



مقدمه‌ای بر بیوانفورماتیک
نیم‌سال اول ۰۵-۰۴
دکتر شریفی - کوهی

فراز دعاگوی تهرانی ۴۰۲۱۰۵۹۹۸

گزارش تمرین عملی سوم

سوالات بخش اول

۱. زیرا که داده های ماتریس سریال پردازش شده و یا خلاصه است و مناسب نیست.
۲. تعداد سطر ها تعداد ژن ها و ستون ها تعداد نمونه هاست که به ترتیب ۲۸۸۸۹ و ۳۶ تا اند.
۳. به دلیل اینکه rRNA ها معمولاً حذف میشوند و به خوبی توالی یابی نمیشوند.
۴. به دلیل گسترده بودن دامنه و غیر صحیح بودن آن، normalized هستند.
۵. توجه نکردن به آن، میتواند باعث بایاس شدن خروجی شود و به یک نمونه بیشتر توجه کرد.

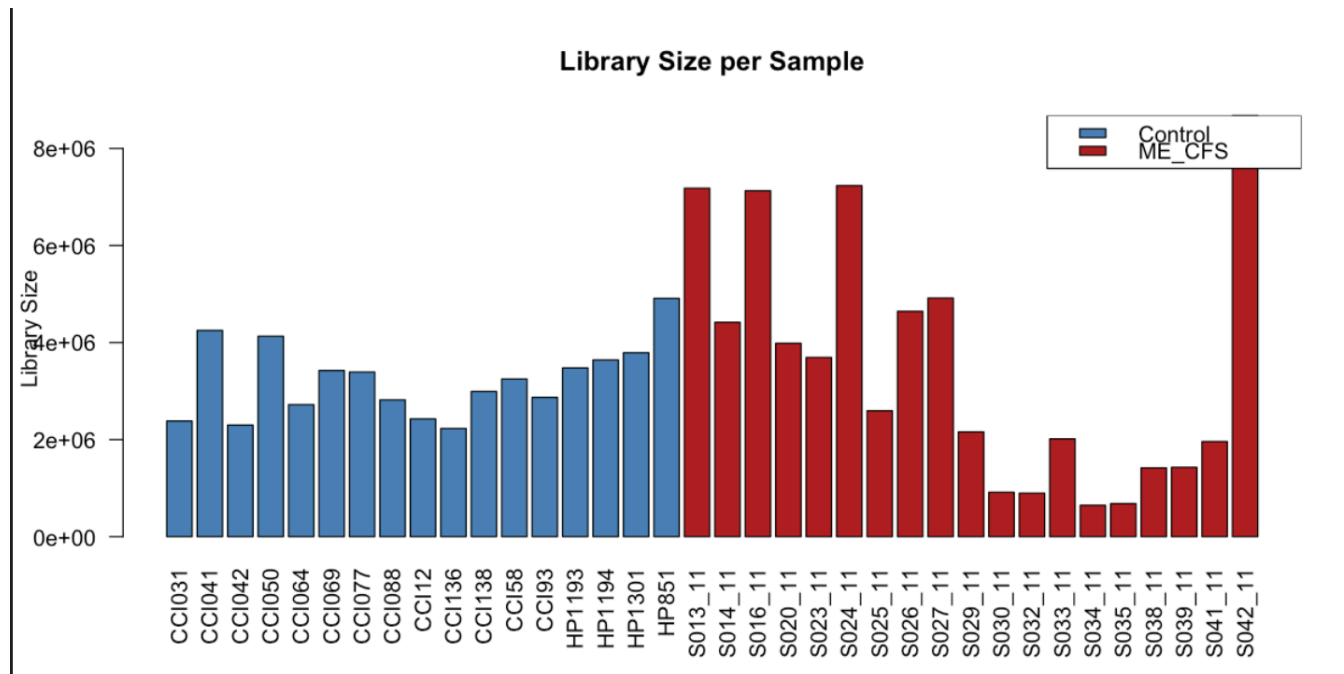
سوالات بخش دوم

۱. ستون tile بیشترین کمک را کرد چرا که در آن مشخص بود که نوع health control بوده یا برای بیماری.
۲. از گروه control ۱۷ تا و از دیگری ۱۹ تا داریم که تقریباً برابرند و تعادل خوبی را نشان میدهد.
