

THE DISCORDANT METHOD: A NOVEL APPROACH FOR DIFFERENTIAL CORRELATION

Charlotte Siska, Russell Bowler, Katerina Kechris; 2015

By Soel Micheletti · JQ Lab Journal Club · June 2023

Condition 1



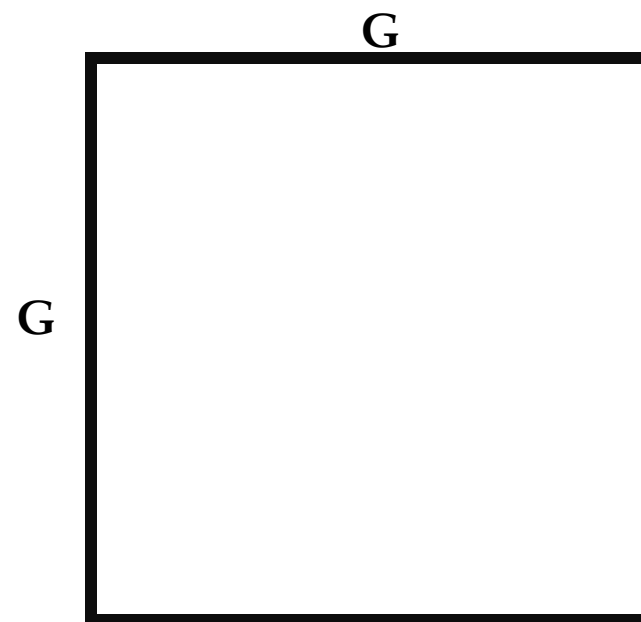
Condition 2



Condition 1



Condition 2



Condition 1

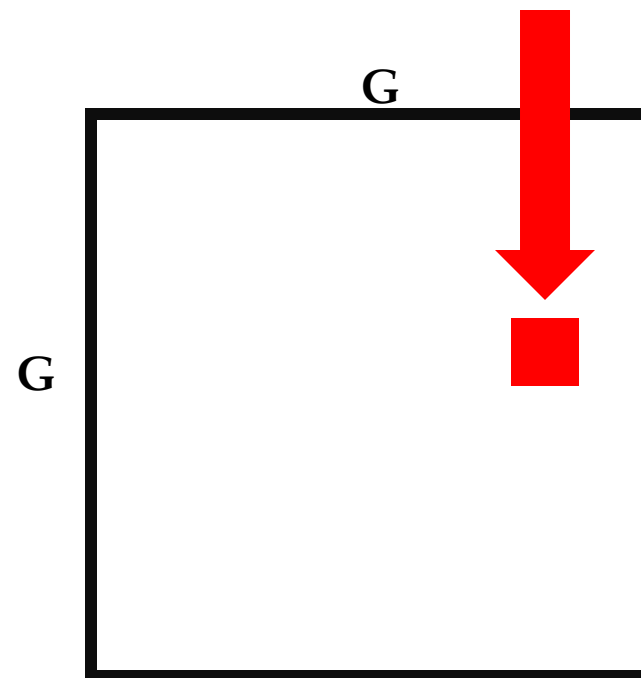


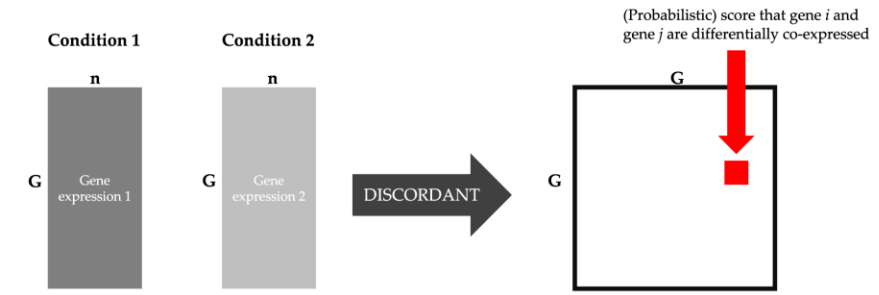
Condition 2



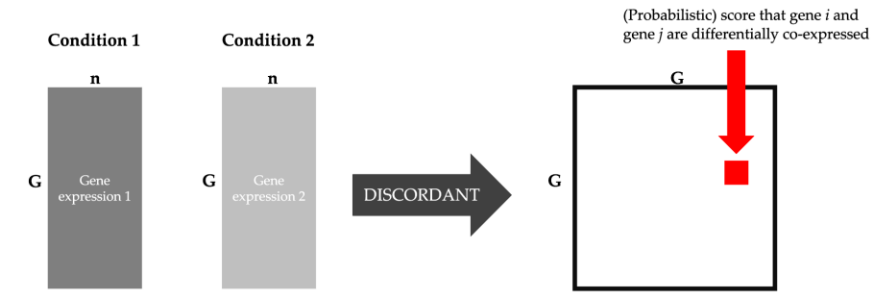
DISCORDANT

