

# Identifying Relevant Genes Related Atopic Dermatitis to using Transcriptomic

Creel, Kristin  
ISPED  
University of Bordeaux  
kristin.creel@etu.u-bordeaux.fr

Dong, Larry  
Department of Epidemiology, Biostatistics  
and Occupational Health  
McGill University  
larry.dong@mail.mcgill.ca

Nama Ravi, Sneha Keerthi  
ISPED  
University of Bordeaux  
sneha-keerthi.nama-ravi@etu.u-bordeaux.fr

## 1. INTRODUCTION

Atopic dermatitis (AD) or atopic eczema is an itchy, inflammatory skin condition characterised by poorly defined erythema with edema vesicles, and weeping in the acute stage and lichenification in the chronic stage. The global prevalence is 15 – 20% in children and 1 – 3% in adults, posing a significant burden on health-care resources and patients' quality of life [9]. The etiological factors associated with the initiation and progress of the disease are known to be genetic, environmental and immunological that affects the epithelial barrier-immunity interplay [10].

Clinical investigations and discoveries in molecular medicine have positively identified 46 genes linked to AD. Mutations in filaggrin (FLG) genes (influencing intermediate filament protein filaggrin expression) are most common in the AD diseased population, it affects 10 – 50% of AD patients worldwide. Few additional barrier genes encoded by the epidermal differentiation complex (EDC) locus chromosome 1q21, including claudins, loricrin (LOR), involucrin (IVL), SPINK5, AND tmem79/matt, are also associated with AD. The genes of innate immune system such as NOD1, NOD2, TLR2, CD14, and DEFB1, all of which encode the integral factors in cutaneous immunologic response to non-specific antigens may also experience mutations and cause AD [6]. Studies have identified FLG gene mutation to be the most significant risk factor for AD, followed by the genes in the type 2 T helper lymphocyte (Th2) signalling pathways. Additionally, gene profiling assays demonstrated AD is associated with decreased gene expression of epidermal differentiation complex genes and elevated Th2 and Th17 genes. Hypomethylation of TSLP and FCER1G in AD were also reported; and miR-155, which targets the immune suppressor CTLA-4, was found to be significantly over-expressed in infiltrating T cells in AD skin lesions [2], [6].

This study attempts to examine the degree of dissimilarity between the transcriptomic expression between le-

sional and non-lesional samples in participants with atopic dermatitis. Data-driven methods are employed to identify relevant genes that may cause AD in better understanding the molecular mechanisms behind the occurrence of this disease.

## 2. MATERIALS AND METHODS

### 2.1. Data Source

The data at hand comes from the Microbes in Allergy and Autoimmunity Related to the Skin (MAARS). The purpose of this multidisciplinary research consortium is to better understand the characteristics surrounding two major chronic inflammatory skin diseases: atopic dermatitis and psoriasis [?]. In the MAARS study, these two illnesses are investigated from different scientific perspectives and one of the research directions involved exploring the transcriptomic data of lesional and non-lesional samples in patients with atopic dermatitis, psoriasis (PSO) and neither of the two studied skin diseases.

*Data Management:* Demographic data for all patients were available alongside their known allergies, concomitant medication and transcriptomic data for sequenced samples. A total of 1317 patients with AD, PSO or neither illness were examined and 618 samples were sequenced. However, because it was only of interest to assess RNA expression levels in samples from patients with atopic dermatitis, two sets of data were discarded prior to conducting subsequent analyses: participants with AD who have no sequenced samples and genetic information for samples from healthy individuals and ones with psoriasis. The normalization procedure for transcriptomic data was already performed prior to obtaining the data for this analysis. No genetic information was missing, so statistical methods for imputation were not needed for this study.