۳. اگر یکی نباشند تطابق و انتصاب نمونه ها به گروه ها درست انجام نمیشود و ادامه تحلیل دچار خطا میشود.
۴. برخی نتایج برعکس میشوند و پیش بینی ها را خراب میکنند.
۵. معمولاً سطح مرجع Control انتخاب می شود تا logFC بیان تغییرات بین بیمار و سالم محاسبه شود. انتخاب Con-trol به عنوان مرجع باعث می شود مقدار مثبت logFC نشان دهنده افزایش بیان در بیمار نسبت به سالم باشد و تفسیر ساده و مستقیم داشته باشد.

سوالات بخش سوم

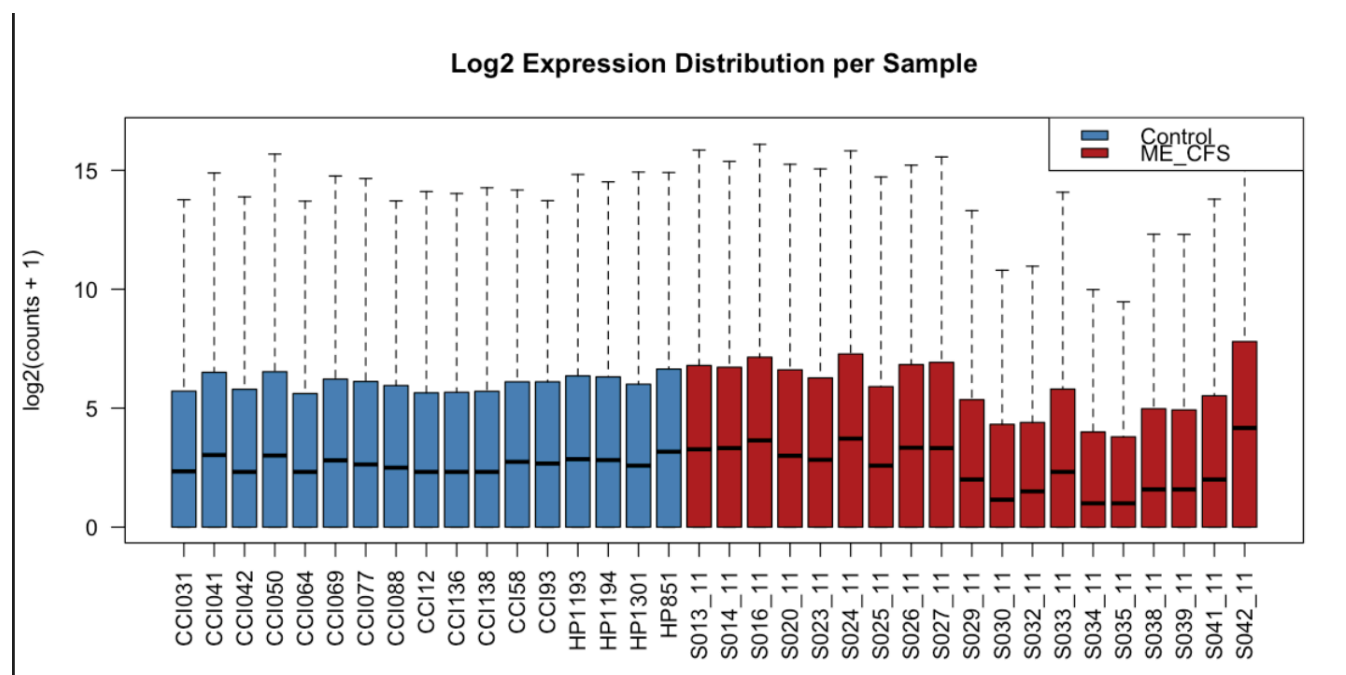
۱. با توجه به خروجی اولین سلول این بخش، کمترین مقدار برابر ۰، و بیشترین مقدار برابر ۷۲.۹۵۷۲۶ میباشد. همچنین میانه و میانگین به ترتیب برابر ۶۲.۴ و ۹۵.۱۱۶ میباشد. نبود اعداد منفی و وجود اعشار نشان میدهد که داده ها شمارشی خام نیستند و نرمال شده هستند.