(Probabilistic) score that gene i and gene j are differentially co-expressed

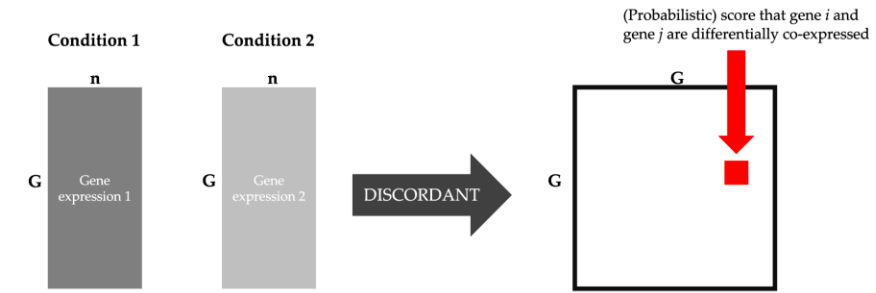




- Compute co-expression matrix for both conditions



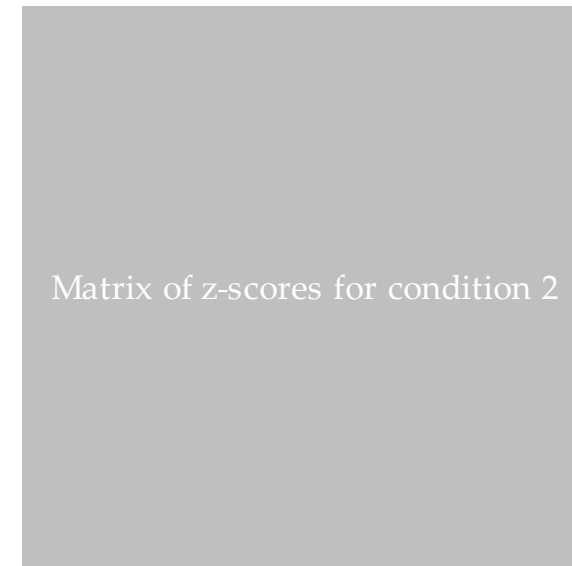
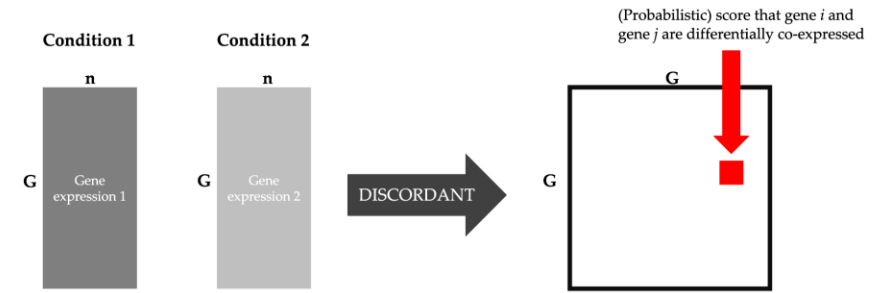
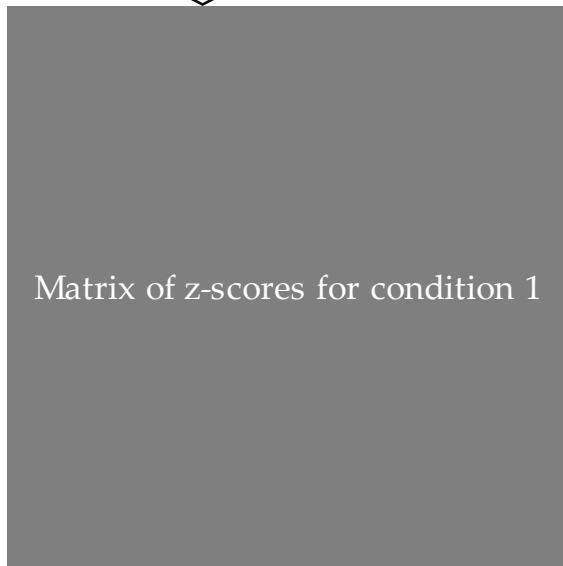
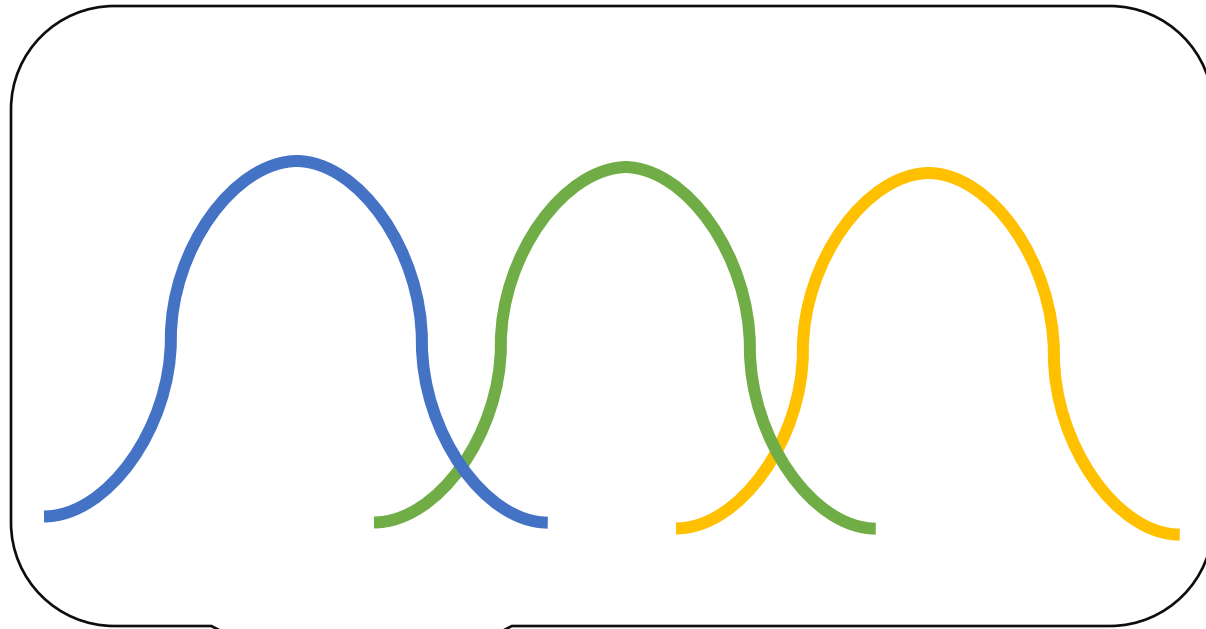
- Compute co-expression matrix for both conditions
- Apply Fisher's transformation to convert Pearson's correlation coefficients into z-scores

