

JAMES K. PRINGLE
140.674 Stat Theory
Dr. Constantine Frangakis
Assignment 6
May 12, 2014

Assignment 6

The microarray experiment

1. Let Z be a Bernoulli random variable with probability $\text{pr}(Z = 1) = \pi$; let Y be such that, if $Z = 1$, then Y is a draw from a standard normal distribution; but if $Z = 0$, then Y is a draw from a normal distribution with mean μ and variance 1.

(i) Find $E(Y \mid \mu, \pi)$.

Proof. Let ϕ be the density of a standard normal distribution. Then

$$\text{pr}(Y = y \mid Z = z, \mu, \pi) = \phi(y)^z \phi(y - \mu)^{1-z}$$

Therefore,

$$\begin{aligned} \text{pr}(Y = y \mid \mu, \pi) &= \int \text{pr}(Y = y \mid Z = z, \mu, \pi) \text{pr}(Z = z \mid \mu, \pi) dz \\ &= \int \phi(y)^z \phi(y - \mu)^{1-z} \pi^z (1 - \pi)^{1-z} dz \\ &= \phi(y)\pi + \phi(y - \mu)(1 - \pi) \end{aligned}$$

Finally, calculating,

$$\begin{aligned} E(Y \mid \mu, \pi) &= \int y(\phi(y)\pi + \phi(y - \mu)(1 - \pi)) dy \\ &= \pi \int y\phi(y) dy + (1 - \pi) \int y\phi(y - \mu) dy \\ &= (1 - \pi)\mu \end{aligned}$$

□

(ii) Find $E(Y^2 \mid \mu, \pi)$.

Proof. Calculating,

$$\begin{aligned}
 E(Y^2 \mid \mu, \pi) &= \int y^2(\phi(y)\pi + \phi(y - \mu)(1 - \pi))dy \\
 &= \pi \int y^2\phi(y)dy + (1 - \pi) \int y^2\phi(y - \mu)dy \\
 &= \pi + (1 - \pi)(1 + \mu^2) \\
 &= 1 + (1 - \pi)\mu^2
 \end{aligned}$$

□

2. One type of experiment of microarray technology is summarized as follows:

- (a) A chip has a very large number of spots, say, $i = 1, \dots, n = 10000$ spots. In each different specific spot i we place a different specific gene i .
- (a) In each spot i , approximately equal abundance of cancerous and non-cancerous tissue are then placed. We can then compare the relative abundance of gene i 's DNA that binds to the cancerous versus to the non-cancerous tissue. Basically, for each gene i , this comparison leads to constructing a statistic Y_i .

For this problem set, assume that (a) a gene i is either equally expressed in both types of tissue (in which case we say $Z_i = 1$) or it is not (in which case we say that $Z_i = 0$); (b) $Z_i, i = 1, \dots, 10000$ are i.i.d. Bernoulli trials with probability $\text{pr}(Z = 1) = \pi$; (c) Y_i will be a draw from $N(0, 1)$ if gene i is equally expressed in both types of tissue, but Y_i will be a draw from $N(\mu, 1)$ (with $\mu \neq 0$) if gene i is not equally expressed in the two types of tissue; (d) (Z_i, Y_i) are independent vectors across different genes i ; (e) we observe Y_i but do not know Z_i, π or μ . We are interested in estimating $(1 - \pi)$, the proportion of genes not equally expressed, and ultimately in studying further those genes that we think are not equally expressed.

- (i) You are sent by email 10000 observations Y_i of a chip. By using problem 1 above, find moment estimates of the proportion, $1 - \pi$, of genes that are not equally expressed, and for μ .

Proof. From (i) and (ii) above, we have

$$\begin{aligned}
 E(Y^2 \mid \mu, \pi) &= 1 + [(1 - \pi)\mu]\mu \\
 &= 1 + E(Y \mid \mu, \pi)\mu \\
 \frac{E(Y^2 \mid \mu, \pi) - 1}{E(Y \mid \mu, \pi)} &= \mu
 \end{aligned}$$

and then from (i) and the solution for μ above,

$$\begin{aligned} E(Y \mid \mu, \pi) &= (1 - \pi)\mu \\ &= (1 - \pi) \frac{E(Y^2 \mid \mu, \pi) - 1}{E(Y \mid \mu, \pi)} \\ \frac{E(Y \mid \mu, \pi)^2}{E(Y^2 \mid \mu, \pi) - 1} &= 1 - \pi \end{aligned}$$

The appropriate method of moments estimates are

$$\begin{aligned} E(Y \mid \mu, \pi) &= \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y} \quad \text{and} \\ E(Y^2 \mid \mu, \pi) &= \frac{1}{n} \sum_{i=1}^n Y_i^2 = \overline{Y^2} \end{aligned}$$