۲. با توجه به نمودار و خروجی بخش‌های قبلی، کمترین library size مربوط به S034_11 با مقدار ۶۴۷۲۸۸/۵ و بیشترین مقدار مربوط به S042_11 با مقدار ۸۶۷۸۴۶۹ می‌باشد. با توجه به خروجی‌های سلول دوم این بخش، تفاوت چشم‌گیری بین میانگین‌ها وجود ندارد اما مقدار ماکسیمم در صورت اختلاف زیاد، به سمت یکی بایاس پیدا می‌کنیم و در ادامه پیش‌بینی‌هایمان درست نمی‌شود.



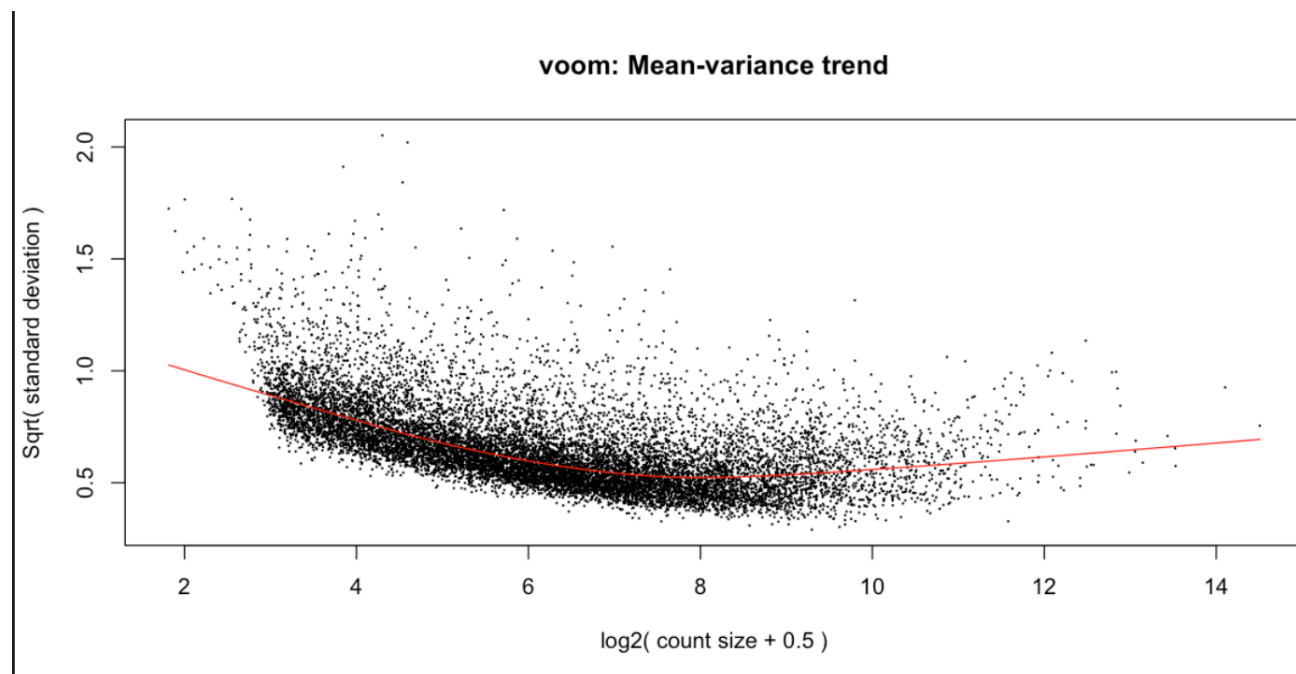
شکل ۱: نمودار library size ها به تفکیک گروه

۳. در کل، همانطور که در نمودار پایین می‌بینید، تقریباً همه مقادیر میانه و کلیت آنها با هم برابرند هرچند که برخی نمونه‌ها مثل راست‌ترین نمونه بزرگتر از بقیه و برخی دیگر نیز کوچک‌ترند. دلیل این تفاوت می‌تواند در کیفیت بد، داده پرت یا library preparation متفاوت باشد.



شکل ۲: بیان لگاریتمی به تفکیک گروه

۴. با کاهش میانگین بیان، واریانس زیاد میشود و با افزایش آن، پایدار تر میشود. در این مدل واریانس های بیشتر وزن کمتری میگیرند تا پایداری حفظ شود.