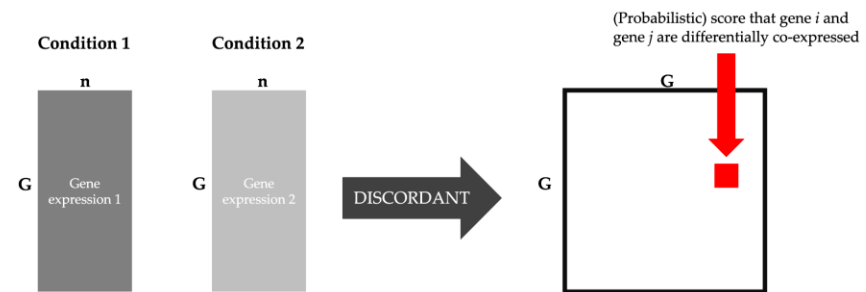
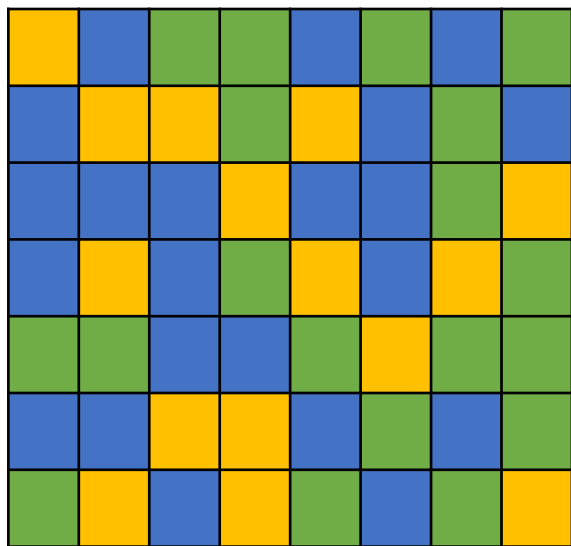
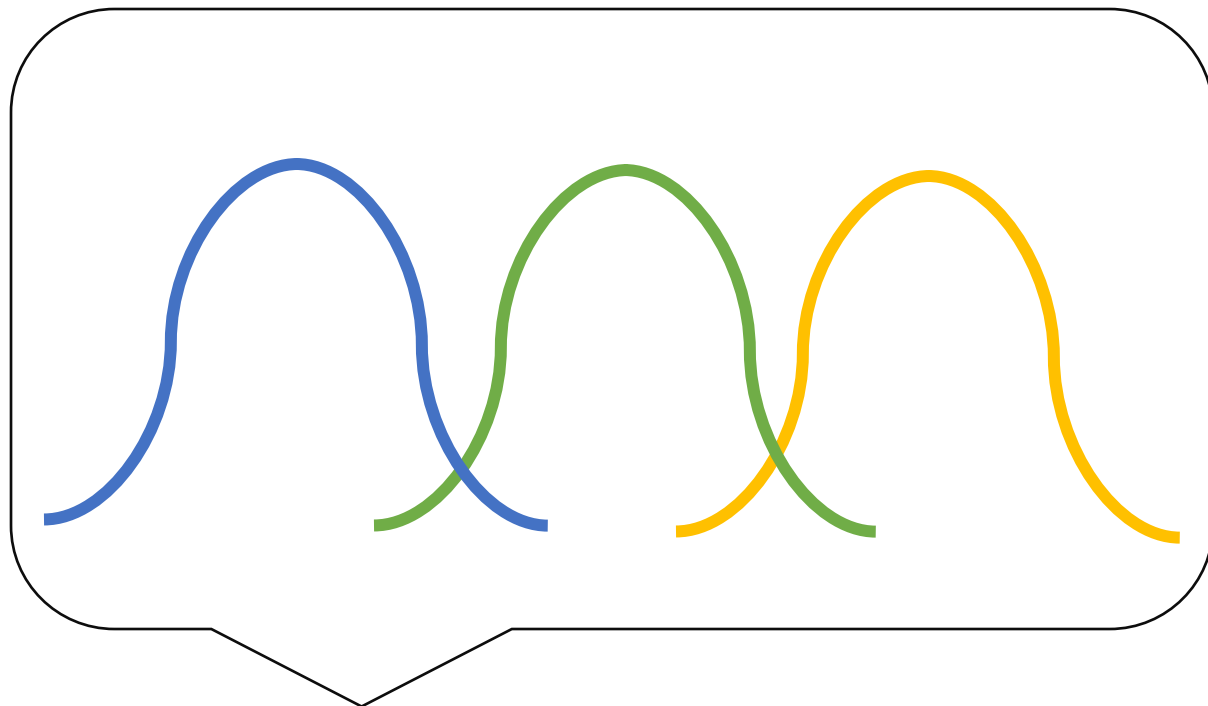