Using the method of moments estimates in the equations above, estimates for the parameters of interest are

$$\begin{aligned} \mu^{MOM} &= \frac{\overline{Y^2} - 1}{\bar{Y}} \quad \text{and} \\ (1 - \pi)^{MOM} &= \frac{\bar{Y}^2}{\overline{Y^2} - 1} \end{aligned}$$

□

- (ii) Assuming π is in $(0, 1)$ but in a set bounded away from 0 or 1, find the large sample distribution of the estimator for $\log(\pi/(1 - \pi))$ resulting from (i). Hence, report an approximate 95% CI for the proportion $1 - \pi$ of genes that are not equally expressed.

Proof. Notice from question 2, part (i),

$$\begin{aligned} \frac{E(Y \mid \mu, \pi)^2}{E(Y^2 \mid \mu, \pi) - 1} &= 1 - \pi \\ \pi &= 1 - \frac{E(Y \mid \mu, \pi)^2}{E(Y^2 \mid \mu, \pi) - 1} \\ \pi &= \frac{E(Y^2 \mid \mu, \pi) - E(Y \mid \mu, \pi)^2 - 1}{E(Y^2 \mid \mu, \pi) - 1} \end{aligned}$$

Thus

$$\frac{\pi}{1 - \pi} = \frac{E(Y^2 \mid \mu, \pi) - E(Y \mid \mu, \pi)^2 - 1}{E(Y \mid \mu, \pi)^2}$$

Simplify notation by setting $\mu_1 = E(Y \mid \mu, \pi)$ and $\mu_2 = E(Y^2 \mid \mu, \pi)$. Then

$$\log \left(\frac{\pi}{1 - \pi} \right) = \log \left(\frac{E(Y^2 \mid \mu, \pi) - E(Y \mid \mu, \pi)^2 - 1}{E(Y \mid \mu, \pi)^2} \right) = \log \left(\frac{\mu_2 - \mu_1^2 - 1}{\mu_1^2} \right)$$

Calculating, $\hat{\gamma} = s_{11} - \text{var}(y_s q) s_{21} < -\text{cov}(y, y_s q) s_{12} < -\text{cov}(y_s q, y) s_{22} < -\text{var}(y) \text{sigma}_h \text{at} < -\text{matrix}(c(s_{11}, s_{21}, s_{12}, s_{22}), \text{nrow} = 2, \text{ncol} = 2) @ \text{Thus, by the Delta Method}$
 $N(0, \nabla g(\mu_2, \mu_1)^T \Sigma \nabla g(\mu_2, \mu_1))$

By Slutsky's Theorem, a consistent estimator for $\nabla g(\mu_2, \mu_1)^T \Sigma \nabla g(\mu_2, \mu_1)$ is $\hat{\Gamma} := \nabla g(\bar{Y}^2, \bar{Y})^T \hat{\Sigma} \nabla g(\bar{Y}^2, \bar{Y})$. Calculating, $\hat{\gamma} = \text{gamma}_h \text{at} < -\text{drop}(\text{gradient}_h \text{at gamma}_h \text{at} @ \text{Thus by Slutsky's Theorem})$
 $N(0, 1)$ Thus an approximate 95% confidence interval for $\log \frac{\pi}{1-\pi}$ is

$$\left[\log \left(\frac{\bar{Y}^2 - \bar{Y}^2 - 1}{\bar{Y}^2} \right) - 1.96 \sqrt{\frac{\hat{\Gamma}}{n}}, \log \left(\frac{\bar{Y}^2 - \bar{Y}^2 - 1}{\bar{Y}^2} \right) + 1.96 \sqrt{\frac{\hat{\Gamma}}{n}} \right]$$

Since

$$1 - \frac{\exp \left(\log \frac{\pi}{1-\pi} \right)}{1 + \exp \left(\log \frac{\pi}{1-\pi} \right)} = 1 - \frac{\frac{\pi}{1-\pi}}{1 + \frac{\pi}{1-\pi}} = 1 - \frac{\frac{\pi}{1-\pi}}{\frac{1}{1-\pi}} = 1 - \pi$$

then applying the function

$$h(x) = 1 - \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(x)}$$

to the confidence interval above will give a 95% confidence interval for $1 - \pi$. The desired confidence interval is

$$\left[\left(1 + \frac{(\bar{Y}^2 - \bar{Y}^2 - 1) \exp \left\{ 1.96 \sqrt{\frac{\hat{\Gamma}}{n}} \right\}}{\bar{Y}^2} \right)^{-1}, \left(1 + \frac{(\bar{Y}^2 - \bar{Y}^2 - 1) \exp \left\{ -1.96 \sqrt{\frac{\hat{\Gamma}}{n}} \right\}}{\bar{Y}^2} \right)^{-1} \right]$$