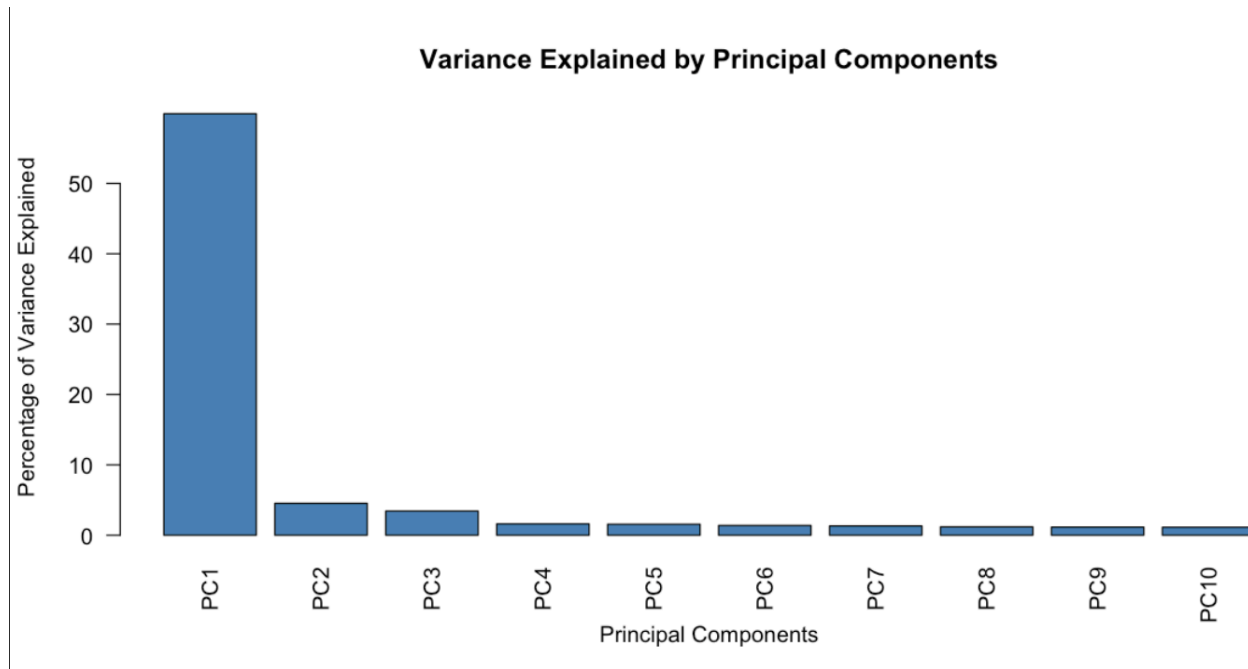
شکل ۳: نمودار voom

۵. در روش DESeq2 فرض میشود که داده ها صحیح اند اما در limma-voom نیازی به این فرض نیست و بر داده های

تقریباً نرمال هم درست عمل میکند.