- Compute co-expression matrix for both conditions
- Apply Fisher's transformation to convert Pearson's correlation coefficients into z-scores

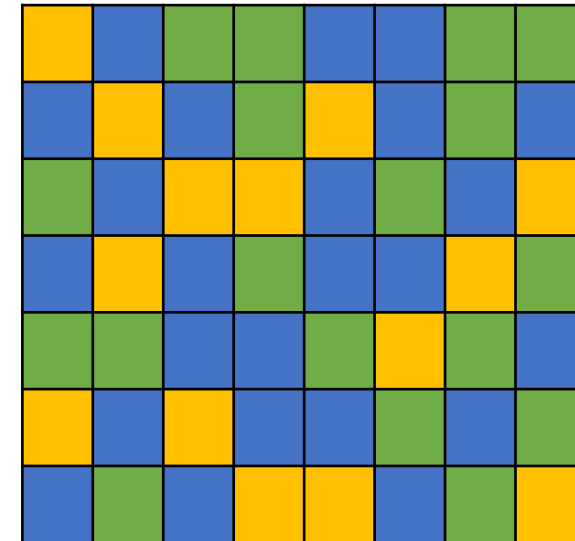
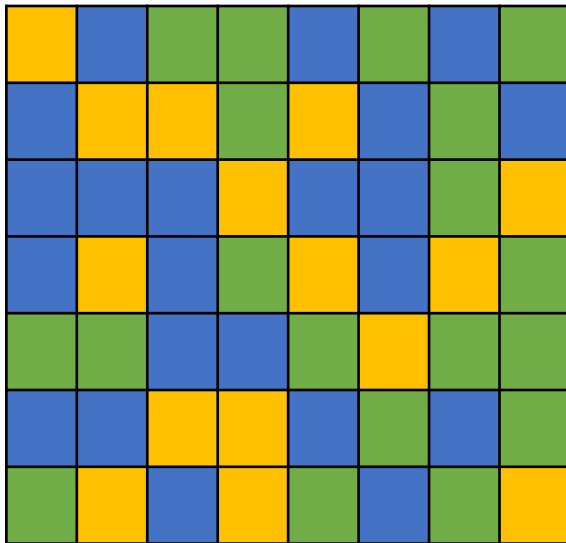
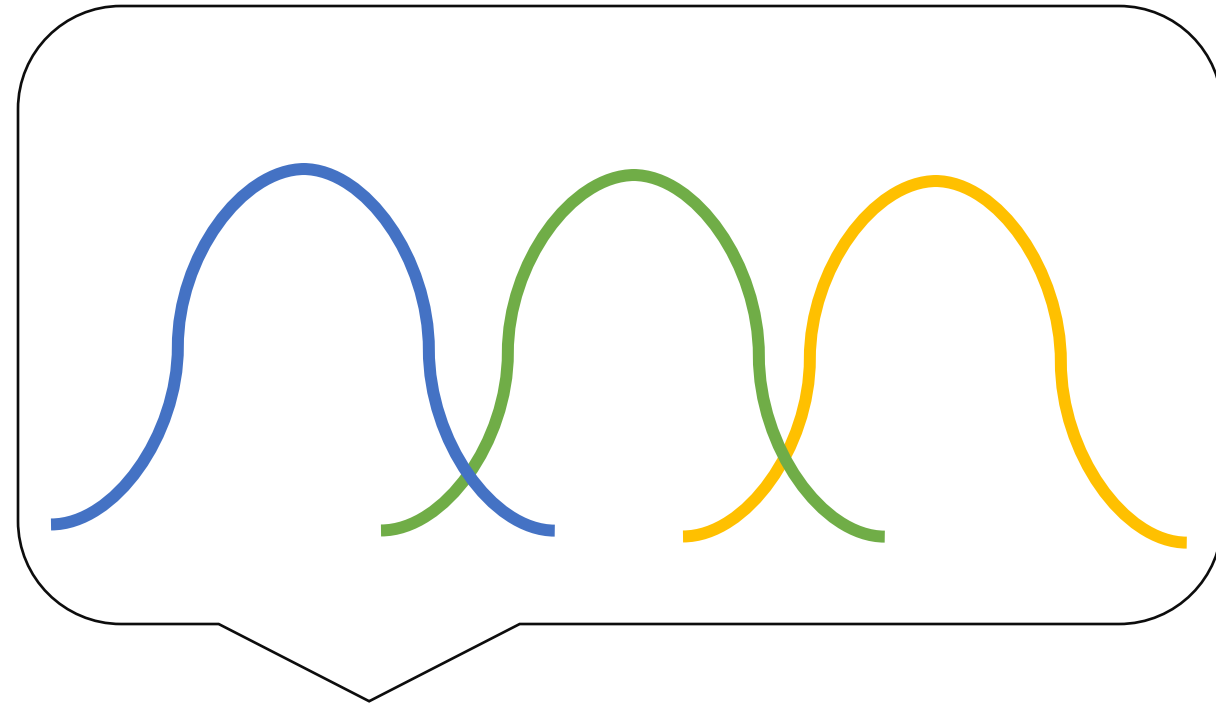
Matrix of z-scores for condition 1