Calculating, $\hat{\gamma} = t_{1-\pi} - (y_{2bar} - y_{bar} * y_{bar} - 1) / (y_{bar} * y_{bar}) \text{leftbound} < -1 / (1 + t_{1-\pi} * \exp(qnorm(0.975) * \sqrt{\text{gamma}_h \text{at} / n})) \text{rightbound} < -1 / (1 + t_{1-\pi} * \exp(qnorm(0.025) * \sqrt{\text{gamma}_h \text{at} / n})) \text{c(leftbound, rightbound)} @ \text{Hence an approximate 95\% confidence interval for } 1 - \pi \text{ is}$

$$[\text{leftbound} \quad \text{rightbound}]$$

□

Allowing π to be in $[0, 1]$, can your CI of part (ii) be taken as evidence that there are some genes that are differentially expressed?

Proof. The above confidence interval cannot contain 0 by construction (the function $h(x)$ is positive). This confidence interval is not appropriate for testing $1 - \pi = 0$, or equivalently, $\pi = 1$, that none of the genes are differentially expressed. □

State one (alternative ?) way of testing whether or not there are some genes that are differentially expressed. Report the result from this test.

Proof. An alternative way to test that some genes are differentially expressed is to test if each gene came from the standard normal distribution. This needs multiple testing correction, since there are $n = 10000$ tests and the null hypothesis is rejected if one of those tests is significant. If the significance value is $p_0 = 0.05$, then using the Bonferroni correction, the significant p -value under multiple testing is $p_0/n = 5e-6$. $p_{critical} < -p_0/n \cdot \text{left2} < -qnorm(p_{critical}) \cdot \text{left2} < -qnorm(1-p_{critical}) \cdot \text{left2}$. Thus, any Y_i values outside of $[\text{left2}, \text{right2}]$ would be considered differentially expressed under this hypothesis test. Those values are $y[\text{which}(y \notin [\text{left2}, \text{right2}])]$. \square

Find the maximum likelihood estimate (MLE) of $1 - \pi$ and μ (e.g., use `optim` in R; in the likelihood function, parametrize π by its logit; use starting values equal to the moment estimates). If we allow π to be in $[0, 1]$, are all the conditions in (A.1) to (A.4) in p. 516 of the text satisfied? Discuss whether we should worry or not worry about this issue with our data.

Proof. The likelihood function for Y is the same as (5). Thus

$$L(\theta; y) = L(\mu, \pi; y) = \prod_{i=1}^n \text{pr}(y_i | \mu, \pi)$$

Thus, the log-likelihood function is

$$\ell(\theta; y) = \log(L(\theta; y)) = \sum_{i=1}^n \log(\phi(y_i)\pi + \phi(y_i - \mu)(1 - \pi))$$

We want to maximize the log-likelihood over the parameter space. Since we maximize numerically, we want the parameters to take values in \mathbb{R} . Hence we map $\pi \in (0, 1)$ to $\log(\pi/(1 - \pi)) = x$. The inverse mapping is $x \mapsto \frac{e^x}{1 + e^x} = \pi$ for $x \in \mathbb{R}$. With this parameterization, we maximize

$$\sum_{i=1}^n \log \left(\phi(y_i) \frac{e^x}{1 + e^x} + \phi(y_i - \mu) \frac{1}{1 + e^x} \right) \quad (1)$$

over $x, \mu \in \mathbb{R}$. This gives MLEs for x, μ . By the invariance of MLE, using the inverse mapping for x , we get the MLE for π , and hence can get the MLE for $1 - \pi$. \square

Log likelihood function $\text{getLoglik} <- \text{function}(param, y = y) \{ t1 <- \text{dnorm}(y) * \exp(param[2]) / (1 + \exp(param[2])) \}$
 $\text{t2} <- \text{dnorm}(y - param[1]) * \exp(-param[2]) / (1 + \exp(param[2]))$
 $\text{optout} <- \text{optim}(c(mu_mom, \log(pi_mom/(1-pi_mom))), \text{getLoglik}, y = y, \text{control} = \text{list}(fnscale = -1), \text{hessian} = \text{TRUE})$
 $\text{mu_hat} <- \text{optoutpar}[1]$
 $\text{pi_hat} <- \exp(\text{optoutpar}[2]) / (1 + \exp(\text{optoutpar}[2]))$
 $c(\text{mu_hat}, \text{pi_hat})$ @ Hence, the MLEs are

Next, we check the conditions for the consistency of MLEs found on pg. 516 of Casella and Berger. Do the conditions hold?