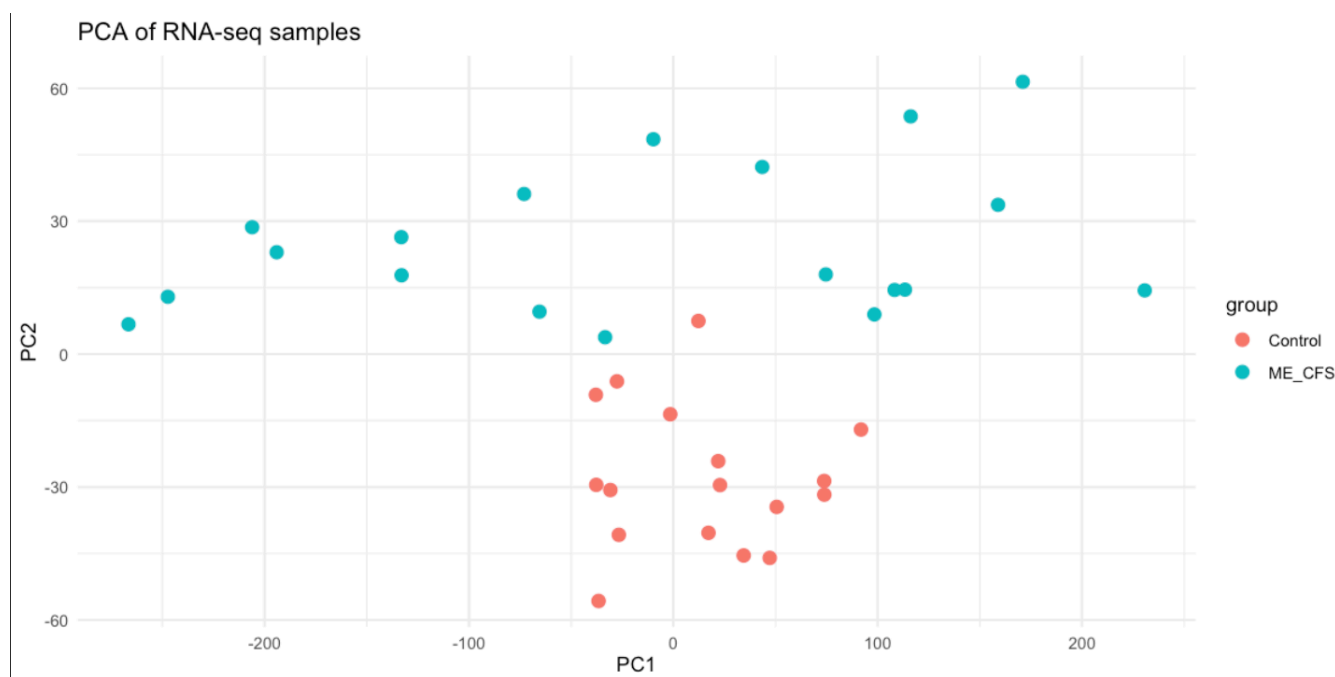
سوالات بخش چهارم

۱. بنا به خروجی سلول سوم این بخش، حدود ۶۰ درصد واریانس از PC ۱ که این مسئله، نشان دهنده غلبه این بعد است و آن را توضیح میدهد.



شکل ۴: درصد واریانس توضیح داده شده توسط اولین ۱۰ PC

۲. در این نمودار، داده‌ها تقریباً به صورت کامل با یک خط جدا میشوند که این مسئله بیانگر این است که داده‌های این دو گروه تفاوت‌های معناداری دارند و در فضای کم‌بعد، قابل تفکیک اند.



شکل ۵: داده‌ها در فضای PC1-PC2

۳. در نمودار بالا، نقاطی که فاصله بسیار زیادی به دسته خود دارند، مانند یک نقطه ای که قرمز است ولی در بخش آبی قرار گرفته است، به عنوان داده پرت در نظر گرفته میشود. دلیل این تفاوت میتواند موارد گوناگونی از جمله کیفیت پایین نمونه، خطا در sequencing، آلودگی نمونه، وضعیت بالینی خاص و یا دیگر موارد باشد.

۴. زیرا این نمونه‌ها، اطلاعات خاصی را به ما اضافه نمیکند و در هر کجا مقدار نسبتاً ثابتی دارند و وجود آنها تنها باعث نویز میشود.

۵. اگر از $scale = TRUE$ استفاده میکردیم، ژن‌های کم‌واریانس باقی مانده که بیشتر آنها نویز هستند، وزن بیشتری میگرفتند و موثرتر واقع میشدند که باعث ایراد در تصمیم‌گیری ما میشد.

۶. ژن‌ها با بیان بالاتر کاملاً غالب میشدند و کل ساختار با به خود بایاس میکردند. با استفاده از \log این داده‌ها متعادل‌تر میشوند و استفاده از آنها منطقی‌تر میشود.

سوالات بخش پنجم

۱. ژن‌هایی که بیان کمی دارند و اطلاعات کافی ندارند را حذف میکند. چرا که به پیشبرد و دسته‌بندی کمکی نمیکند و تنها باعث نویز میشوند.

۲. این نمودار که در بخش ۳ به آن پرداختیم، نشان میدهد که با افزایش بیان، واریانس کاهش میابد. وزن دهی بر همین اساس در این روش انجام میگیرد و ضرورت آن به دلیل برقرار ماندن فرض همسانی واریانس میباشد.

