

转录组分析报告

1 数据简介

小鼠胎盘组织的 RNA-seq 数据，敲了x基因，含2个时间点 (E9.5， E11.5)。

E9.5 (E10) ： 6个 (3个KO, 3个WT)

E11.5 (E12) ： 4个 (2个KO, 2个WT)

前期数据已经经过SOAOnuke过滤软件对原始数据进行质控，接下来的分析则使用处理过后的cleandata。其中原始数据所测的为双端150bp，因此使用STAR软件构建相应的参考基因组。

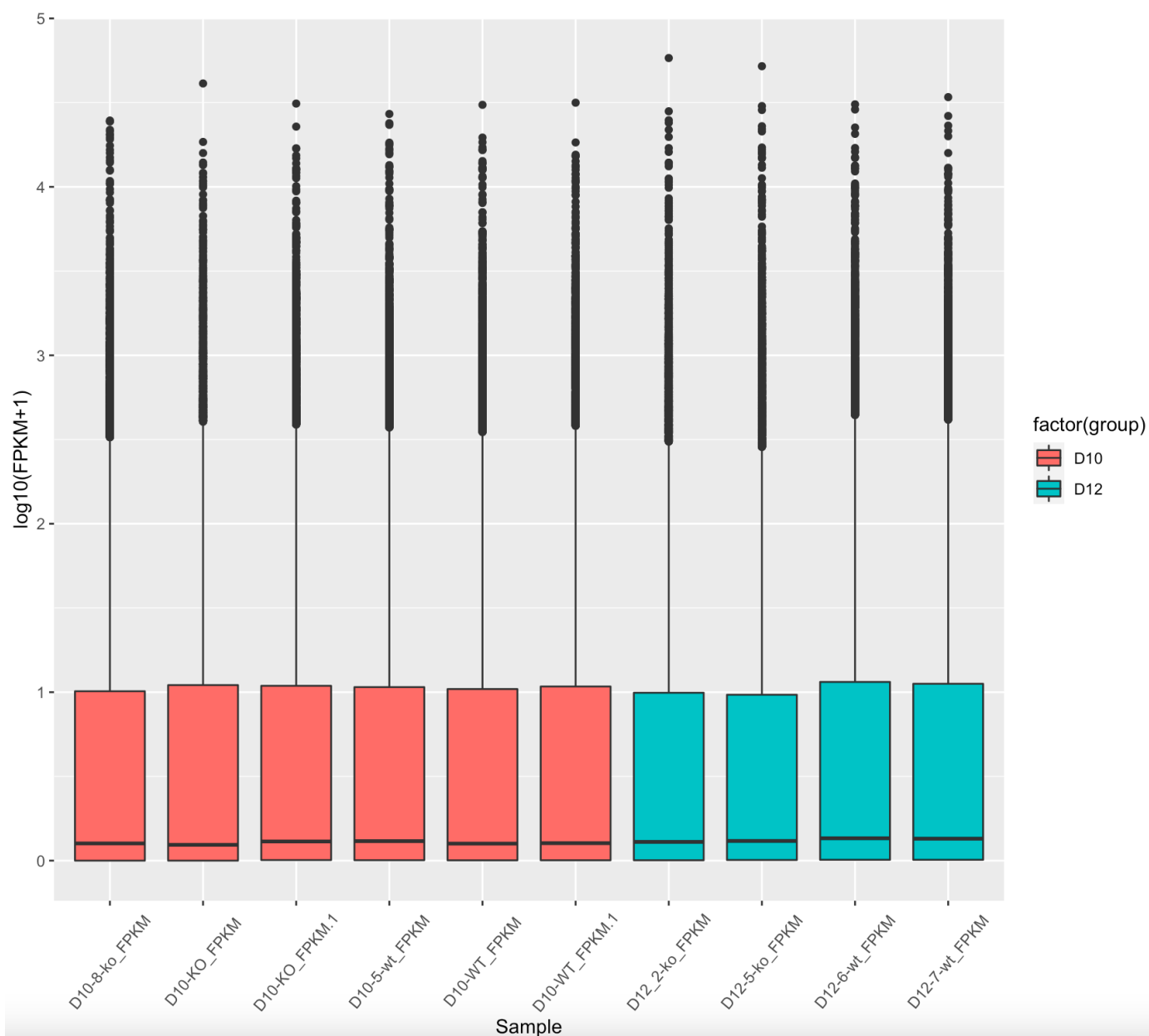
Sample
0001_D10-5-wt
013-02-D10-WT
013-10-D10-WT
0001_D10-8-ko
013-4-D10-KO
013-7-D10-KO
0003_D12-5-ko
0003_D12-6-wt
0003_D12-7-wt
0005_D12_2-ko

