

**Differential gene Expression Analysis to distinguish Psoriatic
Arthritis in patients with Cutaneous Psoriasis: A Machine Learning Approach**

Prof: Dr.Gurjit Randhawa

Jesse J Annear

145429

Conflicts of Interest:

The individual who performed this analysis has a diagnosis of psoriasis and is working with a rheumatologist on a diagnosis of Psoriatic Arthritis. This may make this analysis more prone to bias.

Terminology

PSO - Psoriasis

PSA - Psoriatic arthritis

Background:

Psoriatic arthritis (PSA) is a chronic condition that affects around 30% of individuals with psoriasis, an inflammatory skin disease[1]. It's important to note that 90% of patients with psoriatic arthritis have psoriasis and if they don't have psoriasis there is usually a family member who does, which highlights the genetic nature of psoriasis and psoriatic arthritis. While both psoriasis and PSA worsen over time, diagnosing PSA when psoriasis is present has been a challenge. A study showed that only 23% of patients were diagnosed with PSA at symptom onset, with 45% diagnosed after 2 years[2]. Early diagnosis and treatment are crucial in improving long-term outcomes and preventing joint damage and poor functional outcomes. Additionally, PSA can present similarly to other autoimmune diseases, and according to arthritis.org PSA patients have double the risk of developing cardiovascular disease, 43% more likely to have or develop heart disease and a 23% increased risk of developing conditions that affect blood flow to the brain[3]. Dna expression research has led to effective biologic medications like taltz(ixekizumab), which lower the expression of il-17 in patients with psoriasis, PSA, and other diseases. Further research is needed to better define the differences between psoriasis, PSA, and other autoimmune diseases, and to isolate biomarkers that have the potential to revolutionize medical care.

In our analysis, we are exploring the use of gene expression data to develop a machine learning model that can differentiate between psoriatic arthritis and psoriasis(only). We are utilizing differential gene expression analysis to improve the accuracy of our model by reducing the number of dimensions in the data. Developing an accurate diagnostic tool is critical as early courses of biologics have been shown to greatly reduce patient outcomes (management of psoriatic arthritis: early diagnosis, monitoring of disease severity and cutting edge therapies). If a patient can get onto suitable treatment to reduce inflammation early it can help stop aggressive bone, cartilage and tendon damage that otherwise would be irreversible. Identifying specific biomarkers associated with PSA can also help us better understand the condition and develop more targeted treatments. The potential for personalized treatment plans using isolated biomarkers could also be significant. We hope that this analysis will help show the use of gene expression on top of rheumatological diagnosis tools like CASPAR[1]. CASPAR alone may not be ideal as it typically is a lagging indicator as it can only measure inflammatory damage after it has persisted for a good amount of time. That's why gene expression analysis could aid patients in getting a more timely diagnosis and reduce their risks of developing other mental disorders like somatic symptom disorder that commonly goes along with psoriatic patients who take an average of 2.5 years[1] after symptom development to finally get their diagnosis with PSA.

Although some research has been done using differential gene expression analysis in order to identify important biological pathways, there has been little to no research done on

using machine learning models on top of differential gene expression analysis in order to detect psoriatic arthritis in patients with psoriasis. We believe that this kind of analysis may provide a solution to the problem of late diagnosis of psoriatic arthritis. We suggest that the use of an 'all encompassing model' which includes differential gene expression data of multiple cell types, patient medical history, CASPAR scores from a rheumatologist, MRI and ultrasound technology may provide a solution to the late diagnosis problem. Unfortunately, even with the recently popularized use of machine learning and big data, there has been little to no research done on the use of machine learning to aid in the diagnostic process of psoriatic arthritis.

This may be because most of the research done is focused almost exclusively on finding new biological pathways in which biologics can be made to treat and profit immensely from. Currently the biologics space is a multi-billion dollar industry, with it estimated to reach a total industry market value of 719.84 billion by 2023[4]. This can be observed as companies like Abbvie (Stock ticker ABBV) have seen tremendous growth and profits from drugs like humira which can be solidified by constant stock growth due to increasing profit margins and a 'wall of patents' which allowed ABBV to have 165 patents on a single drug, effectively pushing back the opportunity for biosimilars in an effort to 'game' the biologics industry[5]. This isn't to say that biologics aren't important because they are some of the greatest technology advancements that the medical field has observed, but it's important to encourage more research around early diagnosis so patients can start biologics treatments earlier and have improved patient outcomes and reduced health care costs.

