

Technical advance

Open Access

Application of microarray outlier detection methodology to psychiatric research

Carl Ernst¹, Alexandre Bureau² and Gustavo Turecki^{* 1,3}

Address: ¹McGill Group for Suicide Studies, McGill University, Montreal, Canada, ²Centre de recherche Université Laval Robert-Giffard and Department of social and preventive medicine, Université Laval, Canada and ³Douglas Hospital Research Centre, Pavilion Frank B Common, Rm. F-3125, 6875 LaSalle, Blvd., Verdun, Montreal, Quebec, H4H 1R3, Canada

Email: Carl Ernst - carl.ernst@mail.mcgill.ca; Alexandre Bureau - Alexandre.Bureau@msp.ulaval.ca; Gustavo Turecki* - gustavo.turecki@mcgill.ca

* Corresponding author

Published: 23 April 2008

Received: 29 February 2008

BMC Psychiatry 2008, 8:29 doi:10.1186/1471-244X-8-29

Accepted: 23 April 2008

This article is available from: <http://www.biomedcentral.com/1471-244X/8/29>

© 2008 Ernst et al; licensee BioMed Central Ltd.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/2.0>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Abstract

Background: Most microarray data processing methods negate extreme expression values or alter them so that they do not lie outside the mean level of variation of the system. While microarrays generate a substantial amount of false positive and spurious results, some of the extreme expression values may be valid and could represent true biological findings.

Methods: We propose a simple method to screen brain microarray data to detect individual differences across a psychiatric sample set. We demonstrate in two different samples how this method can be applied.

Results: This method targets high-throughput technology to psychiatric research on a subject-specific basis.

Conclusion: Assessing microarray data for both mean group effects and individual effects can lead to more robust findings in psychiatric genetics.

Background

Currently used psychiatric nosology is based on a compilation of clinical symptoms into categories based primarily on symptom clustering and course. Diagnostic systems such as the current version of the DSM, allow for certain flexibility in the definition of diagnostic categories, with no assumption that each category of mental disorder is a completely discrete entity. As such, individuals diagnosed under a certain diagnostic class are not clinically homogeneous, there are no clear boundaries between classes, and different classes are not mutually exclusive. It is, therefore, unrealistic to expect that all subjects diagnosed with a given disorder will share a common

psychopathological process, which would be associated with a common underlying biological process.

Most research efforts in psychiatry are directed towards the identification of group effects, negating the fact that significant etiological heterogeneity may exist. This limitation is particularly true for microarray research in psychiatry, where gene expression from different brain areas has been assessed comparing all affected subjects to non-affected subjects. A possible solution would be to carry out studies aiming at the identification of biologically meaningful effects focusing on single individuals or subgroups. This approach would mimic fruitful efforts in the identification of genetic factors underlying heterogeneous

conditions such as, among others, spinocerebellar ataxia [1] and Alzheimer's disease [2].

Microarray data from psychiatric subjects can be investigated for individual or subgroup effects that may be of genuine biological significance. Specifically, we hypothesize that specific subgroups can be identified through microarray data screening for extreme expression values. Three previous studies have described how microarray data can be investigated on a subject-specific basis when analyzing data from cancer studies [3-5]. We suggest here that microarray experiments using brain-gene expression levels from psychiatric experiments (e.g. schizophrenia group vs non-schizophrenia group) can utilize microarray data not only for group mean effects (i.e. standard microarray analysis) but also, whenever possible, should evaluate expression levels by individual subjects.

Most microarray projects in psychiatry involve examining more than one neural region [6-9]. Specifically, researchers studying gene expression in brain tend to analyze more than one brain region on more than one array. This leaves researchers with gene expression data from multiple brain regions for each subject. This offers the possibility of using data from different arrays as confirmations of findings which may appear to be outliers. Any outlier on one chip that is also an outlier on a different chip may represent a valid finding.

Human brain has been categorized in two main ways: either by gross anatomical structure, the preference of imaging specialists, or by Brodmann region, the preference of neuro-anatomists. Irrespective of how the human brain is categorized anatomically, what is less obvious is whether gene expression varies between neighboring regions. Two recent, replicating studies suggest that brain gene expression of samples from the same individual, while non-identical, are biologically-related [9,10]. This sharing of similar expression patterns across samples allows for the exclusion of extreme values in the microarray data due to noise. This provides a potential to validate microarray data, particularly for variables that are extreme values, across chips. Those extreme values present across chips for the same probe set and the same individual may represent a true biological effect.

We have designed a method that can assess extreme values and utilizes expression data across chips from the same individual. The method will allow for the detection of any subjects that have probe set values that differ drastically from a mean and outside of a certain threshold (e.g. Standard deviations from a mean). This method, termed Extreme Values Analysis (EVA), takes into account the complex and heterogeneous nature of psychiatric diseases. We illustrate this approach in two different situa-

tions. First, in a publicly available sample where extreme values were simulated, and second, in a sample of subjects who died by suicide and sudden death controls. EVA functions to screen microarray data individual-by-individual in search for any extreme values that may signify some abnormality.