### 2.2. Exploratory Data Analysis

Exploratory data analysis was conducted to better understand the data at hand from a data-driven perspective. To this end, an MA diagram was plotted alongside several clustering methods which were used to help depict possible underlying patterns from the MAARS transcriptomic

data: principal component analysis (PCA) and hierarchical clustering [?]. In statistical machine learning, clustering techniques falls under unsupervised learning, which is a paradigm that attempts to perform inference on data without response variables or “labels” [5]. In the recent years, this area of machine learning has gathered much attention from researchers; data separability without the use of labels is known as disentangling in the autoencoder literature [8] [3]. Here, the sample labels are the lesional status of the skin sample and the primary purpose of performing such unsupervised analyses was to visually assess the separability of the samples using their transcriptomic information.

### 2.3. Statistical Analysis

**2.3.1) Differential Gene Expression Analysis:** The overarching goal of the statistical analysis was the determine the relevant genes from a transcriptomic dataset that are predictive of AD. The analysis plan was performed two-fold; the identification of genes highly associated with AD was first performed using differential gene expression analysis within a multiple hypothesis testing framework. The transcriptomic data was available as a continuous variable since the microarray data has been normalized.

Let  $Y_{ij} \in \mathbb{R}$  denote the gene expression level of gene  $j$  in sample  $i$  and let  $X_i \in \{0,1\}$  represent its lesional status. The following linear model was fitted:

$$\begin{aligned} Y_{ij} &= \beta_{0j} + \beta_{1j}X_i + \sum_{k: k \in \mathcal{P}} \beta_k \mathbf{1}_{pat(i)=k} + \epsilon_{ij} \\ &= \beta_{0j} + \beta_{1j}X_i + \beta_{pat(i)} \mathbf{1}_{pat(i)} + \epsilon_{ij} \end{aligned}$$

where the  $\mathcal{P}$  is the set of all patients and  $pat(i) \in \mathcal{P}$  is used to represent the patient to which belongs sample  $i$ . The coefficient  $\beta_k$  accounts for the intersample correlation for samples which are obtained from the same patient. By construction,  $\mathbf{1}_{pat(i)=k} = 1$  for exactly one  $k \in \mathcal{P}$  and  $|\mathcal{P}| < n$ .

Originally developed for two-color microarray experiments, it is now commonplace to employ an empirical Bayes method to analyze differential gene expression with linear models like the one above [11]. In a hypothesis framework, the following distributional assumptions for  $\beta_{1j}$  are assumed for all genes  $j$ :

$$\begin{aligned} H_0 : \quad & \beta_{1j} | \sigma_j^2 \sim \mathcal{N}(0, \sigma_j^2) \\ H_1 : \quad & \beta_{1j} | \sigma_j^2 \not\sim \mathcal{N}(0, \sigma_j^2) \end{aligned}$$

where obtained  $p$ -values were adjusted using the Benjamini-Hochberg procedure to reduce the error in false discovery rate prone to be inflated due to performing multiple parallel statistical tests [1]. A significance level of 0.05 was used in determining the genes highly associated with the outcome of interest, the lesional status of sequenced samples; a volcano plot was

used to highlight potential significant genes.

An ensemble learning prediction model was then implemented to assess predictability of sample lesional status; this binary model was chosen due to its strengths in obtaining better predictions as it weights multiple models according to how well each perform on this specific task [5]. **MENTION VARIABLE SELECTION.** The receiving operating characteristic (ROC) curve is often used to evaluate data-driven binary prediction models and the area under this curve (AUROC) was used here to quantify the goodness of the prediction model [7]. Internal validation was performed using leave-one-out cross-validation method whereby the AUROC was evaluated at every iteration [4].

## 3. RESULTS

In Fig. 1, a simple 2D visualization of the results obtained from PCA is available, whereby lesional samples are displayed in black and non-lesional samples are displayed in red.

