



## مسئله ۱. نرمال سازی ( ۵۰ نمره)

در این بخش، می‌خواهیم به مقایسه روش‌های مختلف نرمال سازی داده‌های RNA-seq بپردازیم. در ضمیمه صورت سوالات دو فایل GSE60450\_Lactation-GenewiseCounts.txt (که شامل اطلاعات تعداد readها و طول ژن برای هر نمونه است) و SampleInfo.txt (که شامل اطلاعات مربوط به گروه هر نمونه است) ارائه گردیده. این مجموعه داده شامل ۶ گروه ۲ تایی از نمونه‌هاست.

الف) شما باید داده‌های ارائه شده را به کمک سه روش TPM، DeSeq2 و نرمال سازی کرده و سپس به کمک روش‌های PCA و خوشه‌بندی سلسله‌مراتبی، اقدام به نمایش و خوشه‌بندی داده‌های نرمال شده نمایید. در نتیجه در این بخش باید ۶ نمودار (خروجی PCA و خوشه‌بندی سلسله‌مراتبی برای هر کدام از روش‌های نرمال سازی عنوان شده) ارائه گردد.

ب) در این بخش نیاز است تا به مقایسه روش‌های نرمال سازی بپردازید و نقاط ضعف و قوت هر کدام را به صورت جداگانه بیان فرمایید. لازم به ذکر است که ارجاع‌دهی به مراجع مورد استفاده ضروری است.

نکته ۱: شما می‌توانید جهت پیاده‌سازی از زبان‌های R و Python استفاده نمایید.

نکته ۲: جهت دریافت اطلاعات بیشتر در ارتباط با مجموعه داده ارائه شده می‌توانید به [Galaxy Training](#) مراجعه فرمایید.

## مسئله ۲. یادگیری ساختار گرافی ( ۵۰ نمره)

در این بخش می‌خواهیم به بررسی روش Glasso بپردازیم. ما یک مدل گرافی گوسی ایجاد کرده‌ایم و از آن نمونه گرفته‌ایم. ماتریس precision مدل مربوطه در فایل PrecisionMatrix.csv و نمونه‌های گرفته شده در فایل Samples.csv قرار دارد.

الف) ابتدا با روش نمونه‌برداری بدون جایگزین از نمونه‌های موجود در فایل Samples.csv، دسته‌های ۱۰ تایی، ۱۰۰ تایی و ۱۰۰۰ تایی از نمونه‌ها را ایجاد کنید. سپس به کمک روش Glasso و بر اساس مقادیر مختلف پارامتر تنظیم کننده، اقدام به شناسایی ساختار گراف نمایید. یال‌های گراف به دست آمده در حالات مختلف (با تعداد نمونه‌های مختلف) را با یال‌های گراف اصلی مقایسه نمایید. در ادامه منحنی‌های مربوط به نتایج به دست آمده برای حالات مختلف را در نموداری که محور عمودی آن  $TPR^1$

<sup>1</sup> True Positive Rate

و محور افقی آن  $FPR^2$  می‌باشد، رسم نمایید. نهایتاً سطح زیر نمودار را برای هر کدام از حالت‌ها به صورت جداگانه محاسبه و گزارش فرمایید.

خروجی: نموداری که محور عمودی آن  $TPR$  و محور افقی آن  $FPR$  است و در آن سه منحنی مربوط به حالات استفاده از ۱۰، ۱۰۰ و ۱۰۰۰ نمونه، وجود دارد + اندازه ناحیه زیر منحنی برای هر سه منحنی

ب) تابع  $f(x) = e^x$  را بر روی داده‌های مورد بررسی در بخش قبل اعمال کرده و بررسی‌های عنوان شده در بخش "الف" را مجدداً بر روی داده‌های جدید انجام دهید.

خروجی: نموداری که محور عمودی آن  $TPR$  و محور افقی آن  $FPR$  است و در آن سه منحنی مربوط به حالات استفاده از ۱۰، ۱۰۰ و ۱۰۰۰ نمونه، وجود دارد + اندازه ناحیه زیر منحنی برای هر سه منحنی

ج) با استفاده از داده‌های مربوط به گره‌های ۱ تا ۱۰ موجود در فایل Samples.csv و به کمک روش‌های انتخاب مدل  $AIC^3$  و  $BIC^4$  به منظور شناسایی مقدار مناسب پارامتر تنظیم‌کننده مدل Glasso، ساختار گرافی را در این حالات شناسایی کرده و نمایش دهید.

خروجی: دو ساختار گرافی، یکی برای حالت استفاده از روش  $AIC$  برای تنظیم پارامتر و دیگری برای حالت استفاده از روش  $BIC$

---

<sup>2</sup> False Positive Rate

<sup>3</sup> Akaike Information Criterion

<sup>4</sup> Bayesian Information Criterion