Methods

We used a publicly available data set to first evaluate EVA. The data set comprises 9 subjects screened over 20 regions of the CNS and can be found here [9]. In this dataset, one of the authors (AB) inputted simulated extreme values for two subjects across all CNS regions for one randomly selected probe set each. The expression values were multiplied by $4^{(1+0.25Z)}$ for one of the probe sets and by $0.25^{(1+0.25Z)}$ for the other, where Z is a standard normal random variable which was generated independently for each CNS region. Another of the authors (CE), blinded to the experimental manipulation, applied the method to detect the inputted value(s). The rationale for this experiment is to determine if EVA can detect an artificially generated extreme value in one probe set from > 11 million different data points (10 subjects X 20 regions X ~55,000 probe sets).

We assessed EVA in a second sample that comprised a group of suicide completers and sudden death controls. Information on the subjects, clinical variables, and microarray data quality of the suicide and sudden death controls can be found in Sequeira et al., [6].

EVA can be applied under a control:experimental design (suicide and sudden death controls example) or in a one sample design (CNS screening example). We describe the control:experimental setting, although the description applies also to the one sample design. In the one sample design, all individual values are compared to the group to which they belong.

The mean and standard deviation (SD) of \log_2 -transformed expression level is computed in the experimental group for all probe sets in every region. In our example, this was done in 2 cortical brain regions from suicide subjects. Log transformation stabilizes the variance, allowing comparison of SD across probe sets. After this step, the probe sets with the highest SD values were selected for further analysis. We used only those probe sets in the top 5% of SD values. We reasoned that these probe sets likely have individual values that are extreme, which accounts for a high SD value.

To buffer against detecting mathematical artifacts, EVA selects only those probe sets with high SD values in all regions. In our example, we selected probe sets that were common across both cortical regions. Next, we assess

whether the same subject is responsible for the high SD value across brain regions. We set as criterion for an extreme expression a value of ± 3 fold greater than the mean expression level of the specific probe set among the control group (in our example, the sudden death controls). This approach operates on the assumption that neighboring brain regions are not discrete units and that gene expression should not vary widely from one cortical region to another. Even if brain region-specific expression is more common, it is not expected that a subject that is an outlier in one region is necessarily an outlier in a neighboring region. In other words, extreme values that are detected across multiple brain regions are more likely to represent real biological phenomena. We note that this method is conservative.

Individual expression values also have to be outside of 1.5 SD's of the control group, after having met the above criteria. While we selected 1.5 SD's from the mean of the opposite group, this number can be changed depending on the false discovery level acceptable to the experimenter. Manipulating the SD threshold establishes the false discovery rate (FDR) of the experiment.

The statistical significance of each identified outlier can be assessed by computing the p-value of the subject's expression values for a probe set in the multiple brain regions compared to the multivariate distribution of the expression values in the control group. The null distribution of the \log_2 -transformed probe set-specific expression is estimated by fitting a normal mixed model where the subject effect is random. Letting, X_{ij} be the probe set-specific expression of the i^{th} subject in the j^{th} brain region, and $Y_{ij} = \log_2(X_{ij})$, this model has the form:

$$Y_{ij} = \mu_j + a_i + e_{ij}, \quad a_i \sim N(0, \tau^2), \quad e_{ij} \sim N(0, \sigma^2)$$

where μ_j is the region-specific mean expression, a_i is the subject random effect and e_{ij} is the residual. We fit such a model by restricted maximum likelihood (REML) using the maanova package [11] for the R statistical software [12]. The subject random effect captures the expected correlation between expression in different brain regions of the same subject. The p-value for the observed deviation of the \log_2 -transformed expression level of the i^{th} subject from the mean of the group of reference $\hat{\mu}_j, j = 1, \dots, J$ (or observed fold change on the original scale) is given by

$$P(|Y_i - \mu_1| > |y_{i1} - \mu_1|, \dots, |Y_j - \mu_j| > |y_{ij} - \mu_j|)$$

which we compute using a multivariate t-distribution with the covariance matrix estimated under the normal mixed model.

Results

EVA in partially simulated data

We tested EVA in a sample data set that included 20 different CNS regions [9]. This dataset was selected because A) we could test how the method works with the RMA algorithm and B) we could demonstrate the method in a one-sample case.

We began by computing the standard deviation (SD) for three of the 20 CNS regions described in this data set. The probe sets in the top 5% of SD values was selected for each of three regions and those probe sets that were common to all regions were selected. Five hundred forty-five probe sets were common to all three regions. Next, we screened for any individual values that lay outside of ± 1.5 SD's and was three-fold different from the mean. There were 14 genes that were found to be 3-fold greater than the mean and outside of $+1.5$ SD's and 245 values that were three fold below the mean and outside of -1.5 SD's. Each of these values was then cross-referenced across all 20 CNS regions. Two probe sets were found that met all criteria (1 above the mean for one subject and one below the mean for another subject). These were the probe sets that had been artificially altered (Figure 1).