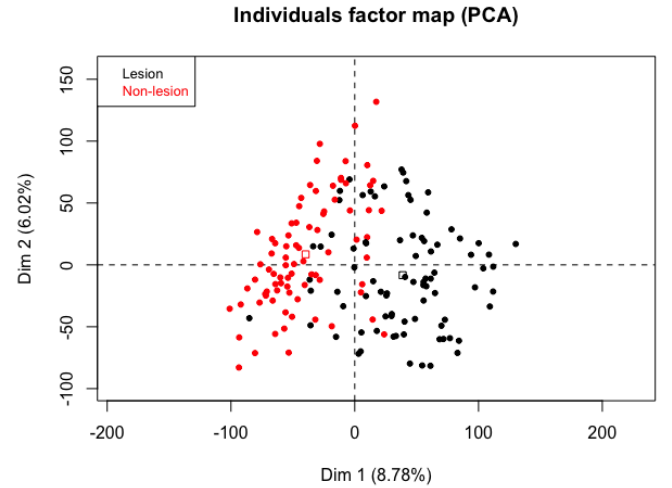


Fig. 1: testing

The horizontal axis represents the first component obtained through PCA and the vertical axis the second; they respectively explain 8.78% and 6.02% of the variance in the transcriptomic data.

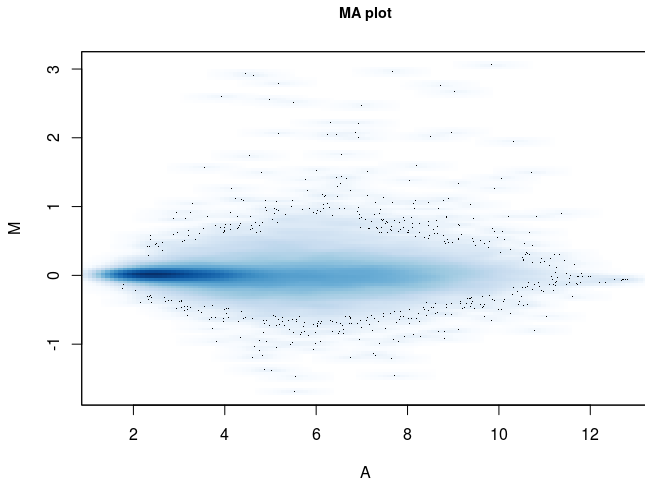
A Differential Gene Analysis (DGE) was conducted using a linear model for paired data. The results indicated 81 genes that differed significantly between the lesioned and non-lesioned skin samples: 65 that are up-regulated and 16 that down-regulated. A heatmap illustrating the normalized gene expression levels is displayed below.

*insert heatmap here*

talk about clustering methods

An MA plot was used to visually represent the genes and determine which are up- and down-regulated. In the plot, each dot represents a gene. The x-axis is the average expression of the gene over the mean of normalized counts;

the y-axis is the log2 fold change between the lesional statuses. The MA plot below indicates some genes that are up-regulated and fewer that are down-regulated.



Additionally, an sPLS analysis was conducted to identify genes that can “best explain” the variability in the lesional status of the sequenced samples. Information regarding the ten first components were retained whereby a different number of genes was selected for each component; the number of variables selected was determined using a 5-fold cross-validation method. This analysis resulted in the identification of 32 genes, 7 of which were found to be common in the two conducted analyses: DGE and sPLS. Genes KRT16, PRSS27, S100A9, S100A12, PI3, GJB2, ARK1B10 were found common in both analyses. Gene annotations using GSA<sub>n</sub> (Gene Set Annotation) CITATION, which includes a search on the Gene Ontology Database (GOA) as well, were performed for the seven genes that were identified as significant in both the DGE and the sPLS analyses. Using both the gene identifiers and their synonyms, GSA<sub>n</sub> retained 19 terms, covering 13 out of 16 genes and 6 of them are synthetic terms.

Figure 2 provides information about the annotated genes within (GOA), the genes covered by GSA<sub>n</sub>, and the group-wise similarity between them. The first gauge shows excellent gene coverage of 100% within the GOA file. The second gauge shows that 81% of the genes were covered with the GSA<sub>n</sub> analysis. The third gauge shows that 44% of the genes share a term in both datasets.