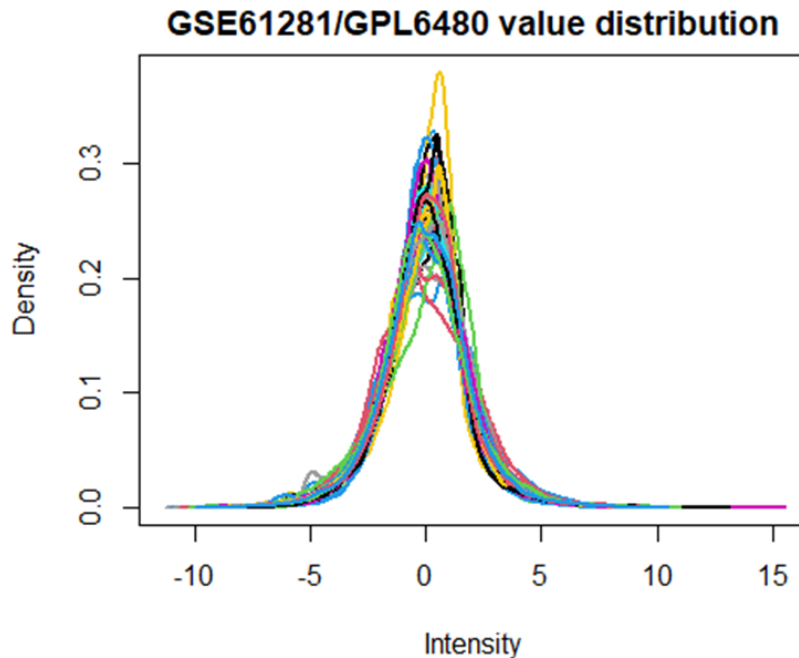
Dataset:

The dataset is public and was posted on September 24, 2014 as a GEO dataset published on the NIH website under GSE61281[6]. It involves human whole blood samples and compares transcriptional profiles between psoriatic arthritis (PSA), cutaneous psoriasis without arthritis (PSO), and unaffected controls. The experiment was conducted using expression profiling by array. The overall design includes three conditions: PSA, PSO, and unaffected controls. There were 20 biological replicates each for PSA and PSO and 12 for controls. Although sampling of inflammatory cells would be more ideal, for example synovial fluid cells as this is a primary area of disease activity, blood cells could provide more information on the high comorbidity rate that is apparent in both psoriasis and psoriatic arthritis patients with cardiovascular diseases and diabetes. On top of an increased rate when compared to controls of these diseases it has also been shown that smoking can increase the chance of developing both of these diseases[7]. So, it is definitely of interest for more studies to be down with cardiovascular specific cells like blood, or even lung cells.

Data exploration:

Initially, we conducted an exploratory analysis to identify any significant observations using various visualization techniques. A scatter plot was generated to illustrate the distribution of gene expression values, which indicated a range between -7 and 7. Based on our observations, it appears that the data had been transformed using a $\log_2(x)$ method. Although it is not always advisable to assume data transformation, $\log_2(x)$ is a common method used in gene expression data analysis, and the range and distribution of our dataset suggests that it has been log transformed. This is beneficial since it eliminates the need for data scaling and normalization. Furthermore, a density plot was constructed to assess the distribution of

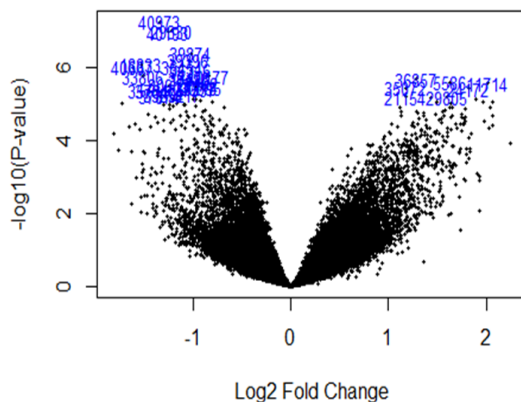
individual patient data, which revealed a relatively consistent distribution amongst patients. This suggests that the data had been previously scaled and normalized before the data was published.