EVA in real microarray data

To demonstrate this technique, we used a sample that included 6 control subjects and 8 suicide completers with microarray data from BA 8/9 and BA 11. We first screened all expression values for MAS 5.0 present/absent call leaving 14,896 probe sets in BA 8/9 and 14,412 probe sets in BA 11. We next calculated a standard deviation for all probe sets from suicide subjects. This was done using \log_2 -transformed expression values. We then selected the probe sets with the highest SD values (top 5%) from both BA 8/9 and BA 11.

Any probe sets that was identical to both BA 8/9 and BA 11 after SD filtering was selected. There were 180 probe sets that were common to both regions. Next, to account for the variability of expression in control values, we searched the data for any suicide data point greater than 3-fold from the control mean and outside of 1.5 SD's. We reasoned that an extreme value across all regions for the same subject(s) could represent a biologically relevant event.

Beginning in BA 8/9, we filtered the 180 probe sets for those probe sets from suicide completers outside of 3-fold from the control mean. There were 20 probe sets where X_{ij} (a particular expression value from a particular subject in a given brain region) was not outside of 3-fold from the control mean. From the 160 remaining probe sets, 108 probe sets were also outside of 1.5 SD's in BA 8/9. Probe sets from BA 11 were then filtered for these probe sets.

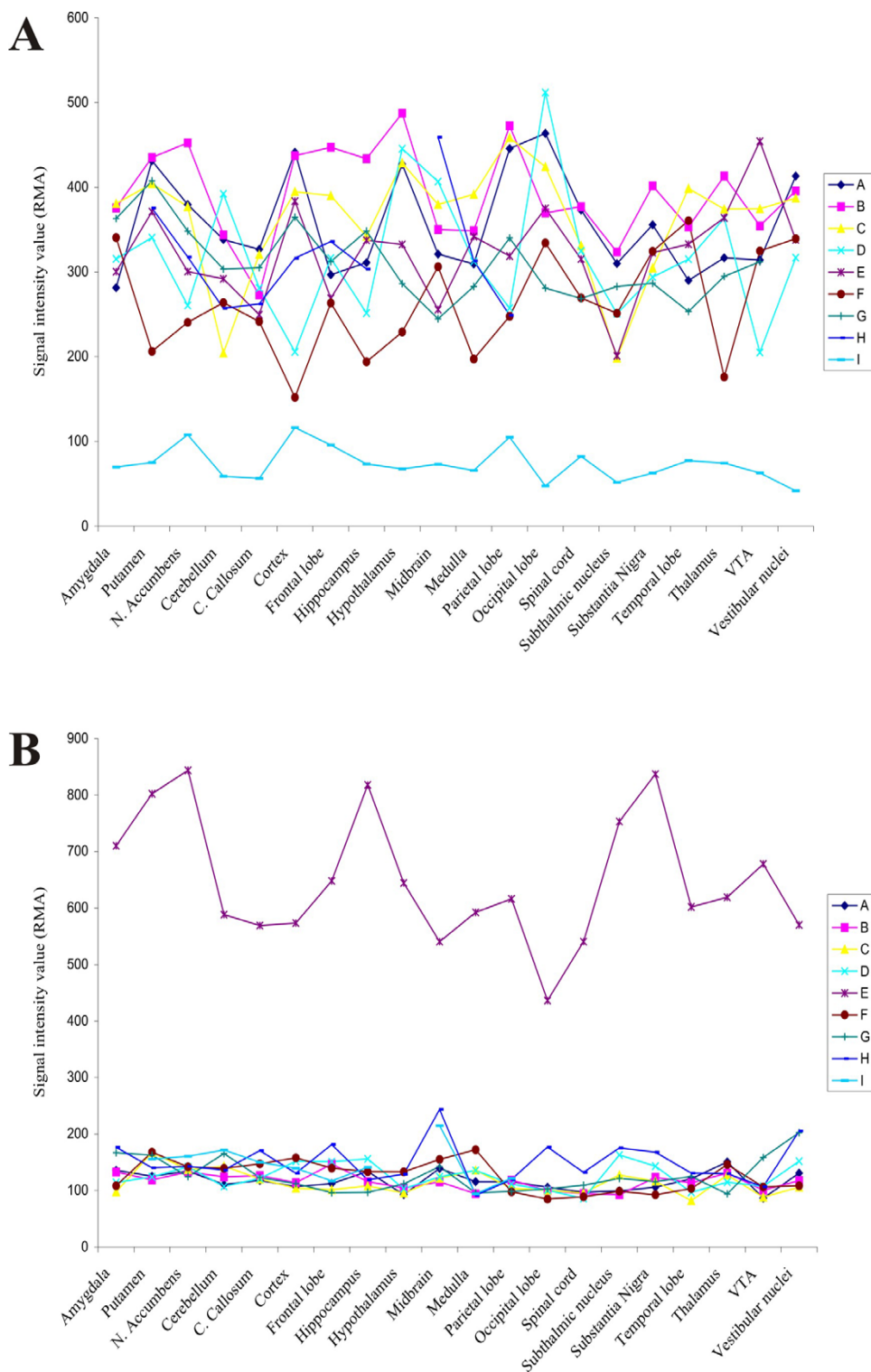


Figure 1
Microarray expression values demonstrating two subjects who passed all EVA criteria in a one sample case. A) Extreme low expressor (light blue trace) compared to other subjects in the same sample for one probe set identified across multiple CNS regions. Each subject is represented by a letter (A, B, C...). **B)** Extreme high expressor (purple trace) compared to other subjects in the same sample for one probe set across multiple CNS regions.

Table 1: HG-UI33 plus 2 probe sets that met all EVA criteria. Numbers represent p-values generated for each probe set across each subject (S). Subjects who met EVA criteria have p-values underlined. Note that p-values are generated from the number of SD's from the mean, therefore some subjects with very small p-values may be outside of a given number of SD's but < 3-fold different than the mean.