Figure 3 displays the information content and the gene coverage of the synthetic terms. Functional annotation categories from lesional AD samples include serine endopeptidase activity, protein metabolic process, antimicrobial humoral response, cytoskeleton and toll-like receptor binding.

Figure 4 provides more detailed information about the representative terms. Further exploration of the tree will

provide additional information on the terms sharing the same informative ancestor or the genes annotated by more than one term. Using this information, information can be obtained on the biological role of these genes in the occurrence of lesional AD.

#### 4. DISCUSSION

Differential gene expression analysis helps in understanding the patterns of expressed genes in the diseased population which can be used as biomarkers to differentiate the diseased from the controls. The current analysis considered all atopic dermatitis samples and compared lesional and non-lesional subtypes from the same patients. The current analysis was unable to detect much variation in transcriptomic expression between lesional and non-lesional AD skin samples. While some differentiation was identified (up to 25% of the variation explained by five principal components), lesional state does not account for much of the differences in the samples.

Unsupervised analysis like PCA and Dendrogram were used for data exploration. Principal Components Analysis (PCA) is a statistical technique used to explore data during microarray analysis. It is performed to reduce the dimensionality of data by transforming the correlated variables into a smaller number of uncorrelated variables called components. PCA can determine the key variables in the data that best explain the differences in the observations. The matrix of the data used shows rows as samples (individuals) and columns as genes (variables). Each eigenvector defines a principal component and their corresponding eigenvalue is proportional to its explained variance. It is the variance of the component over all the genes. The eigenvectors with large eigenvalues are the ones that contain most of the information. In our study PC1 and PC2 have an eigenvalue of approximately 2856 and 1965 and it contains the highest variance and information about the genes.

The set of seven genes that were identified in both the DGE and sPLS processes were used for gene annotation. The functional annotation of these seven genes shows their roles in inflammatory process (S100A9/A12), lipid metabolism (AKR1B10, AKR1B11) and skin barrier pathways (KRT16, PRSS27, GJB2, PI3).

- KRT16 gene regulates innate immunity by producing signals in response to any epidermal barrier breach and mutations in this gene contributes to the initiation and exacerbation of AD.
- Peptidase Inhibitor 3 (PI3) is an antimicrobial peptide and prevents elastase-mediated tissue proteolysis. It is increased during the onset of AD.
- Gap Junction Protein Beta 2(GJB2) expression increases the number of keratinocytes and causes thickness of epidermis which is observed in AD patients. Alto-keto reductase (AKR1B10 and AKR1B11) regulates keratinocyte differentiation. Alterations in AKR1B10 expression is seen in atopic dermatitis.

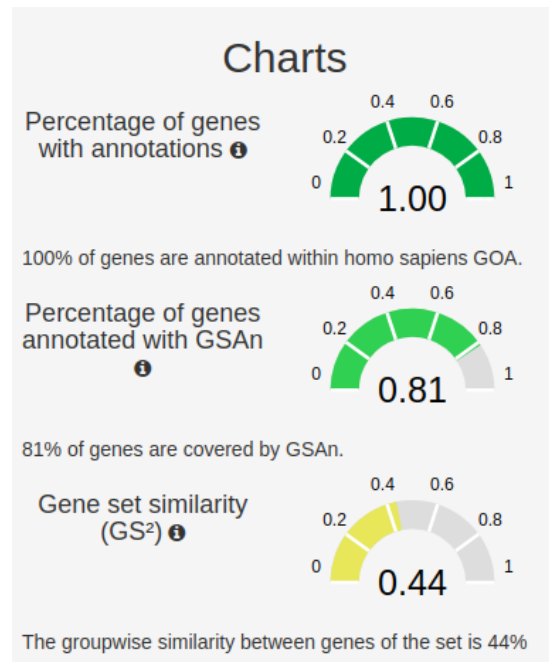


Fig. 2: Descriptive statistics of results obtained with GSAn on the 7 genes of interest.