2 定量分析

比对软件使用的为STAR，基因定量为subread中featurecounts函数。

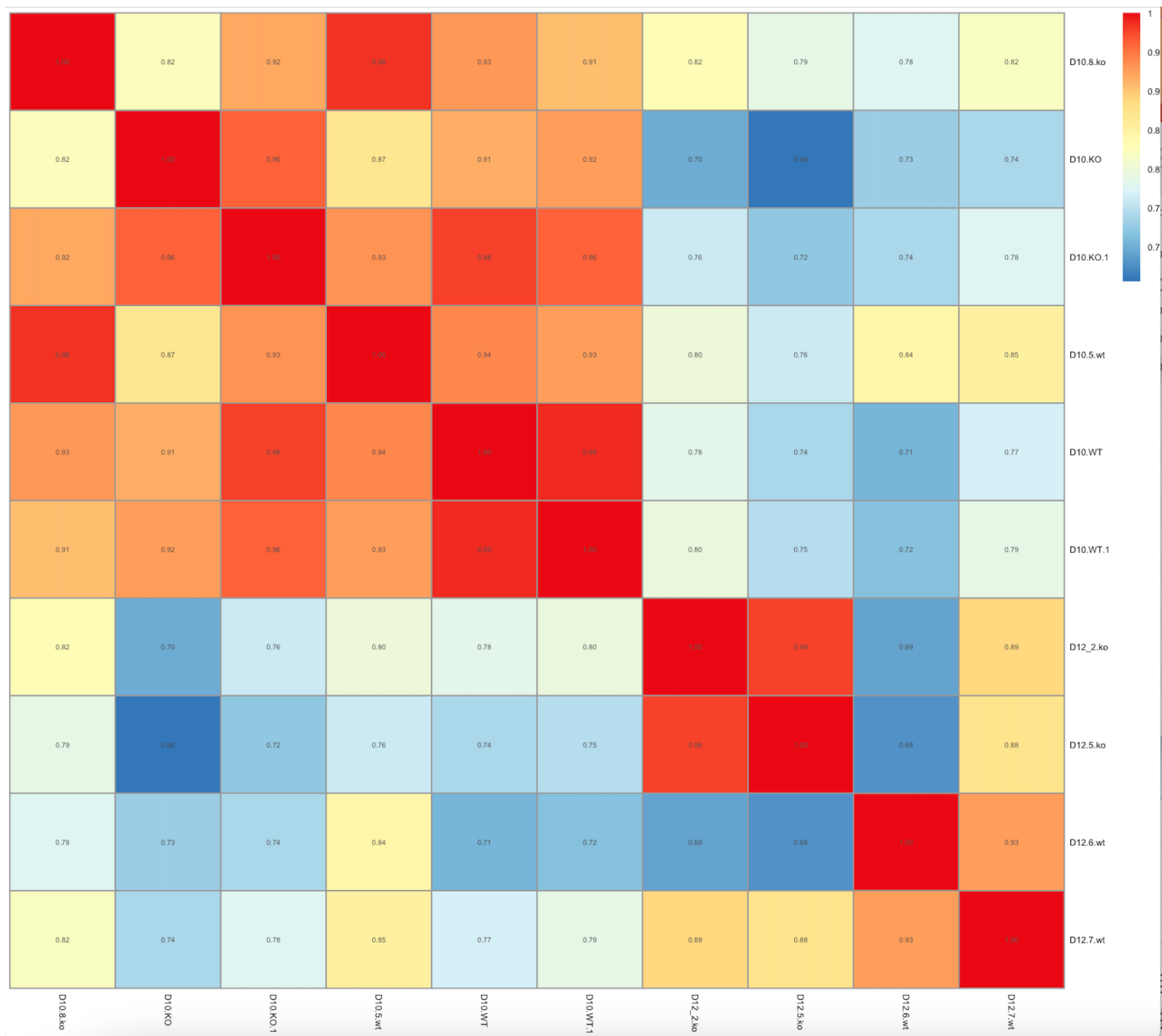
all_feature.txt代表所有样本比对得到的counts文件。

- Genename：基因名
- sample：各样本定量所得的原始read count值
- chr：基因所在的染色体名称
- start：基因所在染色体的起始位置
- end：基因所在染色体的终止位置
- strand：基因所在染色体的正负链信息
- length：基因长度，基因起始到终止所有exon非重叠区域的总和



