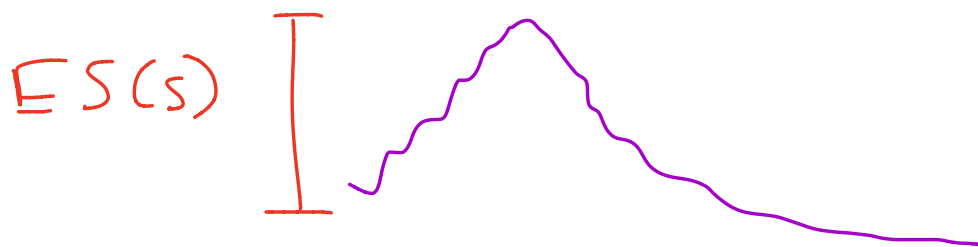# Randomization and Permutation

- Compensate for lack of power

- $S$ gene set
- $Z_s = \{Z_i : i \in S\}$  $\quad m = |s|$

$$S(z_s) = \sum_{i=1}^{s} \frac{S(z_i)}{m}$$

then we use KS testing on the estimated CDFs.

$$ES(s) \quad$$



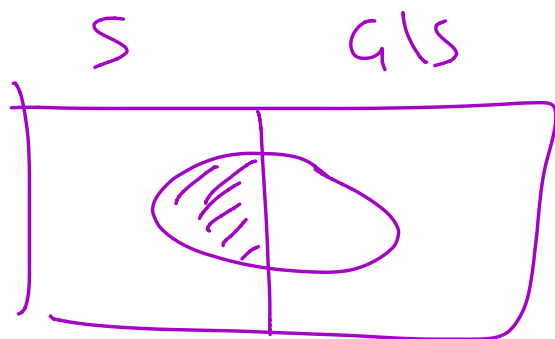How to asses significance of the enrichment statistic?

$$|S| = m \qquad J = n$$

$$|G| = N$$

$$A = \{i : |z_i| > c\} \quad |A| = a$$

$$|S \cap A| = h$$



$$h \sim \text{Hypergeometri}(m, N-m, a)$$

$$P(X = h) = \frac{\binom{m}{k}\binom{N-m}{a-h}}{\binom{N}{a}}$$

So our test statistic is given by

$$\sum_{k \geq \#obs.} \mathbb{P}(x=k)$$

In general we don't just use counts we use some score.

| Group 1 | Group 2 |
|---|---|
| gene 1 $\longrightarrow$ | |
| $\vdots$ | |
| gene N $\longrightarrow$ | |

→ Column permutation to calculate

$$z_*^1, \ldots, z_*^B$$

$$S_b^* = S(z_s^*)$$

$$p\text{-val} = \#\{S^{*b} > S\}$$

$$\frac{\qquad}{B}$$

- Row randomization

  - Select subset of size $m$

  $$Z_7^b = \{z_i : i \in I_7^b\}$$

  $$S_7^b = (S_{Z_7^b})$$

  $$p\text{-val} = \frac{\#\ \{S_7^b > S\}}{B}$$

$$S(z_s) = \sum_{z_{ic}} \frac{S(z_i)}{m}$$

$$m_s = \frac{1}{N} \sum_{i=1}^{n} s(z_i)$$

$$Sd_s = \sqrt{\frac{1}{N} \sum_{i=1}^{N} (S(z_i) - m_s)^2}$$

then    from the permutations

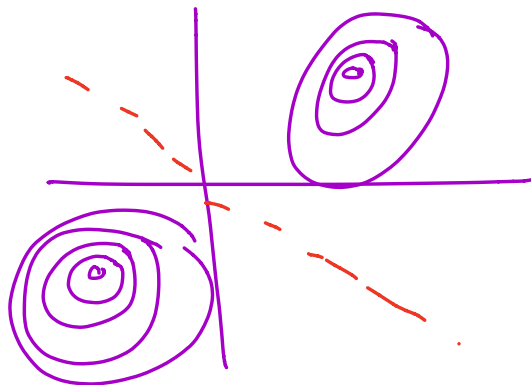$$S_{**}^{b} = m_S + S d_S \left\{ \frac{S_*^{b} - m_S^{*}}{s d_*} \right\}$$