| Probe set | S1 | S2 | S3 | S4 | S5 | S6 | S7 | S8 |
|-------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|-----------------|
| 233814_at | <u>0.00011</u> | 0.262901 | 0.060973 | 0.119425 | 0.340339 | 0.065993 | 0.291037 | 0.044823 |
| 225440_at | <u>0.003777</u> | 0.041856 | 0.549775 | 0.617546 | <u>0.000754</u> | 0.044895 | 0.52789 | <u>0.000355</u> |
| 203638_s_at | <u>0.008775</u> | 0.076607 | 0.099598 | 0.154338 | <u>0.000159</u> | 0.089298 | 0.01575 | 0.357839 |
| 214680_at | <u>2.89E-05</u> | <u>0.010459</u> | <u>0.007155</u> | <u>0.001488</u> | <u>0.00034</u> | <u>0.007818</u> | <u>0.000768</u> | 0.148031 |
| 37170_at | <u>0.001734</u> | <u>0.002016</u> | 0.095708 | 0.010663 | <u>0.000198</u> | <u>0.000259</u> | <u>0.003458</u> | 0.140702 |
| 227556_at | <u>0.00074</u> | 0.397533 | <u>0.000597</u> | 0.010661 | 0.291994 | <u>7.80E-07</u> | 0.100161 | 0.010813 |
| 229917_at | <u>0.00895</u> | 0.170125 | 0.421083 | 0.044559 | 0.401754 | <u>0.000248</u> | 0.055834 | <u>0.003505</u> |
| 227330_x_at | 0.018916 | 0.412934 | 0.378829 | 0.053889 | <u>0.002547</u> | 0.576196 | 0.037467 | 0.46594 |
| 231804_at | 0.411319 | 0.07522 | 0.566597 | 0.115634 | 0.485137 | 0.489949 | 0.28554 | <u>0.012676</u> |
| 200904_at | 0.066736 | 0.080689 | 0.236431 | 0.504115 | <u>0.009175</u> | 0.255437 | 0.466512 | 0.105319 |
| 230141_at | 0.041214 | 0.157873 | 0.408831 | 0.09835 | <u>0.013597</u> | 0.067947 | 0.333108 | <u>0.004167</u> |
| 222020_s_at | <u>6.92E-05</u> | 0.071294 | 0.239087 | 0.018151 | 0.271669 | 0.043433 | 0.044077 | 0.428814 |
| 213812_s_at | <u>0.004793</u> | <u>0.002408</u> | 0.151259 | 0.121657 | 0.15422 | 0.030286 | 0.022571 | 0.02784 |
| 201505_at | 0.316616 | 0.570798 | 0.160905 | 0.60401 | 0.000337 | 0.19324 | 0.157064 | <u>0.00128</u> |
| 240467_at | 0.103838 | <u>0.007113</u> | 0.076142 | <u>0.001703</u> | 0.050236 | 0.332126 | 0.148029 | 0.199109 |
| 229861_at | <u>0.005266</u> | 0.047216 | 0.055746 | <u>0.000202</u> | 0.028745 | 0.028233 | 0.019457 | <u>0.00037</u> |
| 225872_at | <u>0.001483</u> | 0.528714 | 0.628325 | 0.427906 | <u>0.000266</u> | 0.396505 | 0.019428 | 0.194352 |
| 241758_at | 0.625246 | 0.629225 | 0.435153 | 0.452194 | <u>0.005473</u> | 0.282962 | 0.223662 | <u>0.001718</u> |
| 214449_s_at | 0.400041 | 0.272717 | 0.30737 | 0.020469 | <u>0.000229</u> | 0.040009 | 0.340049 | <u>0.000813</u> |
| 200648_s_at | 0.011284 | 0.10939 | 0.371707 | 0.02763 | <u>0.013746</u> | 0.07564 | 0.02506 | 0.264293 |
| 221795_at | <u>0.000617</u> | <u>0.006302</u> | 0.070084 | <u>0.008107</u> | <u>0.005545</u> | <u>0.039509</u> | <u>0.01825</u> | 0.284639 |
| 204379_s_at | 0.046084 | 0.083475 | 0.391619 | 0.237209 | <u>0.025689</u> | 0.079905 | 0.073448 | 0.22916 |
| 215172_at | 0.109908 | 0.139031 | 0.199485 | 0.152187 | 0.155703 | 0.203207 | 0.144057 | <u>0.001802</u> |
| 203324_s_at | <u>0.005691</u> | 0.092966 | 0.150765 | <u>0.000126</u> | <u>0.00021</u> | 0.056637 | 0.205541 | 0.183977 |
| 202800_at | 0.077088 | 0.491247 | 0.411825 | 0.162812 | <u>0.013449</u> | 0.102756 | 0.028244 | 0.452363 |
| 236223_s_at | 0.030208 | 0.196886 | 0.209618 | <u>0.000686</u> | 0.25837 | 0.13088 | 0.373756 | 0.268783 |
| 209023_s_at | <u>5.10E-05</u> | <u>0.000708</u> | 0.030401 | <u>0.001934</u> | 0.251739 | <u>0.002706</u> | <u>2.10E-05</u> | 0.054205 |
