

Building genetically inspired constraints with the Parsimony framework

11 avril 2014

Résumé

1 Les données et les contraintes

1.1 Websters et résultats univariés

Données décrites dans les papiers [2, 4]. Voir aussi le rapport de V Guillemot sur l'étude que nous avons faite pour se caler sur le papier séminale de 2009 concernant l'association avec le GE.

Les données contiennent des SNPs, de l'expression de gène et le statut clinique pour 380 (TOFIX) personnes. pour l'expression, des données considérées sont les celles de l'article qui supprime les effets confondants (statut Apo, TOFIX). Pour la regression logistique on considère le statut (Alz, Ctl). Pour la régression linéaire on considère l'expression du gène KIF1B qui ressort dans le papier comme gène s'associant avec des SNPs en univarié lorsque la maladie est avérée : l'interprétation exacte du papier [4] est très difficile !

1.2 Contraintes

Les contraintes considérées sont de plusieurs natures :

1.3 Contraintes de groupes

Rendues par du GroupLasso. On considère les informations :

- GeneOntology (group avec overlap massif car en tant qu'ontologie il s'agit de groupe organisés hiérarchiquement),
- pathway Keggs ou autre base de collection de gènes : c'est aussi du group avec overlap. Cette base est utilisée par Montana et coll. [3] et aussi par Chen et coll. [1]

On pourrait considérer la prise en compte d'une contrainte de type fused lasso (flou pour ce qui me concerne).

1.4 Contraintes de régions recombinantes

Il s'agit de rendre compte de la plus ou moins grande probabilité de recombinaison observée sur le génome. Les zones à faible probabilité pourraient être traitées avec une approche Group-TV : de la cohérence serait demandée dans de telles régions.

La description des données de probabilité de recombinaison est donnée dans

2 Travaux en cours

2.1 Approches GL

Les pénalités travaillées sont Ridge + L1 dans le cas de la regression logistique.

- Montage des contraintes repose sur inst_bioresource OK pour gene et snp tedious pour les pathways
- Mise en place d'une visu des résultats
- Mise en place d'un cadre pour le choix des hyperparametres. Map Reduce et usage de Gabriel
- passage a l'échelle pour le nombre de pathways
- Choix des poids a appliquer à chacun des pathway. ici il y a une ambiguïté explicitée dans la Figure 1

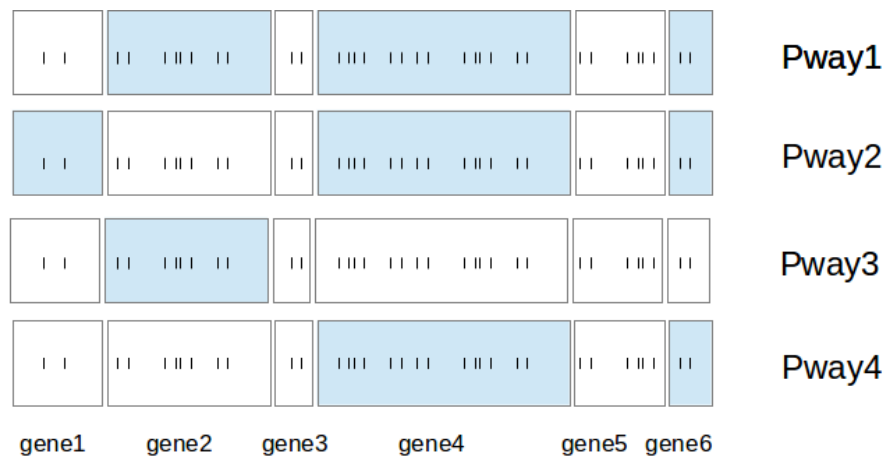


FIGURE 1 – Illustration de la difficulté des choix du poids pour un pathway donné. Les régions bleues dans un pathway sont les gènes qui sont gardés. Les SNPs sont les petits bâtons. Pour un pathway de même longueur il peut être constitué de gènes plus ou moins longs. Il semble que les pathways contenant des gènes longs sont favorisés.

Références

- [1] Lin S Chen, Carolyn M Hutter, John D Potter, Yan Liu, Ross L Prentice, Ulrike Peters, and Li Hsu. Insights into colon cancer etiology via a regularized approach to gene set analysis of GWAS data. *American journal of human genetics*, 86(6) :860–71, June 2010.
- [2] Amanda J Myers, J Raphael Gibbs, Jennifer A Webster, Kristen Rohrer, Alice Zhao, Lauren Marlowe, Mona Kaleem, Doris Leung, Leslie Bryden, Priti Nath, Victoria L Zismann, Keta Joshipura, Matthew J Huentelman, Diane Hu-Lince, Keith D Coon, David W Craig, John V Pearson, Peter Holmans, Christopher B Heward, Eric M Reiman, Dietrich Stephan, and John Hardy. A survey of genetic human cortical gene expression. *Nature genetics*, 39(12) :1494–9, December 2007.

- [3] Matt Silver, Eva Janousova, Xue Hua, Paul M Thompson, and Giovanni Montana. Identification of gene pathways implicated in Alzheimer's disease using longitudinal imaging phenotypes with sparse regression. *NeuroImage*, 63(3) :1681–1694, November 2012.
- [4] Jennifer A Webster, J Raphael Gibbs, Jennifer Clarke, Monika Ray, Weixiong Zhang, Peter Holmans, Kristen Rohrer, Alice Zhao, Lauren Marlowe, Mona Kaleem, Donald S McCorquodale, Cindy Cuello, Doris Leung, Leslie Bryden, Priti Nath, Victoria L Zismann, Keta Joshipura, Matthew J Huentelman, Diane Hu-Lince, Keith D Coon, David W Craig, John V Pearson, Christopher B Heward, Eric M Reiman, Dietrich Stephan, John Hardy, and Amanda J Myers. Genetic control of human brain transcript expression in Alzheimer disease. *American journal of human genetics*, 84(4) :445–58, April 2009.