Finding differentially expressed genes:

Differential gene analysis is a powerful tool for gaining insights into the molecular basis of various diseases and disorders, and for identifying genes that are dysregulated in affected individuals. In this study, we utilized the R package limma to isolate genes that were significantly upregulated or downregulated in psoriasis patients compared to controls. We employed the $\log_2(\text{fold change})$ method to identify significant genes, which is a widely-used technique for detecting differentially expressed genes. We generated a volcano plot to visualize the

distribution of our data and identified the top 30 genes that were dysregulated in psoriasis patients.



by the heterogeneity of whole blood. However, they also noted that many of the differentially expressed genes between psoriatic arthritis and psoriasis without arthritis had smaller changes in expression in PSO patients compared to controls, suggesting that the more severe phenotype of PSA is an exacerbation of small gene expression changes that occur in PSO.

This finding raises an interesting question about the relationship between psoriasis and psoriatic arthritis. It is possible that psoriatic disease represents a spectrum of disease severity rather than discrete entities, and that the dysregulated genes that we failed to identify in psoriasis patients are still differentially expressed but to a lesser extent. Indeed, some patients with psoriatic arthritis do not show signs of psoriasis, and it is possible that many of these cases go undiagnosed by rheumatologists who do not consider the diagnosis without skin symptoms.

In conclusion, our analysis using limma failed to identify any significantly dysregulated genes in psoriasis patients compared to controls. This unexpected result raises interesting questions about the relationship between psoriasis and psoriatic arthritis and highlights the need for further studies to elucidate the molecular mechanisms underlying these disorders.

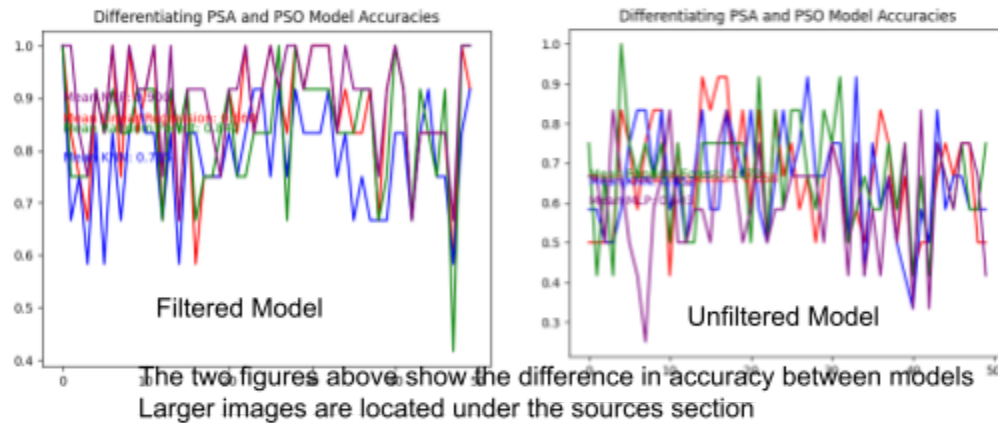
Enrichment Analysis:

Following the identification of significantly expressed genes, the next step in gene expression analysis is to perform enrichment/pathway analysis. Enrichment analysis, also known as gene list analysis, involves comparing and modeling an entire gene list or just the significant genes against a database of various biological pathways and important genetic information. The purpose of this step is to gain a deeper understanding of the biological processes involved in the significant genes or the ordered list of genes sorted by their p-values. In this study, we utilized g: profiler to perform enrichment analysis, and identified several interesting findings related to metabolic disorders that are highly comorbid with psoriasis and psoriatic arthritis. A bubble chart comparing the different categories was generated and revealed significant biologic processes and pathways related to protein kinase camp-activated catalytic subunit alpha, which has been implicated in the development of new target drugs for type 2 diabetes, another common comorbidity associated with psoriatic disease. However, parsing through this vast and complex information requires significant time and biological expertise. Therefore, careful examination and interpretation of the data is necessary to fully appreciate the important insights provided by enrichment analysis.