| 213593_s_at | <u>0.005575</u> | <u>0.001263</u> | 0.108932 | <u>0.005416</u> | <u>0.007774</u> | 0.20946 | 0.083709 | 0.036244 |
| 222249_at | 0.101364 | 0.131587 | 0.040535 | <u>0.007024</u> | <u>0.009912</u> | 0.096597 | <u>0.00014</u> | 0.012586 |
| 220460_at | 0.018979 | 0.539318 | 0.195645 | 0.120094 | <u>0.009469</u> | 0.067115 | 0.013541 | 0.346847 |
| 201656_at | <u>0.000481</u> | 0.13031 | 0.064596 | 0.010041 | <u>0.000455</u> | 0.018925 | 0.14539 | 0.117989 |
| 235775_at | <u>0.000373</u> | 0.079169 | 0.040214 | 0.11224 | <u>0.000676</u> | 0.222721 | 0.208163 | 0.049906 |
| 204516_at | <u>9.40E-05</u> | 0.051213 | 0.096458 | 0.027735 | <u>0.003274</u> | 0.277745 | 0.208683 | 0.17238 |
| 201843_s_at | 0.047846 | 0.096131 | 0.228376 | 0.016614 | 0.010238 | 0.104263 | <u>0.005888</u> | 0.155651 |
| 204712_at | 0.024623 | 0.417385 | 0.556647 | 0.127233 | <u>0.010306</u> | 0.182687 | 0.121824 | 0.39341 |
| 224736_at | <u>0.002159</u> | 0.02277 | 0.118888 | 0.114733 | 0.126571 | 0.333625 | 0.187404 | 0.021966 |
| 214203_s_at | 0.018997 | 0.038021 | 0.062828 | 0.074293 | <u>0.00133</u> | <u>0.017613</u> | <u>0.003911</u> | 0.047244 |
| 200914_x_at | 0.010589 | <u>0.005141</u> | 0.075284 | 0.058115 | 0.042784 | 0.274358 | 0.103149 | 0.103501 |
| 222404_x_at | <u>0.003613</u> | 0.040237 | 0.059019 | 0.220239 | 0.151453 | 0.188272 | 0.358947 | 0.103866 |
| 229553_at | <u>0.008515</u> | 0.022552 | 0.02734 | 0.070236 | 0.14418 | 0.145378 | 0.331622 | 0.223119 |
| 203249_at | 0.139981 | 0.138389 | 0.229283 | 0.100477 | <u>0.004921</u> | 0.094854 | 0.07573 | <u>0.002453</u> |
| 203041_s_at | <u>0.01742</u> | 0.111182 | 0.227792 | 0.402694 | 0.014214 | 0.310238 | 0.14055 | 0.174457 |
| 209292_at | 0.068175 | 0.345992 | 0.507253 | 0.01934 | <u>0.006611</u> | 0.099673 | 0.002438 | 0.275859 |
| 226084_at | <u>0.00373</u> | <u>0.005272</u> | 0.059754 | 0.096792 | 0.076702 | 0.642227 | 0.101137 | 0.153124 |
| 204976_s_at | 0.00165 | 0.083667 | 0.100336 | 0.152992 | <u>0.00468</u> | 0.159244 | 0.449992 | 0.583733 |
| 212368_at | <u>0.013283</u> | <u>0.016864</u> | 0.07334 | 0.086183 | 0.13413 | 0.246171 | 0.126799 | 0.293217 |
| 211962_s_at | 0.00228 | 0.17037 | 0.178376 | 0.026324 | <u>0.0011</u> | 0.017557 | 0.011854 | 0.278736 |
| 226228_at | 0.090147 | 0.423425 | 0.572415 | 0.125902 | <u>0.015242</u> | 0.100574 | 0.077677 | 0.633229 |
| 213954_at | <u>0.002753</u> | <u>0.001918</u> | 0.022841 | 0.013454 | <u>8.47E-05</u> | 0.018911 | 0.075725 | 0.077512 |
| 213922_at | <u>0.004719</u> | <u>0.003408</u> | 0.109555 | 0.024766 | 0.07629 | 0.417034 | 0.042621 | 0.199903 |
| 221517_s_at | <u>0.001167</u> | <u>0.00413</u> | 0.048293 | 0.022654 | <u>0.003271</u> | 0.150882 | 0.080039 | 0.011072 |
| 227099_s_at | 0.047549 | 0.092241 | 0.069821 | 0.333917 | <u>0.000425</u> | 0.00361 | 0.200933 | <u>0.009189</u> |
| 201737_s_at | <u>0.000159</u> | 0.012972 | 0.219393 | 0.110936 | 0.092836 | 0.273891 | 0.078928 | 0.035264 |
| 214279_s_at | <u>0.008594</u> | 0.363441 | 0.365898 | 0.01142 | 0.036218 | 0.014181 | 0.059127 | 0.066575 |
| 205709_s_at | <u>0.007059</u> | 0.021998 | 0.110402 | 0.359529 | 0.08021 | 0.218438 | 0.306235 | 0.206374 |
| 225810_at | 0.039054 | 0.256396 | 0.159142 | 0.452173 | <u>0.001525</u> | 0.010676 | 0.017581 | 0.130546 |
| 226435_at | 0.150996 | 0.248593 | 0.432525 | <u>0.004589</u> | 0.032368 | 0.12243 | 0.072016 | 0.216067 |