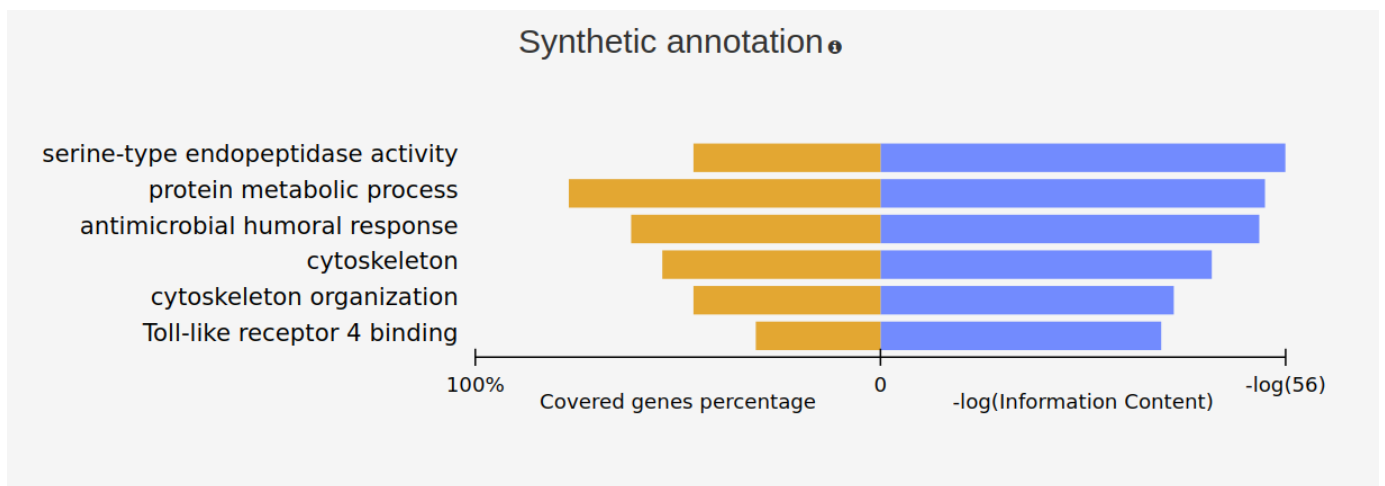


Fig. 3: Summary statistics of information and role of the genes of interest.

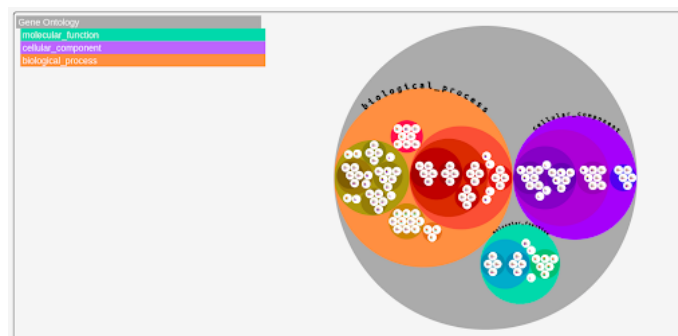


Fig. 4: Diagram showing the groupings of genetic information in the identified genes into three overarching categories: biological processes, cellular components and molecular functions.

- Expression of S100A9 causes abnormal differentiation and hyper proliferation of the epidermal cells; the upregulation of this gene can occur due to stress and certain drugs.
- S100A12 is known to exert antimicrobial activity and the skin barrier genes involved in AD enhances the expression of S100A12 by prevents it from exerting its activity.

#### 4.1. Strengths and Limitations

The current study successfully answered the research question and accepts the null hypothesis that there is no indication of genetic differentiation between lesional and non-lesional skin samples obtained from individuals with atopic dermatitis. One strength of the study is that all of the genes identified in the analysis were found in previous literature to play an important role in pathogenesis of AD.

Therefore, our results are consistent with previous findings. A second strength is that the study used two types of analysis to determine the final list of genes. Additionally, the sample groups in this study were different from each other and showed no correlation. A high number of upregulated genes in DGE analysis were identified.