Modeling phase:

We created a model to differentiate between psoriatic arthritis and psoriasis(only) by identifying significant genes using gene expression data. Since we weren't interested in differentiating between psoriasis and controls we left the controls portion of the data out. In the modeling phase, we tested various machine learning models such as logistic regression, KNN classifier, random forest classifier, and MLP neural network. We found that the MLP neural network had the highest accuracy compared to other models, which was further improved by isolating only the significantly expressed genes. The accuracy of the logistic regression model was measured to be an average of 90%. These results are much better than our model where we did not filter out genes that weren't significant as the MLP only had an accuracy of 67% in that case which is 23% less accurate. Therefore, this shows that differential gene expression analysis is a suitable technique to reduce the dimensionality in our gene expression models

where the data has a high amount of dimensionality and a low amount of data. Our results suggest that the MLP neural network is a suitable method for modeling and classifying psoriatic arthritis in patients with psoriasis. However, to further improve on model efficacy it would be beneficial to combine with information obtained by health care specialists like CASPAR, MRI'S, ultrasounds of tendons and more medical history in order to find a more complete model that can better diagnose psoriatic arthritis in patients with psoriasis. This type of 'all encompassing model' has not been suggested as far as we are concerned and it would likely reduce the time it takes to receive a diagnosis in patients with psoriatic arthritis.



Discussion:

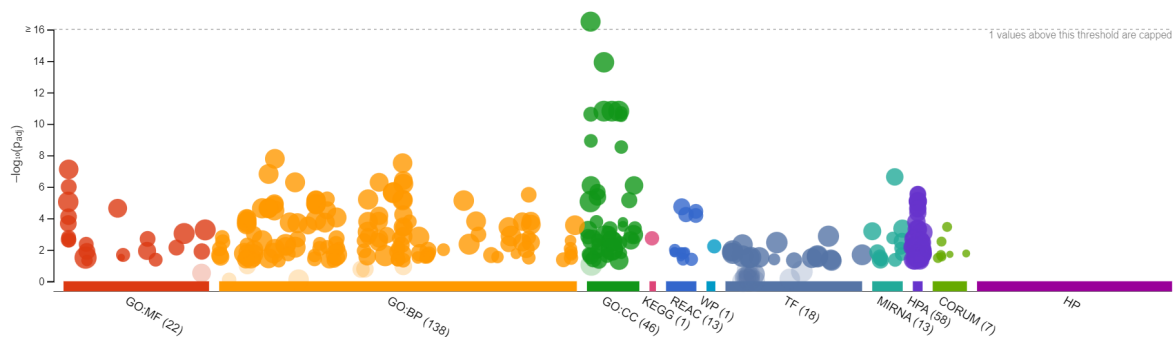
In conclusion, this study aimed to use gene expression data to develop a machine learning model that can differentiate between patients with only psoriasis and those with psoriatic arthritis as well. Early diagnosis of PSA in patients with psoriasis is crucial for improving long-term outcomes and preventing joint damage and poor functional outcomes. The dataset used in this study involved human whole blood samples and compared transcriptional profiles between PSA, cutaneous psoriasis without arthritis (PSO), and unaffected controls. However, the analysis did not identify any genes that were significantly dysregulated in psoriasis patients compared to controls. This finding is unexpected given the increased risk of cardiovascular and metabolic disorders in psoriasis patients. The potential use of isolated biomarkers for personalized treatment plans could be significant, and further research is needed to better define the differences between psoriasis, PSA, and other autoimmune diseases. Although sampling of inflammatory cells would be more ideal, blood cells could provide more information on the high comorbidity rate that is apparent in both psoriasis and PSA patients with cardiovascular diseases and diabetes. Moreover, using an MLP neural network model after identifying significantly expressed genes could potentially aid in the diagnosis of PSA in patients with psoriasis. The model could be used as a tool to predict the likelihood of a patient having PSA based on their gene expression profile. This could potentially help healthcare professionals make more accurate and timely diagnosis, leading to better treatment outcomes for patients.