Table 1: HG-U133 plus 2 probe sets that met all EVA criteria. Numbers represent p-values generated for each probe set across each subject (S). Subjects who met EVA criteria have p-values underlined. Note that p-values are generated from the number of SD's from the mean, therefore some subjects with very small p-values may be outside of a given number of SD's but < 3-fold different than the mean. (Continued)

| | | | | | | | | |
|-------------|-----------------|-----------------|----------|-----------------|-----------------|-----------------|----------|-----------------|
| 226364_at | 0.310524 | 0.026062 | 0.011704 | 0.02649 | 0.026889 | 0.178983 | 0.048683 | <u>0.005222</u> |
| 240482_at | 0.318335 | 0.475263 | 0.195512 | 0.214718 | 0.151663 | 0.390872 | 0.003578 | <u>0.001316</u> |
| 204881_s_at | <u>0.00017</u> | 0.068969 | 0.075294 | 0.117202 | 0.020759 | 0.039673 | 0.310769 | 0.084319 |
| 203841_x_at | <u>0.000201</u> | 0.005823 | 0.089297 | <u>0.006843</u> | 0.052791 | 0.114169 | 0.037911 | 0.052824 |
| 212677_s_at | 0.021412 | <u>0.003193</u> | 0.033988 | 0.022312 | 0.03778 | 0.128097 | 0.090665 | 0.060285 |
| 201502_s_at | <u>0.001372</u> | <u>0.006379</u> | 0.167335 | 0.041531 | 0.035423 | 0.070657 | 0.263158 | 0.614172 |
| 240299_at | 0.017658 | 0.276529 | 0.197699 | <u>0.008206</u> | 0.0622 | 0.215909 | 0.087937 | 0.025787 |
| 212423_at | 0.024073 | <u>0.015847</u> | 0.140173 | 0.140948 | 0.19045 | 0.020469 | 0.19758 | 0.399876 |
| 228811_at | 0.290584 | 0.349143 | 0.140068 | 0.371578 | 0.50686 | <u>0.002176</u> | 0.066324 | <u>0.002291</u> |
| 224737_x_at | <u>0.000367</u> | 0.288851 | 0.112457 | <u>0.002641</u> | 0.097181 | 0.24341 | 0.219535 | 0.021662 |
| 201019_s_at | <u>0.003446</u> | <u>0.00455</u> | 0.276132 | 0.181832 | 0.064501 | 0.680406 | 0.179361 | 0.099558 |
| 229281_at | 0.090713 | 0.333165 | 0.42729 | 0.241316 | <u>0.000732</u> | 0.019187 | 0.025054 | 0.314121 |

Table 1 lists the probe sets across the eight suicide completers and the individual p-values associated with each subject. From 108 probe sets that passed all EVA criteria in BA 8/9, 69 passed all EVA criteria in BA 11. Included in this list of probe sets are a number of genes that have been linked to suicide before including the FGF family [13], NTRK2 [14], and members of the ubiquitin family [15]. Of note, from the table, is that for a number of probe sets there is more than one subject who has an extreme expression value reaching a significance level.