Matrix of z-scores for condition 2



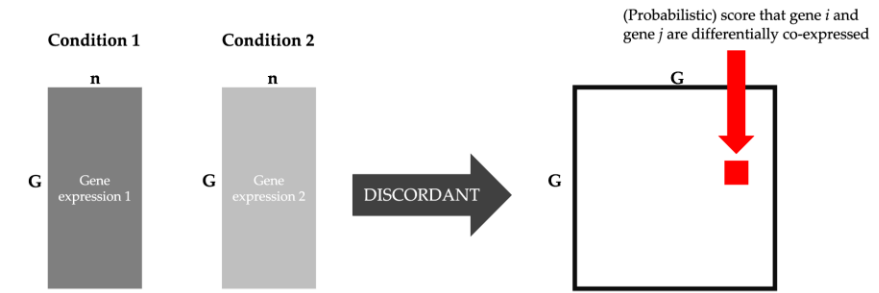
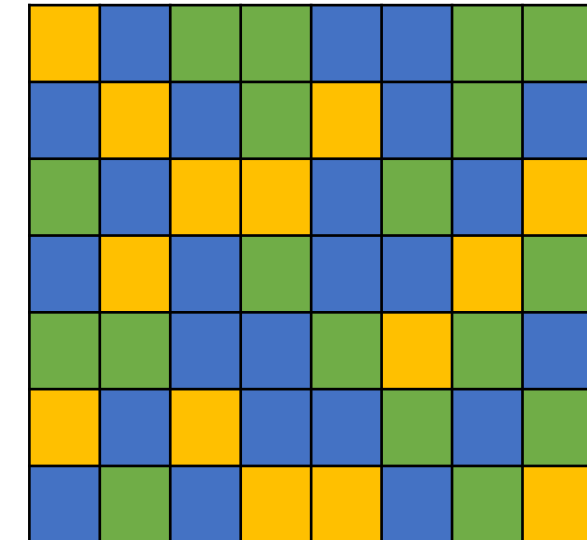
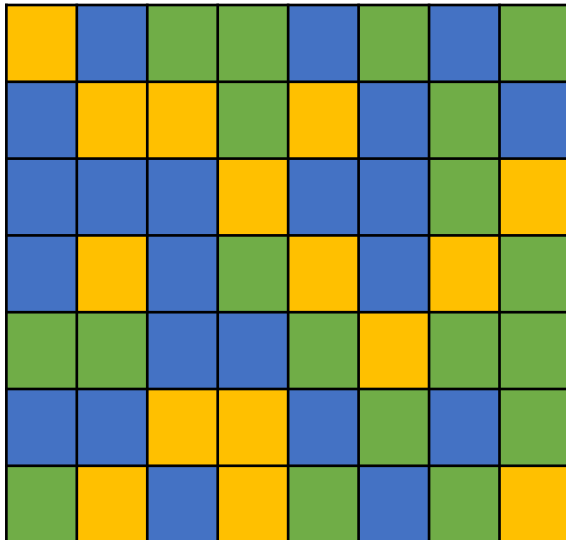


Matrix of z-scores for condition 2



Main idea:

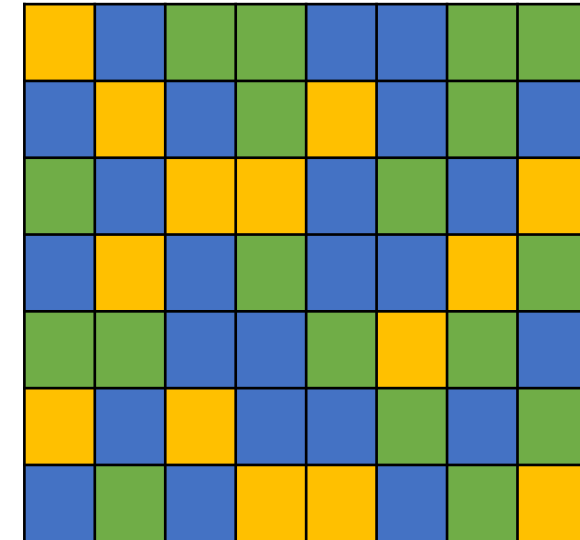
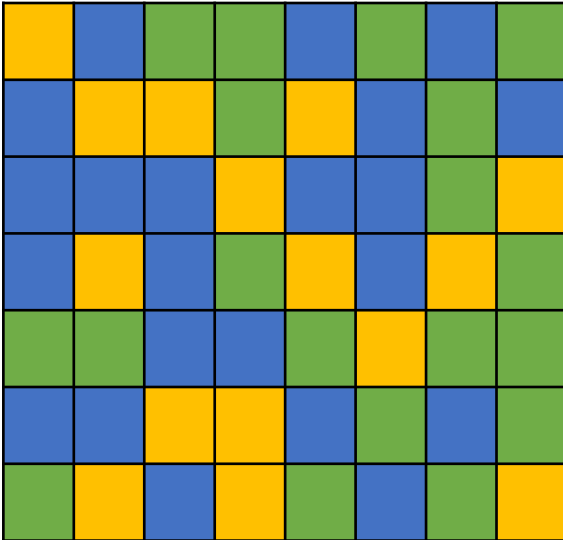
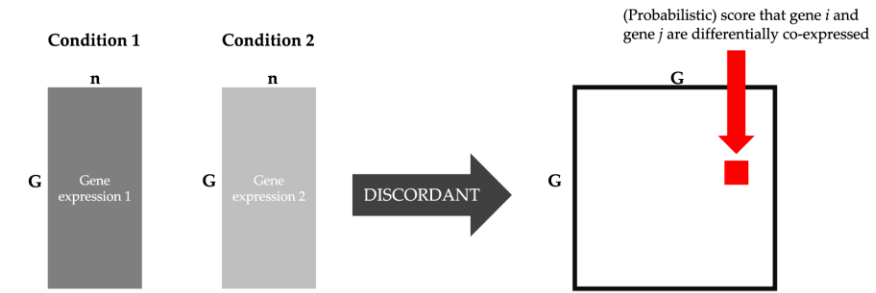
- Correlation smaller than average
- Average correlation
- Correlation larger than average