However, further research is needed to validate the use of such models in clinical settings and likely a larger variety of cell types on top of tests like caspar that a rheumatologist could perform. Further research is needed in order to provide an earlier diagnosis, one where a significant amount of joint damage has occurred. Finding a way to get an earlier diagnosis not only improve PSA patient outcomes but would also lower the burden of a PSA patient on the

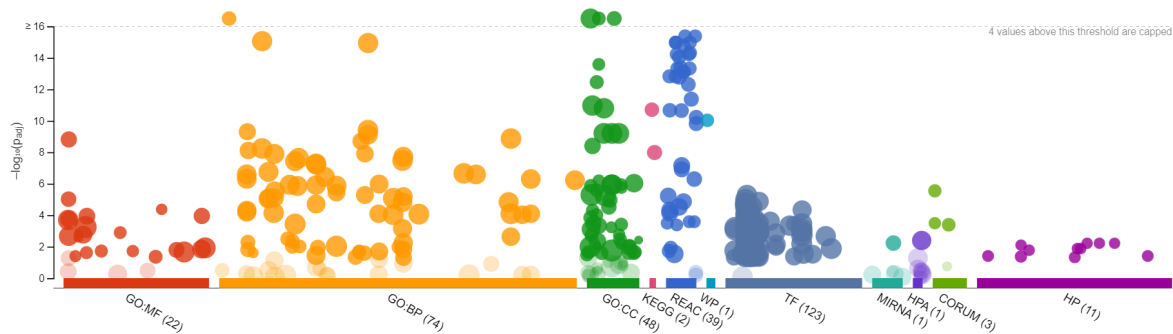
healthcare system by reducing the chance of needing surgery and by having less of a chance to develop comorbidities like metabolic disease, somatic symptom disorder, depression and cardiovascular disease. We suggest that the use of an 'all encompassing model' which includes differential gene expression data of multiple cell types, patient medical history, CASPAR scores from a rheumatologist, MRI and ultrasound technology may provide a solution to the late diagnosis problem. Unfortunately, even with the recently popularized use of machine learning and big data, there has been little to no research done on the use of machine learning to aid in the diagnostic process of psoriatic arthritis and we may be the first to suggest such a model to aid in the diagnosing of psoriatic arthritis.