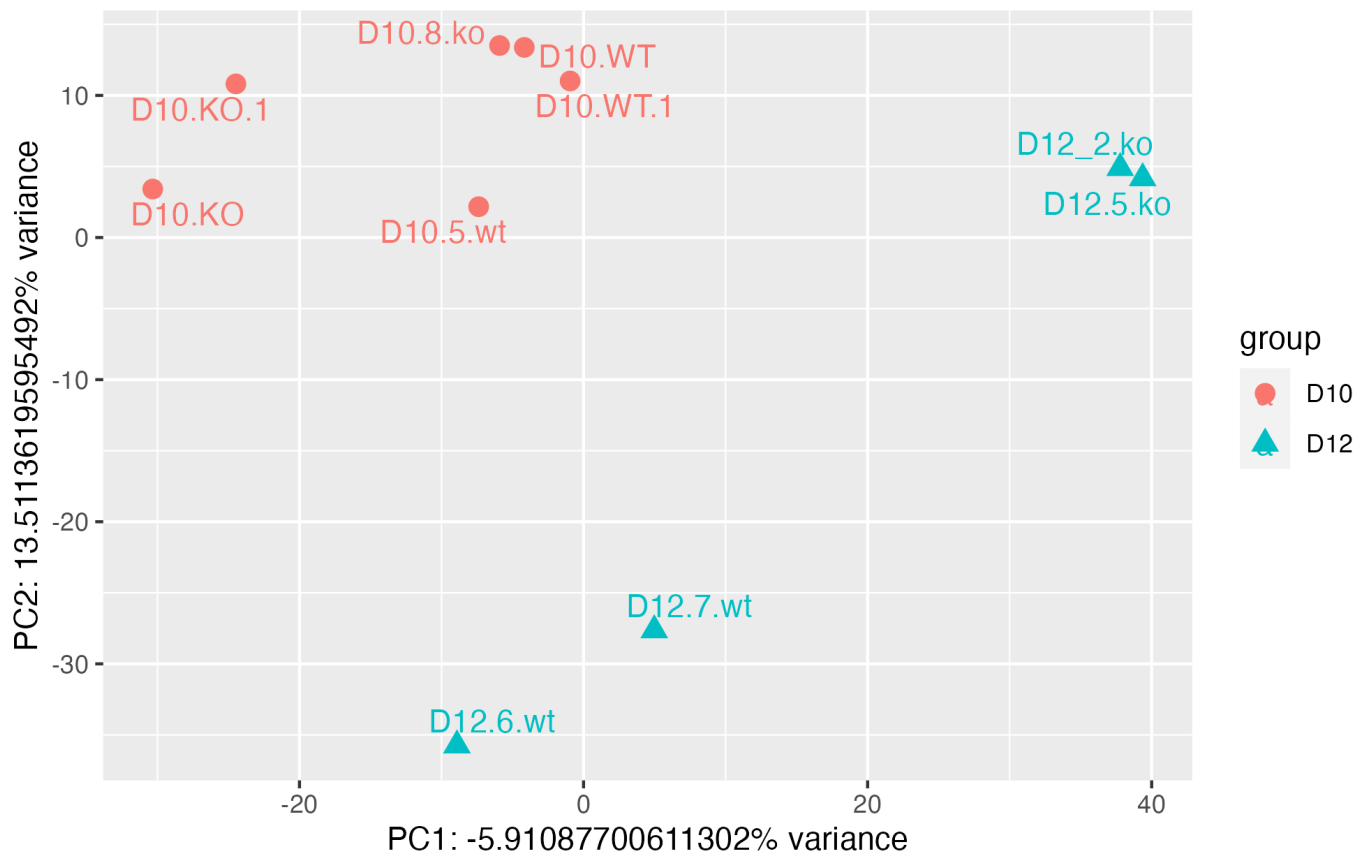
2.2 样本间相关性

生物学重复通常是任何生物学实验所必须的，目前主流期刊也基本要求生物学重复。生物学重复主要有两个用途：一个是证明所涉及的生物学实验操作不是偶然，而是可重复的。另一个是为了确保后续的差异基因分析得到更可靠的结果。样品间基因表达水平相关性是检验实验可靠性和样本选择是否合理的重要指标。相关系数越接近1，表明样品之间表达模式的相似度越高。Encode计划建议皮尔逊相关系数的平方(R²)大于0.92(理想的取样和实验条件下)。具体的项目操作中，我们要求生物学重复样品间R²至少要大于0.8，否则需要对样品做出合适的解释，或者重新进行实验。根据各样本所有基因的FPKM值计算组内及组间样本的相关性系数，绘制成热图，可直观显示组间样本差异及组内样本重复情况。样本间相关性系数越高，其表达模式越为接近，样本相关性热图如下图所示，见结果文件：gene_cor.pdf。



2.3 主成分分析

主成分分析（PCA）也常用来评估组间差异及组内样本重复情况，PCA采用线性代数的计算方法，对数以万计的基因变量进行降维及主成分提取。我们对所有样本的基因表达值（FPKM）进行PCA分析，如下图所示。理想条件下，PCA图中，组间样本应该分散，组内样本应该聚在一起，见结果文件：PCA.png。



3 差异分析

基因表达定量完成后，需要对其表达数据进行统计学分析，筛选样本在不同状态下表达水平显著差异的基因。差异分析主要分为三个步骤。

- 首先对原始的readcount进行标准化（normalization），主要是对测序深度的校正。
- 然后统计学模型进行假设检验概率（pvalue）的计算
- 最后进行多重假设检验校正，得到FDR值（错误发现率，padj是其常见形式）。

针对不同的实验情况，我们选用合适的软件进行基因表达差异显著性分析，具体如下表所示。

表3.5 表达差异分析所用软件及差异基因筛选标准

类型	软件	标准化方法	pvalue计算模型	FDR计算方法	差异基因筛选标准
有生物学重复	DESeq2(Anders et al, 2014)	DESeq	负二项分布	BH	$ \log_2(\text{FoldChange}) \geq 1$ & $\text{padj} \leq 0.05$
无生物学重复	edgeR(Robinson et al, 2010)	TMM	负二项分布	BH	$ \log_2(\text{FoldChange}) \geq 1$ & $\text{padj} \leq 0.05$