Main idea:

- Correlation smaller than average
- Average correlation
- Correlation larger than average

A pair is differentially coexpressed if it belongs to different categories

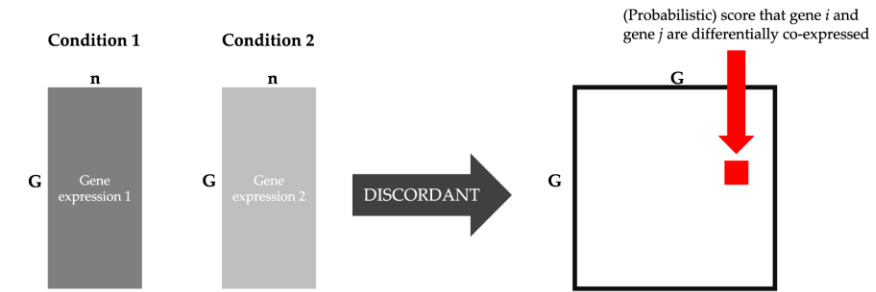
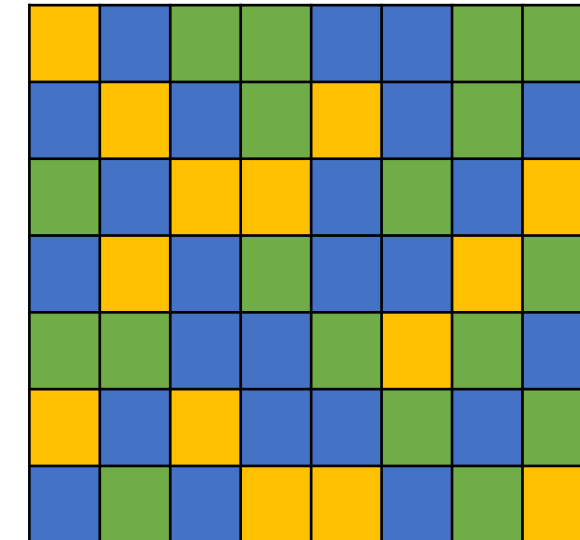
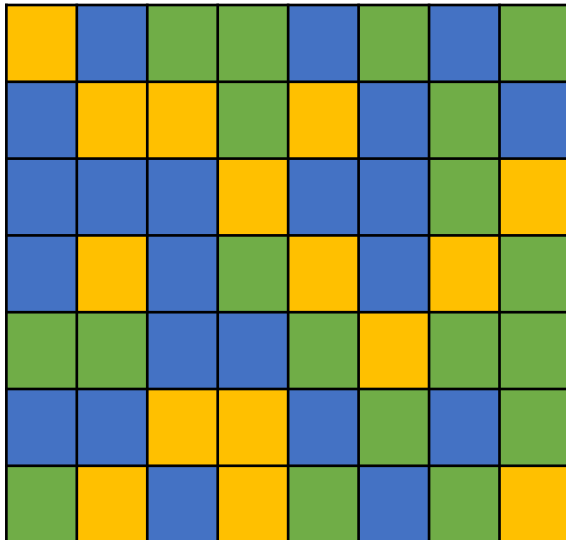


Main idea:

- Correlation smaller than average
- Average correlation
- Correlation larger than average

A pair is differentially coexpressed if it belongs to different categories

For each gene pair, compute $\Pr[\text{blue}]$, $\Pr[\text{green}]$, and $\Pr[\text{orange}]$



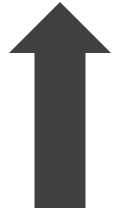
$$\arg \max_{\Theta=[\Theta_1,\Theta_2]} \prod_k Pr \left[z_1(k), z_2(k) | \Theta \right]$$

$$\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta]$$



Gene pairs

$$\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta]$$



Mean, covariance, and weight for
each normal distribution

$$\arg \max_{\Theta=[\Theta_1,\Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] = \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2]$$

$$\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] = \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence})$$

$$\begin{aligned}
 \arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplyfing assumption: independence}) \\
 &= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)})
 \end{aligned}$$

$$\begin{aligned}
\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence}) \\
&= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)})
\end{aligned}$$

Difficult optimization: non-convex with constraints

$$\begin{aligned}
\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence}) \\
&= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)})
\end{aligned}$$

Difficult optimization: non-convex with constraints ==> **EM algorithm**

$$\begin{aligned}
\arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence}) \\
&= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)})
\end{aligned}$$