Bonus Enrichment Analysis

> PSO



> PSA



version
date
organism

e107_eg54_p17_bf42210
2022-12-08, 7:27:21 p.m.
hsapiens

g:Profiler

GO:BP		PSA		pso	
Term Name	Term ID		p_adj		p_adj
cytoplasmic translation	GO:0002181		1.329×10 ⁻¹⁷		1.409×10 ⁻¹
macromolecule biosynthetic process	GO:0009059		4.632×10 ⁻¹⁶		1.717×10 ⁻⁵
cellular nitrogen compound biosynthetic process	GO:0044271		6.122×10 ⁻¹⁶		4.730×10 ⁻⁶
cellular biosynthetic process	GO:0044249		2.423×10 ⁻¹⁰		6.159×10 ⁻⁴
translation	GO:0006412		5.468×10 ⁻¹⁰		3.816×10 ⁻²
organic substance biosynthetic process	GO:1901576		8.430×10 ⁻¹⁰		1.102×10 ⁻³
cellular macromolecule metabolic process	GO:0044260		9.203×10 ⁻¹⁰		1.007×10 ⁻⁴
peptide biosynthetic process	GO:0043043		2.186×10 ⁻⁹		6.591×10 ⁻²
biosynthetic process	GO:0009058		3.629×10 ⁻⁹		2.181×10 ⁻³
regulation of macromolecule biosynthetic process	GO:0010556		7.207×10 ⁻⁹		2.162×10 ⁻⁵
peptide metabolic process	GO:0006518		8.495×10 ⁻⁹		8.445×10 ⁻¹
positive regulation of macromolecule metabolic process	GO:0010604		1.766×10 ⁻⁶		9.685×10 ⁻⁹
regulation of RNA metabolic process	GO:0051252		1.084×10 ⁻⁸		4.145×10 ⁻⁷
amide biosynthetic process	GO:0043604		1.363×10 ⁻⁸		6.778×10 ⁻²
positive regulation of nitrogen compound metabolic process	GO:0051173		1.716×10 ⁻⁵		1.685×10 ⁻⁸
regulation of nucleobase-containing compound metabolic ...	GO:0019219		2.146×10 ⁻⁸		3.479×10 ⁻⁷
regulation of nitrogen compound metabolic process	GO:0051171		2.158×10 ⁻⁸		3.605×10 ⁻⁶
protein metabolic process	GO:0019538		3.016×10 ⁻⁸		6.648×10 ⁻¹
regulation of cellular metabolic process	GO:0031323		3.428×10 ⁻⁸		1.053×10 ⁻⁴
regulation of cellular biosynthetic process	GO:0031326		3.471×10 ⁻⁸		8.250×10 ⁻⁶
positive regulation of metabolic process	GO:0009893		5.099×10 ⁻⁶		9.482×10 ⁻⁸
regulation of biosynthetic process	GO:0009889		1.045×10 ⁻⁷		1.595×10 ⁻⁵
regulation of primary metabolic process	GO:0080090		1.420×10 ⁻⁷		3.965×10 ⁻⁶
DNA-templated transcription	GO:0006351		1.426×10 ⁻⁷		2.412×10 ⁻⁴
nucleic acid-templated transcription	GO:0097659		1.454×10 ⁻⁷		2.437×10 ⁻⁴
positive regulation of RNA metabolic process	GO:0051254		7.220×10 ⁻⁶		1.695×10 ⁻⁷
RNA biosynthetic process	GO:0032774		2.078×10 ⁻⁷		2.969×10 ⁻⁴
positive regulation of nucleobase-containing compound m...	GO:0045935		1.173×10 ⁻⁶		2.354×10 ⁻⁷
regulation of DNA-templated transcription	GO:0006355		2.794×10 ⁻⁷		3.834×10 ⁻⁴
regulation of nucleic acid-templated transcription	GO:1903506		2.847×10 ⁻⁷		3.874×10 ⁻⁴
regulation of RNA biosynthetic process	GO:2001141		3.374×10 ⁻⁷		4.255×10 ⁻⁴
regulation of mRNA metabolic process	GO:1903311		1.000		5.086×10 ⁻⁷
heterocycle biosynthetic process	GO:0018130		6.956×10 ⁻⁷		2.988×10 ⁻⁴
nucleobase-containing compound biosynthetic process	GO:0034654		7.969×10 ⁻⁷		1.451×10 ⁻⁴
aromatic compound biosynthetic process	GO:0019438		8.388×10 ⁻⁷		3.324×10 ⁻⁴
positive regulation of cellular process	GO:0048522		5.514×10 ⁻⁶		1.150×10 ⁻⁶
positive regulation of biological process	GO:0048518		6.511×10 ⁻⁵		1.261×10 ⁻⁶
positive regulation of cellular metabolic process	GO:0031325		1.414×10 ⁻⁶		4.213×10 ⁻⁶
negative regulation of nitrogen compound metabolic process	GO:0051172		3.353×10 ⁻⁴		3.383×10 ⁻⁶

1 to 39 of 39 < > Page 1 of 1 > >|

GO:CC		PSA		pso	
Term Name	Term ID		p_adj		p_adj
nucleoplasm	GO:0005654		1.001×10 ⁻¹⁹		7.499×10 ⁻¹⁹
cytosolic ribosome	GO:0022626		2.756×10 ⁻¹⁹		1.000
ribosomal subunit	GO:0044391		3.759×10 ⁻¹⁹		1.000
nuclear lumen	GO:0031981		9.747×10 ⁻¹²		9.113×10 ⁻¹⁵
cytosolic large ribosomal subunit	GO:0022625		2.914×10 ⁻¹⁴		1.000
large ribosomal subunit	GO:0015934		3.788×10 ⁻¹³		1.000
spliceosomal complex	GO:0005681		1.538×10 ⁻⁴		2.118×10 ⁻¹²
organelle lumen	GO:0043233		4.181×10 ⁻¹⁰		1.325×10 ⁻¹¹
membrane-enclosed lumen	GO:0031974		4.181×10 ⁻¹⁰		1.325×10 ⁻¹¹
intracellular organelle lumen	GO:0070013		4.181×10 ⁻¹⁰		1.325×10 ⁻¹¹
cytosol	GO:0005829		1.592×10 ⁻¹¹		1.804×10 ⁻³
U2-type precatalytic spliceosome	GO:0071005		1.179×10 ⁻⁶		2.159×10 ⁻¹¹
precatalytic spliceosome	GO:0071011		1.838×10 ⁻⁶		3.447×10 ⁻¹¹
catalytic step 2 spliceosome	GO:0071013		4.039×10 ⁻⁴		1.535×10 ⁻¹⁰
U2-type spliceosomal complex	GO:0005684		1.437×10 ⁻⁶		1.377×10 ⁻⁹
ribosome	GO:0005840		4.644×10 ⁻⁹		1.000
catalytic complex	GO:1902494		1.042×10 ⁻⁶		3.339×10 ⁻⁷
cytosolic small ribosomal subunit	GO:0022627		6.711×10 ⁻⁷		1.000