One limitation of the study is that the data provided does not include people of multiple races or age groups which may affect the findings of the study. lesional and non-lesional. Also, a review of relevant literature shows that the FLG gene is a common cause of AD yet this gene did not appear in the results. Perhaps this gene is found in both lesional and non-lesional samples of individuals with AD and that is why it did not appear in the current analyses. More research is required to fully investigate.

Another limitation is computational power: leave-one-out cross validation was have been preferred in the hyper-parameter tuning step of sPLS.

## 5. CONCLUSION

The current study did not find substantial evidence to conclude that the genetic makeup of lesional and non-lesional skin samples are different between same-subject samples from individuals affected by AD. However, the genetic differences that were found may provide some insight into the functionality of healthy and lesioned tissue, mainly surrounding the inflammatory process, skin barrier pathways and lipid metabolism.

The lack of findings in the current study can lead to research questions that may be of more interest. For example, future research may want to focus on the genetic differences of atopic dermatitis as related to race, age, and allergies.

- Race: The literature suggests that there is evidence that Black individuals may be more prone to eczema than individuals of other races, while at the same time

less likely to be treated for the condition. The racial makeup (nearly 90% White) of the sample provided for analysis does not permit genetic comparisons among the races.

- AD often changes over time; either improving or worsening as the individual ages. Therefore, it would be of interest to determine how the genetic structure changes over time in individuals and how those changes impact the progression of the condition.
- Allergies: A substantial portion of the sample provided has allergies to various substances, including pollen (48.2%), food (37.8%), and animals (33.6%). This information may be of importance since people with AD are more likely to develop hay fever and are advised to avoid allergens in an effort to prevent a flare-up of symptoms. Although not the topic of the current paper, investigating the differential gene expression of AD patients with and without various allergies may be of interest.

Further research using a more diverse sample of patients is needed to understand the genetic transgression from healthy to disease state in AD.

## REFERENCES

- [1] Yoav Benjamini. Discovering the false discovery rate. *Journal of the Royal Statistical Society: series B (statistical methodology)*, 72(4):405–416, 2010.
- [2] Lianghai Bin and Donald YM Leung. Genetic and epigenetic studies of atopic dermatitis. *Allergy, Asthma & Clinical Immunology*, 12(1):52, 2016.
- [3] Christopher P Burgess, Irina Higgins, Arka Pal, Loic Matthey, Nick Watters, Guillaume Desjardins, and Alexander Lerchner. Understanding disentangling in  $\beta$ -vae. *arXiv preprint arXiv:1804.03599*, 2018.
- [4] Bradley Efron and Robert J Tibshirani. *An introduction to the bootstrap*. CRC press, 1994.
- [5] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics New York, 2001.
- [6] Emma Guttman-Yassky, Andrea Waldman, Jusleen Ahluwalia, Peck Y Ong, and Lawrence F Eichenfield. Atopic dermatitis: pathogenesis. *Semin Cutan Med Surg*, 36(3):100–103, 2017.
- [7] James A Hanley and Barbara J McNeil. The meaning and use of the area under a receiver operating characteristic (roc) curve. *Radiology*, 143(1):29–36, 1982.
- [8] Irina Higgins, Loic Matthey, Arka Pal, Christopher Burgess, Xavier Glorot, Matthew Botvinick, Shakir Mohamed, and Alexander Lerchner. beta-vae: Learning basic visual concepts with a constrained variational framework. *ICLR*, 2(5):6, 2017.
- [9] Sophie Nutten. Atopic dermatitis: global epidemiology and risk factors. *Annals of Nutrition and Metabolism*, 66(Suppl. 1):8–16, 2015.
- [10] W Peng and N Novak. Pathogenesis of atopic dermatitis. *Clinical & Experimental Allergy*, 45(3):566–574, 2015.
- [11] Gordon Smyth. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.*, 1.