Discussion

The extreme values analysis, or EVA, is a method to detect individual or subsets of outliers for a given probe set in microarray experiments. The rationale for this type of experiment is that psychopathology is not necessarily group specific but more likely sub-group or subject specific. The method outlined here uses log-transformed data to determine which probe sets have the highest variance and screens out those probe sets with little variation. This step is intended to select those probe sets with values that deviate widely from the mean. Next the method compares individual data points to a control mean, and searches for any 3-fold changes. Selected values also have to be outside of 1.5 SD's from the mean. These values were considered extreme expression values. These extreme expression values were next verified in one other cortical region to determine if they were extreme expression values in other cortical regions as well. We reasoned that the use of other cortical regions functioned as replicate experiments and enforced the finding.

We also evaluated this method in a one sample case after inputting artificial values for one probe set across all CNS regions in RMA data. EVA was able to detect the inputted value; the only difference between the control:experimental case and one sample case is the mean value used: In the one sample case the mean used includes the extreme value while in the control:experimental case it does not.

The use of multiple cortical regions as within-subject replicates is a way to detect true extreme expression values in individual subjects. Operating under the assumption the gene expression in one cortical region is generally similar in neighboring cortical regions, we propose that different chips for the same subject can be used as replicate experiments, if probe set outliers on an individual specific basis are being investigated. If an observed outlier is a real biological event, it is very probable that the same subject on the same probe set will also be an outlier in a neighboring region. Consider, for example, the family with a deletion in the MAOA gene [16]. Had this family undergone post-mortem microarray analysis as a part of a larger sample of subjects, EVA would have detected the MAOA decrease in expression whereas microarray analysis using mean group effects would not have. Using multiple brain regions as replicates does undermine the idea that gene expression is different across different brain regions, which it is [10,17]; however, it means that if an effect is detected, it is likely real and robust.

Comparison to PPST method

The PPST method [5] counts the number of subjects in both control and experimental group outside of the 95th percentile of the opposite group. The FDR is therefore controlled by altering the percentile threshold. EVA uses the SD from the opposite group and counts the number of subjects that are outside a given SD value (± 1.5 SDs in this study). Selecting more stringent SD values allows for direct manipulation of the FDR. In this study a liberal cut-off was chosen (outside of 1.5SD's). The FDR among the detected outliers could be estimated from the p-value of the subject's expression values using standard methods such as that of Reiner et al. [18]

Comparison to COPA method

Cancer outlier profile analysis (COPA) is another outlier detection that has proved fruitful in the past[4]. This technique normalizes all probe sets (one sample design) and

calculates the 75th, 90th, and 95th percentiles for each probe set and rank-orders them by percentile score. A prioritized list of probe set with some subjects that have extreme expression values is then investigated. Tibshirani and Hastie [3] introduce the outlier-sum statistic in their paper to improve on the COPA method. Their method differs from COPA by the standardization procedure of each probe set expression level using the median and median absolute deviation.

There are some caveats to be aware of before proceeding with this approach to screen microarray data. Firstly, the method is very conservative and likely has a high beta error rate. It is very likely that there were a number of true positives that were not detected because of the rigidity of the design. Some parameters may need to be adjusted to allow more probe sets to pass filtering (e.g. top 10% of SD values instead of the top 5% being used). Second, this method has the disadvantage of requiring a number of replicates per individual, a component that could be cost-prohibitive. Third, the method can only be used to study genes whose expression levels are similar across brain regions. Finally, we note that all probe-level microarray algorithms dampen extreme values at the scanner. This method is conservative and could only be used to investigate extreme values after initial processing.

Our view for this technique is as another analysis technique to further explore microarray data, in conjunction with more mainstream techniques [19]. This method, termed Extreme Values Analysis, can detect extreme differences in gene expression on a subject-by-subject basis from microarray data across different chips. The method uses high-throughput technology in a non-biased way to understand psychiatric disease for each subject investigated.

Competing interests

The author(s) declare that they have no competing interests.

Authors' contributions

CE conceived of the study, wrote the manuscript, and analyzed data. AB designed the statistical test and wrote the manuscript. GT participated in study design and coordination, and wrote the manuscript. All authors read and approved the final manuscript.

Acknowledgements

This study was funded by grants to GT from the Canadian Institute of Health Research, the Fonds de la Recherche en Sante du Quebec, and the American Foundation for Suicide Prevention. CE is supported by a Natural Science and Engineering Research Council Canada Graduate Scholarship award. AB is supported by a scientist award from the Fonds de la Recherche en Sante du Quebec.