1 to 18 of 18 < > Page 1 of 1 > >|

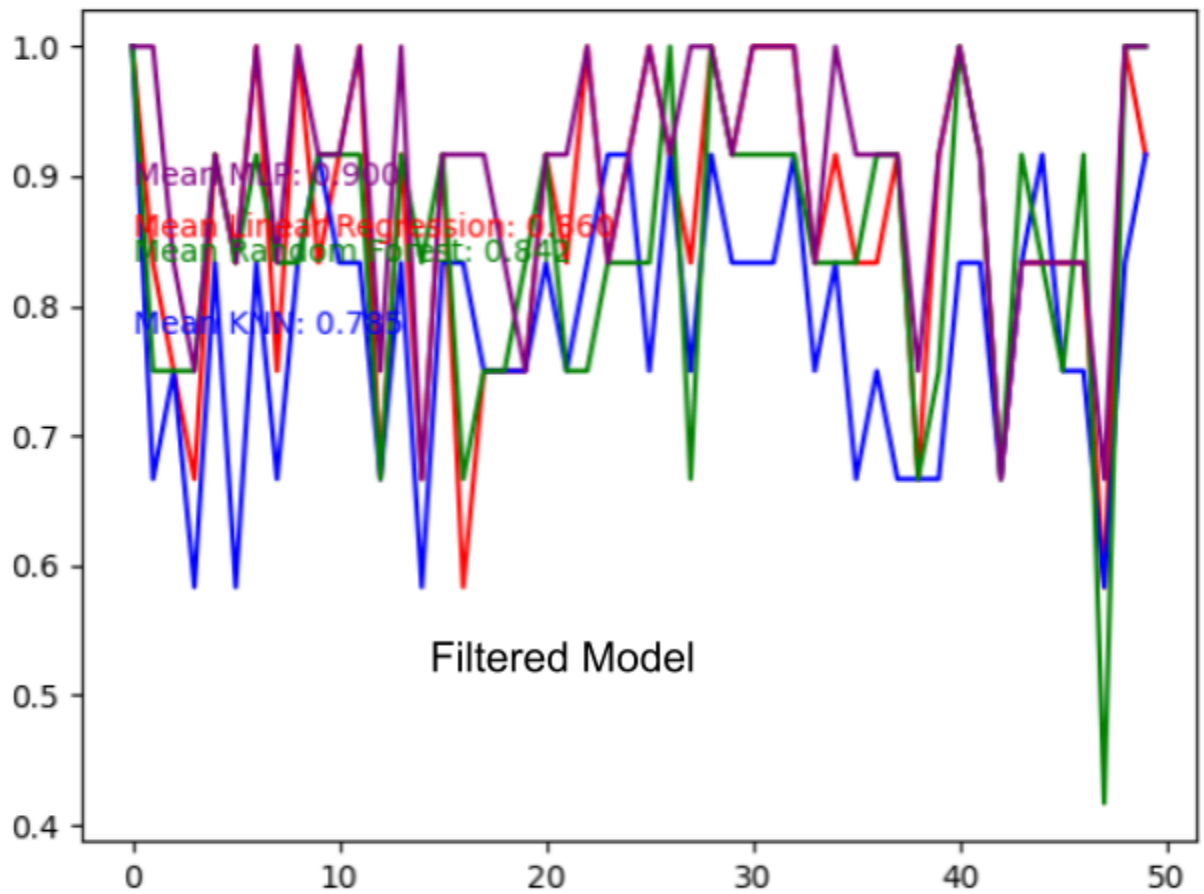
GO:MF		PSA		PSO	
Term Name	Term ID		p_adj		p_adj
structural constituent of ribosome	GO:0003735		1.538×10 ⁻⁹		1.000
RNA binding	GO:0003723		2.165×10 ⁻³		7.685×10 ⁻⁸
mRNA binding	GO:0003729		5.425×10 ⁻²		1.006×10 ⁻⁶
nucleic acid binding	GO:0003676		2.009×10 ⁻⁴		9.039×10 ⁻⁶
transcription coactivator activity	GO:0003713		9.616×10 ⁻⁶		2.276×10 ⁻³
enzyme binding	GO:0019899		5.531×10 ⁻¹		2.300×10 ⁻⁵
mRNA 5'-UTR binding	GO:0048027		4.265×10 ⁻⁵		1.000
transcription coregulator activity	GO:0003712		1.782×10 ⁻⁴		7.904×10 ⁻⁵
DNA-binding transcription factor binding	GO:0140297		1.117×10 ⁻⁴		1.282×10 ⁻²
transcription factor binding	GO:0008134		1.152×10 ⁻⁴		1.263×10 ⁻²
chromatin binding	GO:0003682		3.863×10 ⁻¹		2.153×10 ⁻⁴
heterocyclic compound binding	GO:1901363		1.196×10 ⁻²		5.451×10 ⁻⁴
protein binding	GO:0005515		5.696×10 ⁻⁴		3.166×10 ⁻²
organic cyclic compound binding	GO:0097159		2.183×10 ⁻²		9.270×10 ⁻⁴
nuclear receptor coactivator activity	GO:0030374		1.277×10 ⁻³		1.000
cAMP-dependent protein kinase activity	GO:0004691		1.587×10 ⁻³		1.000
transcription corepressor activity	GO:0003714		1.000		1.633×10 ⁻³
structural molecule activity	GO:0005198		1.735×10 ⁻³		1.000
cadherin binding	GO:0045296		3.313×10 ⁻¹		2.034×10 ⁻³
mRNA 3'-UTR binding	GO:0003730		1.000		2.605×10 ⁻³
cyclic nucleotide-dependent protein kinase activity	GO:0004690		2.720×10 ⁻³		1.000
protein C-terminus binding	GO:0008022		1.000		4.451×10 ⁻³