为了找到更多的基因，将阈值调整为 $|\log_2(\text{FoldChange})| \geq 0$

实验设计：

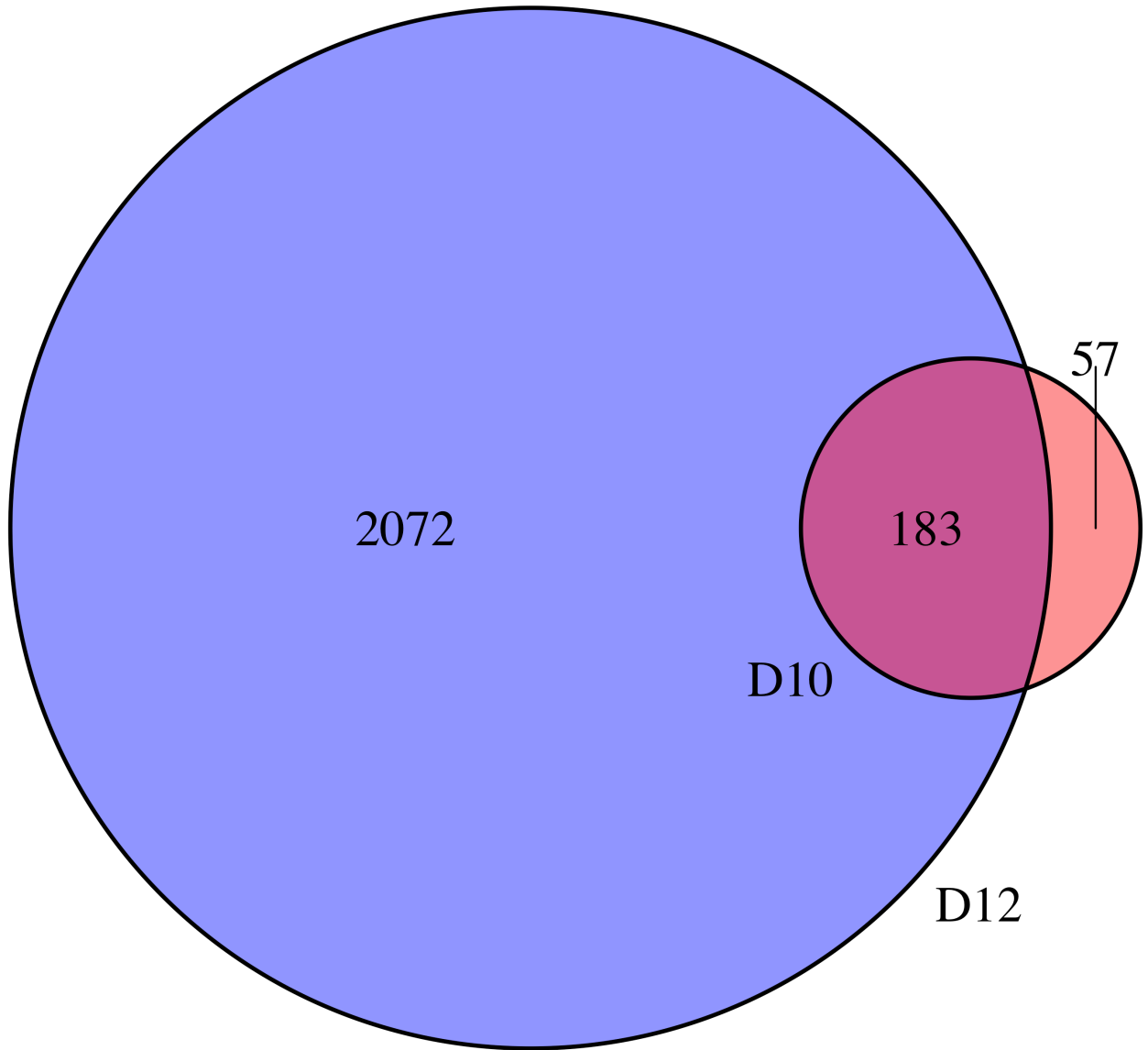
- 1. D10：处理组KO与对照组WT做差异分析
 - 2. D12：处理组KO与对照组WT做差异分析
- **gene_id**：基因编号
 - **basemean**：表示所有样本经过归一化系数矫正的read counts (counts/sizeFactor) 的均值。
 - **log2FoldChange**：处理组与对照组基因表达水平的比值，再经过差异分析软件收缩模型处理，最后以2为底取对数
 - **pvalue**：显著性检验的值
 - **padj**:校正后的p值

