applications. Ce travail compare les techniques statistiques d'extraction du backbone dans le réseau mondial de transport aérien. Les résultats montrent que le Filtre de Vraisemblance Marginale, le Filtre de Disparité et le Filtre LANS accordent plus d'importance aux arêtes à forte pondération. Les autres techniques mettent l'accent à la fois sur les arêtes de faible et de forte pondération. Nous montrons que les filtres basés sur une distribution binomiale sont des techniques conservatrices. En effet, ils conservent toujours une grande proportion de liens et de nœuds. D'autres filtres, à l'exception du Filtre ECM, sont agressifs pour supprimer les arêtes pour des niveaux de signification raisonnables $10^{-2} \leq \alpha \leq 0,05$. Ce comportement se traduit par une grande proportion de nœuds isolés. Le Filtre ECM est plus flexible. Il se situe entre ces deux comportements extrêmes. Les travaux futurs envisageront une enquête approfondie sur divers types de réseaux (synthétiques et réels) pour consolider ces résultats.

FIG. 2 – Le pourcentage d’arêtes, le poids, les nœuds isolés et la taille du plus grand composant connecté pour le squelette extrait en utilisant différents niveaux de signification α .



Références

- Alves, L., A. Aleta, F. Rodrigues, Y. Moreno, et L. Amaral (2020). Centrality anomalies in complex networks as a result of model over-simplification. *New Journal of Physics* 22.
- Chakraborty, D., A. Singh, et H. Cherifi (2016). Immunization strategies based on the overlapping nodes in networks with community structure. In *Int. conf. on computational social networks*, pp. 62–73. Springer.
- Cherifi, H., G. Palla, B. K. Szymanski, et X. Lu (2019). On community structure in complex networks : challenges and opportunities. *Applied Network Science* 4(1), 1–35.
- Coscia, M. et F. M. Neffke (2017). Network backboning with noisy data. pp. 425–436. IEEE.
- Dai, L., B. Derudder, et X. Liu (2018). Transport network backbone extraction : A comparison of techniques. *Journal of Transport Geography* 69, 271–281.
- Dianati, N. (2016). Unwinding the hairball graph : Pruning algorithms for weighted complex networks. *Physical Review E* 93.

- Foti, N. J., J. M. Hughes, et D. N. Rockmore (2011). Nonparametric sparsification of complex multiscale networks. *PLoS ONE* 6, e16431.
- Gemmetto, V., A. Cardillo, et D. Garlaschelli (2017). Irreducible network backbones : unbiased graph filtering via maximum entropy.
- Ghalmane, Z., C. Cherifi, H. Cherifi, et M. E. Hassouni (2021). Extracting backbones in weighted modular complex networks. *Scientific Reports* 11.
- Kumar, M., A. Singh, et H. Cherifi (2018). An efficient immunization strategy using overlapping nodes and its neighborhoods. In *Proc. of the The Web Conference*, pp. 1269–1275.
- Marcaccioli, R. et G. Livan (2019). A pólya urn approach to information filtering in complex networks. *Nature Communications* 10, 745.
- Orman, G. K., V. Labatut, et H. Cherifi (2013). Towards realistic artificial benchmark for community detection algorithms evaluation. *arXiv preprint arXiv :1308.0577*.
- Radicchi, F., J. J. Ramasco, et S. Fortunato (2011). Information filtering in complex weighted networks. *Physical Review E* 83, 046101.
- Rajeh, S., M. Savonnet, E. Leclercq, et H. Cherifi (2020). Interplay between hierarchy and centrality in complex networks. *IEEE Access* 8, 129717–129742.
- Rajeh, S., M. Savonnet, E. Leclercq, et H. Cherifi (2021). Characterizing the interactions between classical and community-aware centrality measures in complex networks. *Scientific reports* 11(1), 10088.
- Rajeh, S., M. Savonnet, E. Leclercq, et H. Cherifi (2022). Modularity-based backbone extraction in weighted complex networks.
- Serrano, M. A., M. Boguna, et A. Vespignani (2009). Extracting the multiscale backbone of complex weighted networks. *PNAS* 106, 6483–6488.
- Yassin, A., H. Cherifi, H. Seba, et O. Togni (2022). Exploring Statistical Backbone Filtering Techniques in the Air Transportation Network. Florence, Italy, pp. 1–8. IEEE.
- Yassin, A., A. Haidar, H. Cherifi, H. Seba, et O. Togni (2023). An evaluation tool for backbone extraction techniques in weighted complex networks. *Scientific Reports* 13(1), 17000.

Summary

Large-scale data collections from biological, social, and technological fields make it challenging to process large networks. Backbone techniques aim at reducing the network size while preserving its essence. We compare prominent statistical methods in the air transportation network. The correlation analysis between the various backbones shows that the Marginal Likelihood Filter (MLF), the Locally Adaptive Network Sparsification Filter (LANS) and the Disparity Filter favor high-weighted edges. Comparing the extracted structures using the size of the largest component, the number of nodes, edges and the total edge weight shows that the MLF and Noise Corrected Filter tend to preserve many edges. In contrast, the Disparity Filter, the Polya Urn Filter, the LANS Filter and the Global Statistical Significance Filter (GloSS) are quite aggressive in filtering edges. These results can guide users in choosing appropriate techniques for their specific applications.