RNA polymerase II-specific DNA-binding transcription facto...	GO:0061629		1.615×10 ⁻²		7.239×10 ⁻³
protein-containing complex binding	GO:0044877		1.000		1.184×10 ⁻²
transcription regulator activity	GO:0140110		1.434×10 ⁻²		3.009×10 ⁻¹
oxidoreduction-driven active transmembrane transporter ac...	GO:0015453		1.909×10 ⁻²		1.000
peptide N-acetyltransferase activity	GO:0034212		1.916×10 ⁻²		1.000
chromatin DNA binding	GO:0031490		1.000		2.123×10 ⁻²
N-acetyltransferase activity	GO:0008080		2.419×10 ⁻²		1.000
US snRNA binding	GO:0030623		1.000		3.056×10 ⁻²
histone acetyltransferase activity	GO:0004402		3.993×10 ⁻²		1.000
protein N-terminus binding	GO:0047485		4.411×10 ⁻²		4.253×10 ⁻²
poly(A) binding	GO:0008143		1.000		4.778×10 ⁻²

1 to 33 of 33 | < > Page 1 of 1 > >|

Sources

1. [Psoriatic Arthritis | NEJM](#)
2. [Diagnostic Delay in Psoriatic Arthritis: A Population-based Study | The Journal of Rheumatology \(jrheum.org\)](#)
3. [Psoriatic Arthritis and Your Heart](#)
4. [Biologics Market Size to Worth Around US\\$ 719.84 Bn by \(globenewswire.com\)](#)
5. [How a Drug Company Made \\$114 Billion by Gaming the U.S. Patent System - The New York Times \(nytimes.com\)](#)
6. [GEO Accession viewer \(nih.gov\)](#)
7. [The relationship between smoking, psoriasis and psoriatic arthritis: Expert Review of Clinical Immunology: Vol 15, No 1 \(tandfonline.com\)](#)

Differentiating PSA and PSO Model Accuracies



Differentiating PSA and PSO Model Accuracies