文件结果：

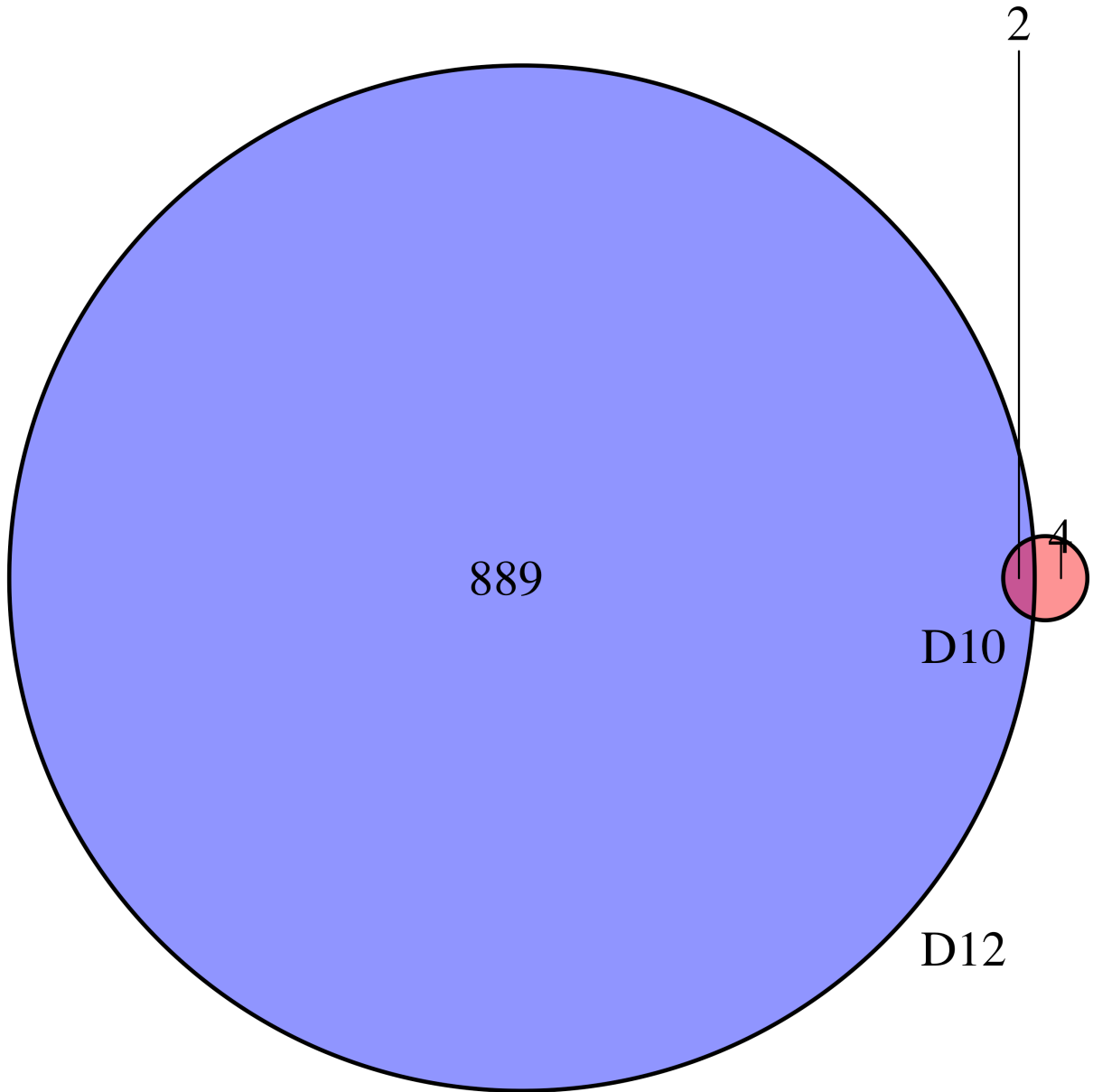
- 1. D10组：
 - D10_DEGs.csv 为padj<0.05的差异基因
 - upDEGsD10.csv为padj<0.05的上调差异基因
 - downDEGsD10.csv为padj<0.05的下调差异基因
- 2. D12组：
 - D12_DEGs.csv 为padj<0.05的差异基因
 - upDEGsD12.csv为padj<0.05的上调差异基因
 - downDEGsD12.csv为padj<0.05的下调差异基因

Venn图

all_D10_D12_KOvs.WT_0.05_venn



up_D10_D12_KOvs.WT_0.05_venn



down_D10_D12_KOvs.WT_0.05_venn