۳. با توجه به خروجی سلول سوم این بخش که در پایین میتوانید آن را ببینید، این سه ژن بیشترین معنا را دارند. همچنین، هرچقدر $\log FC$ بیان بیشتری در گروه ME_CFS نسبت به کنترل داشته باشد، مثبت‌تر میشود و هرچقدر بیان در دیگری بیشتر باشد، منفی‌تر میشود.

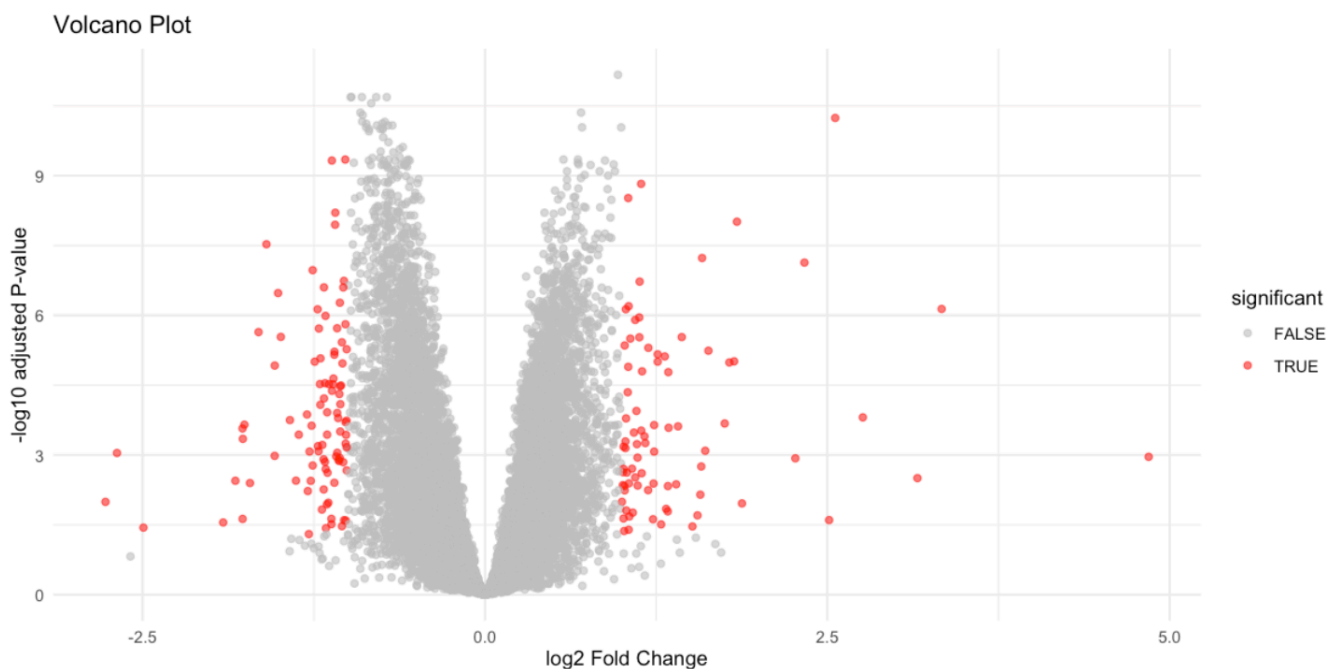
در نتیجه در خروجی پایین اولی در ME_CFS بیان بیشتری دارد و دوتای دیگر در Control.

A data.frame: 3 × 2

	logFC	adj.P.Val
	<dbl>	<dbl>
TAF1D	0.9711999	6.822068e-12
SCAF1	-0.9807698	2.057448e-11
MAP1S	-0.8988431	2.057448e-11