References

- Schols L, Bauer P, Schmidt T, Schulte T, Riess O: **Autosomal dominant cerebellar ataxias: clinical features, genetics, and pathogenesis.** *Lancet Neurol* 2004, **3(5)**:291-304.
- Bertram L, Tanzi RE: **Alzheimer's disease: one disorder, too many genes?** *Hum Mol Genet* 2004, **13 Spec No 1**:R135-41.
- Tibshirani R, Hastie T: **Outlier sums for differential gene expression analysis.** *Biostatistics* 2007, **8(1)**:2-8.
- Lyons-Weiler J, Patel S, Becich MJ, Godfrey TE: **Tests for finding complex patterns of differential expression in cancers: towards individualized medicine.** *BMC Bioinformatics* 2004, **5**:110.
- Kim S, Choi KH, Baykiz AF, Gershenfeld HK: **Suicide candidate genes associated with bipolar disorder and schizophrenia: an exploratory gene expression profiling analysis of post-mortem prefrontal cortex.** *BMC Genomics* 2007/11/14 edition. 2007, **8**:413.
- Karssen AM, Her S, Li JZ, Patel PD, Meng F, Bunney WE Jr., Jones EG, Watson SJ, Akil H, Myers RM, Schatzberg AF, Lyons DM: **Stress-induced changes in primate prefrontal profiles of gene expression.** *Mol Psychiatry* 2007/09/26 edition. 2007, **12(12)**:1089-1102.
- Roth RB, Hevezi P, Lee J, Willhite D, Lechner SM, Foster AC, Zlotnik A: **Gene expression analyses reveal molecular relationships among 20 regions of the human CNS.** *Neurogenetics* 2006, **7(2)**:67-80.
- Sequeira A, Gwady FG, Ffrench-Mullen JM, Canetti L, Gingras Y, Casero RA Jr., Rouleau G, Benkelfat C, Turecki G: **Implication of SSAT by gene expression and genetic variation in suicide and major depression.** *Arch Gen Psychiatry* 2006, **63(1)**:35-48.
- Kerr MK, Martin M, Churchill GA: **Analysis of variance for gene expression microarray data.** *J Comput Biol* 2000, **7(6)**:819-837.
- . *The R Development Core Team 2007; www-r-project.org*.
- Evans SJ, Choudary PV, Neal CR, Li JZ, Vawter MP, Tomita H, Lopez JF, Thompson RC, Meng F, Stead JD, Walsh DM, Myers RM, Bunney WE, Watson SJ, Jones EG, Akil H: **Dysregulation of the fibroblast growth factor system in major depression.** *Proc Natl Acad Sci U S A* 2004, **101(43)**:15506-15511.
- Dwivedi Y, Rizavi HS, Conley RR, Roberts RC, Tamminga CA, Pandey GN: **Altered gene expression of brain-derived neurotrophic factor and receptor tyrosine kinase B in postmortem brain of suicide subjects.** *Arch Gen Psychiatry* 2003, **60(8)**:804-815.
- Ryan MM, Lockstone HE, Huffaker SJ, Wayland MT, Webster MJ, Bahn S: **Gene expression analysis of bipolar disorder reveals downregulation of the ubiquitin cycle and alterations in synaptic genes.** *Mol Psychiatry* 2006, **11(10)**:965-978.
- Brunner HG, Nelen M, Breakefield XO, Ropers HH, van Oost BA: **Abnormal behavior associated with a point mutation in the structural gene for monoamine oxidase A.** *Science* 1993, **262(5133)**:578-580.
- Ernst C, Sequeira A, Klempan T, Ernst N, Ffrench-Mullen J, Turecki G: **Confirmation of region-specific patterns of gene expression in the human brain.** *Neurogenetics* 2007.
- Lein ES, Callaway EM, Albright TD, Gage FH: **Redefining the boundaries of the hippocampal CA2 subfield in the mouse using gene expression and 3-dimensional reconstruction.** *J Comp Neurol* 2005/03/19 edition. 2005, **485(1)**:1-10.
- Reiner A, Yekutieli D, Benjamini Y: **Identifying differentially expressed genes using false discovery rate controlling procedures.** *Bioinformatics* 2003/02/14 edition. 2003, **19(3)**:368-375.
- Tomlins SA, Rhodes DR, Perner S, Dhanasekaran SM, Mehra R, Sun XW, Varambally S, Cao X, Tchinda J, Kuefer R, Lee C, Montie JE, Shah RB, Pienta KJ, Rubin MA, Chinnaiyan AM: **Recurrent fusion of TMPRSS2 and ETS transcription factor genes in prostate cancer.** *Science* 2005, **310(5748)**:644-648.
- Allison DB, Cui X, Page GP, Sabripour M: **Microarray data analysis: from disarray to consolidation and consensus.** *Nat Rev Genet* 2006, **7(1)**:55-65.

Pre-publication history

The pre-publication history for this paper can be accessed here:

<http://www.biomedcentral.com/1471-244X/8/29/prepub>

Publish with **BioMed Central** and every scientist can read your work free of charge

"BioMed Central will be the most significant development for disseminating the results of biomedical research in our lifetime."

Sir Paul Nurse, Cancer Research UK

Your research papers will be:

- available free of charge to the entire biomedical community
- peer reviewed and published immediately upon acceptance
- cited in PubMed and archived on PubMed Central
- yours — you keep the copyright

Submit your manuscript here:
http://www.biomedcentral.com/info/publishing_adv.asp