- Initialize parameters $\Theta_i = \{w_0^{(i)}, w_1^{(i)}, w_2^{(i)}, \mu_0^{(i)}, \mu_1^{(i)}, \mu_2^{(i)}, \Sigma_0^{(i)}, \Sigma_1^{(i)}, \Sigma_2^{(i)}\}$ for both conditions i

$$\begin{aligned} \arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence}) \\ &= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)}) \end{aligned}$$

- Initialize parameters $\Theta_i = \{w_0^{(i)}, w_1^{(i)}, w_2^{(i)}, \mu_0^{(i)}, \mu_1^{(i)}, \mu_2^{(i)}, \Sigma_0^{(i)}, \Sigma_1^{(i)}, \Sigma_2^{(i)}\}$ for both conditions i
- Iterate until convergence

- **E-step**

Assign to each sample a probability of being in 'color i' in the first condition and 'color j' in the second condition

- **M-step**

Solve the MLE using the soft labels

$$\begin{aligned} \arg \max_{\Theta=[\Theta_1, \Theta_2]} \prod_k Pr [z_1(k), z_2(k) | \Theta] &= \prod_k Pr [z_1(k) | \Theta_1] Pr [z_2(k) | \Theta_2] \quad (\text{Simplifying assumption: independence}) \\ &= \prod_k \sum_{i=0}^2 \sum_{j=0}^2 w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)}) w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)}) \end{aligned}$$

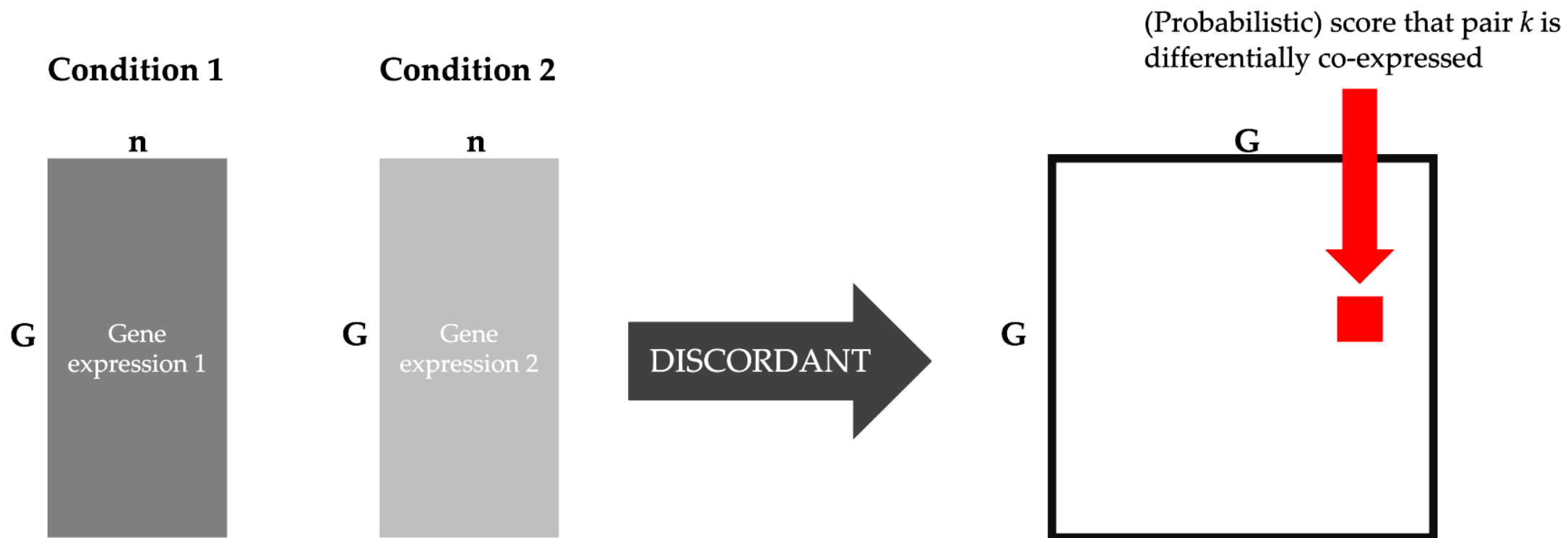
- Initialize parameters $\Theta_i = \{w_0^{(i)}, w_1^{(i)}, w_2^{(i)}, \mu_0^{(i)}, \mu_1^{(i)}, \mu_2^{(i)}, \Sigma_0^{(i)}, \Sigma_1^{(i)}, \Sigma_2^{(i)}\}$ for both conditions i
- Iterate until convergence