$$P_S^{**} = \frac{\# \left\{ S_{**}^{b} > S \right\}}{B}$$

## Chapter 11    Prediction & Others

Fishers linear discrimnt function

$$y = 1 \qquad x \sim N_N (\delta_{,,} \Sigma_{,})$$

$$y = 2 \qquad x \sim N_N (\delta_2, \Sigma)$$

$$P(Y=1|X) \propto P(X|Y=1)P(Y=1)$$

$$N_N(\delta_1, \Sigma)\pi_1$$

$$\log \frac{P(Y=2|x)}{P(Y=1|x)} = \log \frac{P(Y=2)}{P(Y=1)} + \log \frac{P(X|Y=1)}{P(X|Y=2)}$$

$$\log \frac{\exp\left(-\frac{1}{2}(x-\delta_2)^T \Sigma^{-1}(x-\delta_2)\right)}{\exp\left(-\frac{1}{2}(x-\delta_1)^T \Sigma^{-1}(x-\delta_2)\right)}$$

$$\log \exp\left\{ (\delta_2-\delta_1)\Sigma^{-1}x - \frac{1}{2}\delta_2^T \Sigma \delta_2 + \frac{1}{2}\delta_1^T \Sigma \delta_1\right\}$$

So the log ratio is

$$\log \frac{P(y=2|x)}{P(y=1|x)} = \beta_0 + \beta'x$$

$$\beta^T = (\delta_2 - \delta_1)^T \Sigma^{-1}$$

"Best" for Bayesian stuff

- Bad when $\Sigma_1 \neq \Sigma_2$.

- Bad for high dimensional LDA.

## A frequentist model

$$U_i = \frac{X_i - \mu_i}{\sigma_i}$$

$$\vec{u} \sim N\left(\frac{-\delta}{2c_0}, I\right) \text{ for } y=1$$

$$\vec{u} \sim N\left(\frac{\delta}{2c_0}, I\right) \text{ for } y=2$$

$$C_0 = \sqrt{\frac{n_1 n_2}{n}}$$

Construct a prediction function

$$S = \sum \delta_i \cdot u_i \sim N\left(\pm \|\delta\|^2 / 2c_0, \|\delta\|^2\right)$$

Classifier: $S < 0 \quad \hat{y} = 1$

$\qquad\qquad S > 0 \quad \hat{y} = 2$

$$\mathbb{P}_1(S > 0) = \mathbb{P}_1\left( \frac{S + \|\delta\|^2/2c_0}{\|\delta\|} > \frac{\|\delta\|}{2c_0} \right)$$

$$= \Phi\left( -\frac{\|\delta\|^2}{2c_0} \right)$$

Then to estimate

$$\frac{\overline{X_{i_1}} + \overline{X_{i_2}}}{2} \qquad \hat{\sigma_i}^2 = \frac{SS_{i_1} + SS_{i_2}}{n-2}$$

$$\overline{\delta_i} \sim c_0 \frac{\overline{X_{i_2}} - \overline{X_{i_1}}}{\sigma_i} \sim N(\delta_i, 1)$$

We could also use CV to
choose the $\delta$ vector.

Bayes & Empirical Bayes

$$\delta \sim g \qquad z \mid \delta \sim N(\delta, 1)$$

$$f = \int \phi(z-\delta)\, g\, \delta\, d\delta$$

$$g(\delta \mid z) = \exp\left(z\delta - \psi(z)\right)\exp\left(-\frac{\delta^2}{2} g(\delta)\right)$$

$$\psi(z) = \lg\left(\frac{f(z)}{\varphi(z)}\right)$$

in the exp. family.

easy to find mean/variance

$$\mathbb{E}(\delta \mid z) = z + \psi'(z)$$

$$V(\delta \mid z) = 1 + \psi''(z)$$