شکل ۶: سه ژن برتر بر اساس adj.P.Val

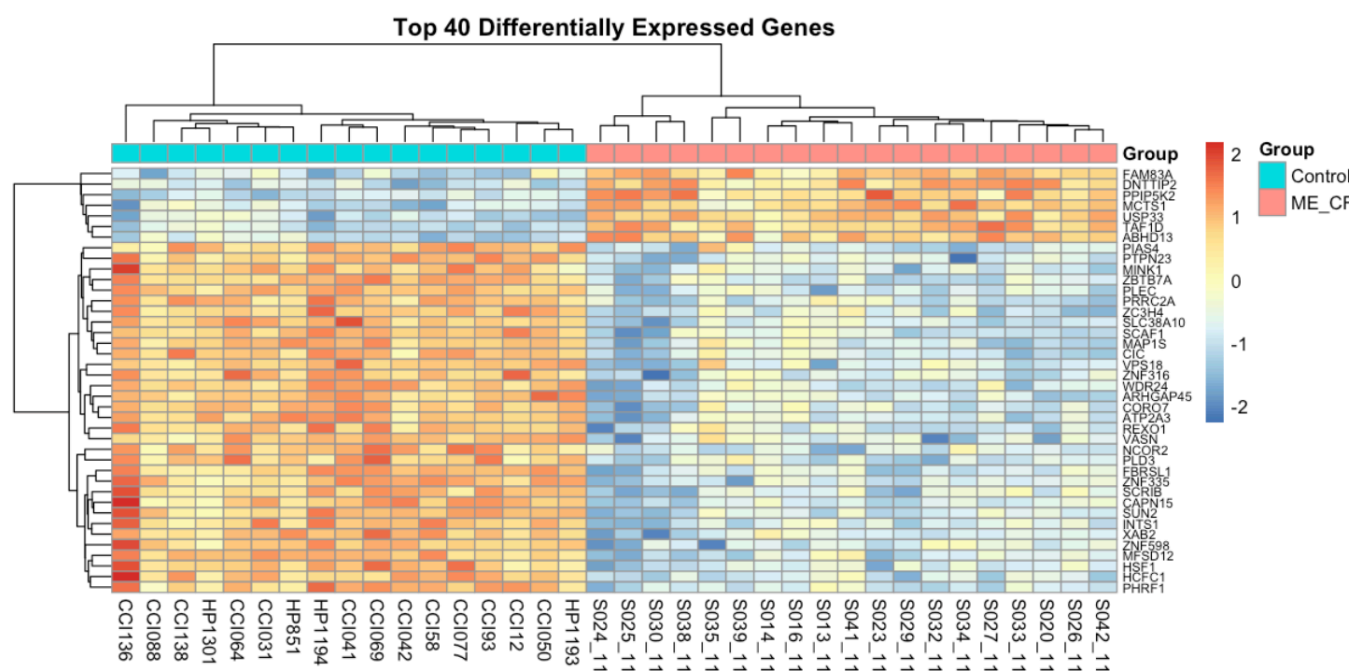
۴. بر اساس این نمودار، بخش‌های سمت چپ محور عمودی که logFC منفی دارند به گروه Control و مثبت‌ها به ME_CFS تعلق دارند. هرچقدر به محور عمودی نزدیک‌تر باشند، از معنای کمتری برخوردارند. هرچقدر مقدار |logFC| بیشتر می‌شود و مقدار adj.P.Val کمتر می‌شود، نقاط از اهمیت بیشتری برخوردارند و در نمودار زیر با رنگ قرمز مشخص شده‌اند. همچنین با توجه به خروجی سلول ۵ این بخش، ۳۳۰۸ تا از ژن‌ها افزایش بیان و ۳۲۹۱ تا کاهش بیان دارند. در نتیجه تقریباً برابرند.



شکل ۷: نمودار volcano

۵. با توجه به خروجی سلول چهارم این بخش، ۱۷۴ ژن این ویژگی را دارند و آن‌ها ژن‌هایی هستند که هم مقدار تغییر زیاد دارند و هم معناداری آماری کافی دارند.

۶. به همانطور که در شکل زیر می‌توانید ببینید نمونه‌های دو دسته در خوشه‌های صحیح قرار گرفتند. اگر نمونه‌ای اشتباه می‌شد، احتمالاً ناشی از خطای بیولوژیکی یا خطای تکنیکی و یا داده پرت بود.



شکل ۸: نمودار volcano

۷. همانطور که بالاتر گفته شد، همچنین، هرچقدر $\log FC$ بیان بیشتری در گروه ME_CFS نسبت به کنترل داشته باشد، مثبت تر میشود و هرچقدر بیان در دیگری بیشتر باشد، منفی تر میشود. در این مدل تفسیر ME_CFS vs. Control اختلاف آنها میشود.