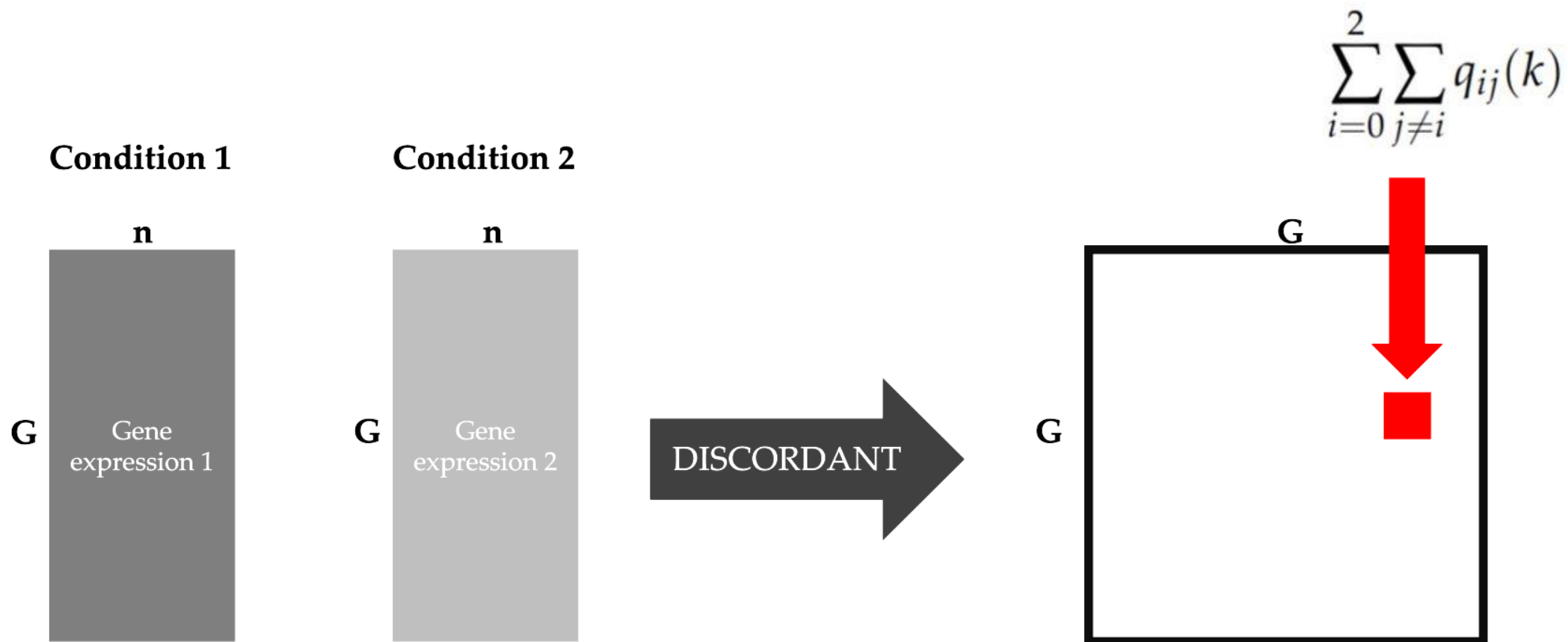
- **E-step**

$$q_{ij}(k) = \frac{w_i^{(1)} \mathcal{N}(z_1(k) | \mu_i^{(1)}, \Sigma_i^{(1)})}{\sum_{\tau=0}^2 w_{\tau}^{(1)} \mathcal{N}(z_1(k) | \mu_{\tau}^{(1)}, \Sigma_{\tau}^{(1)})} \cdot \frac{w_j^{(2)} \mathcal{N}(z_2(k) | \mu_j^{(2)}, \Sigma_j^{(2)})}{\sum_{\tau=0}^2 w_{\tau}^{(2)} \mathcal{N}(z_2(k) | \mu_{\tau}^{(2)}, \Sigma_{\tau}^{(2)})}$$

- **M-step**

Solve the MLE using the soft labels





Biological validation

- Dataset containing miRNA + mRNA expression for 10 healthy subjects and 21 patients with glioblastoma multiforme
- Interested in the differential correlation between miRNA and mRNA pairs

Biological validation

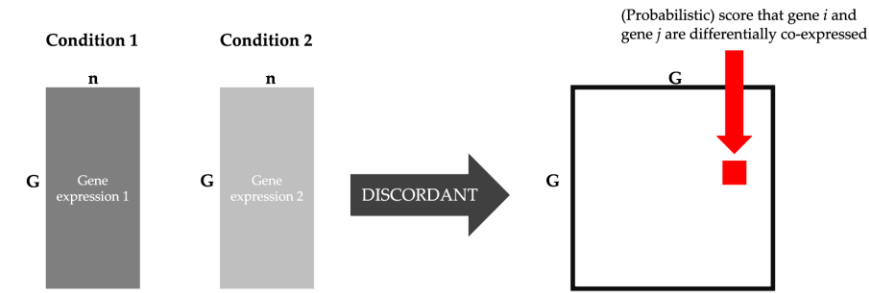
- Dataset containing miRNA + mRNA expression for 10 healthy subjects and 21 patients with glioblastoma multiforme
- Interested in the differential correlation between miRNA and mRNA pairs

Data	Method	Mean	Median
GBM	Discordant	464.75	347.5
	EBcoexpress	815	607
	Fisher	781	801
	Linear (miRNA-independent)	1095	532.5
	Linear (transcript-independent)	2596.5	787.5

My take-aways

- The probability of differential co-expression can be conveniently used in GSEA
- Methodological paper, no thorough validation with biological data
- Idea of a data-driven binning is elegant:
 - It improves interpretability
 - It enables to determine more differentially correlated pairs than other methods, potentially improving power of detecting disrupted interactions
- Assumption that the z-scores form a tri-modal normal distribution is not discussed

My take-aways



- The probability of differential co-expression can be conveniently used in GSEA
- Methodological paper, no thorough validation with biological data
- Idea of a data-driven binning is elegant:
 - It improves interpretability
 - It enables to determine more differentially correlated pairs than other methods, potentially improving power of detecting disrupted interactions
- Assumption that the z-scores form a tri-modal normal distribution is not discussed