(A1) **Yes**, all the Y_i are i.i.d. and the density is $\phi(y)\pi + \phi(y - \mu)(1 - \pi)$.

- (A2) **Yes**, the parameters are identifiable. If $(\mu, \pi) \neq (\mu', \pi')$, then $\phi(y)\pi + \phi(y - \mu)(1 - \pi) \neq \phi(y)\pi' + \phi(y - \mu')(1 - \pi')$ unless $\mu = \mu' = 0$. However, by definition, the genes are differentially expressed, so μ or μ' cannot be 0 in the context of this problem.
- (A3) **Almost**. The densities have common support ($y \in \mathbb{R}$), however, the densities are differentiable in all of the parameter space except where $\pi = 1$ and $\pi = 0$ because they are on the boundary of the allowable values $[0, 1]$. It may be possible to define the derivative at $\pi = 1$ and $\pi = 0$ as the limit of the derivatives as $\pi \uparrow 1$ and $\pi \downarrow 0$, respectively.
- (A4) **No**, the true parameter may not be an element of an open set contained in the parameter space. That is because if $\pi = 0$ or $\pi = 1$, then π is on the boundary, and there is no open set containing the boundary contained in the parameter space.

However, by hypothesis, we assume that $\pi \neq 0$ and $\pi \neq 1$. That is, we assume that there are some genes that are differentially expressed, but not all of them are differentially expressed. Therefore, if it is assumed that $\pi \in (0, 1)$, then the concerns in (A3) and (A4) are resolved. We should not worry about those issues with our data. \square

Find an approximate 95% confidence interval for $1 - \pi$ based on the MLEs in part (v).

Proof. Let $\theta_0 = (\mu_0, \text{logit}(\pi_0))$, the true values of the parameters, and let $\hat{\theta} = (\hat{\mu}, \text{logit}(\hat{\pi}))$, the MLE. The information of one data point is

$$I_1(\theta_0) = -E \left(\frac{\partial^2}{\partial \theta^2} \log \text{pr}(Y | \theta) | \theta \right)$$

By the Law of Large Numbers, and by large sample theory for MLE, a consistent estimator of the information of one data point is

$$\hat{I}_1(\theta_0) = -\frac{1}{n} \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log \text{pr}(y_i | \hat{\theta})$$

According to large sample theory for the MLE,

$$\sqrt{n} (\hat{\theta} - \theta_0) \xrightarrow{D} N(0, I_1(\theta_0)^{-1})$$

Therefore, by Slutsky's theorem,

$$\sqrt{n} \hat{I}_1(\theta_0)^{1/2} (\hat{\theta} - \theta_0) \xrightarrow{D} N(0, I)$$

where $I_1(\theta_0)^{1/2}$ is the Cholesky decomposition of $I_1(\theta)$. The output of `optim` returns the Hessian at the MLE, i.e.

$$H(\hat{\theta}) = \sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log \text{pr}(y_i | \hat{\theta})$$

Notice

$$\sqrt{n}\hat{I}_1(\theta_0)^{1/2} = \left(-\sum_{i=1}^n \frac{\partial^2}{\partial \theta^2} \log \text{pr}(y_i | \hat{\theta}) \right)^{1/2} = (-H(\hat{\theta}))^{1/2}$$

Therefore,

$$(-H(\hat{\theta}))^{1/2}(\hat{\theta} - \theta_0) \xrightarrow{D} N(0, I)$$

and for the purposes of creating confidence intervals,

$$\hat{\theta} - \theta_0 \approx N(0, (-H(\hat{\theta}))^{-1})$$

Thus, an approximate 95% confidence interval for $\log \frac{\pi}{1-\pi}$ is

$$\left[\text{logit}(\hat{\pi}) - 1.96\sqrt{[(-H(\hat{\theta}))^{-1}]_{22}}, \text{logit}(\hat{\pi}) + 1.96\sqrt{[(-H(\hat{\theta}))^{-1}]_{22}} \right]$$

and applying $h(x)$, a 95% confidence interval for $1 - \pi$ is

$$LB = \left(1 + \exp \left\{ \text{logit}(\hat{\pi}) + 1.96\sqrt{[(-H(\hat{\theta}))^{-1}]_{22}} \right\} \right)^{-1}$$

$$UB = \left(1 + \exp \left\{ \text{logit}(\hat{\pi}) - 1.96\sqrt{[(-H(\hat{\theta}))^{-1}]_{22}} \right\} \right)^{-1}$$

where LB and UB are the lower bound and upper bound, respectively. \square

We will set a cutoff c , so that for all genes with $Y_i > c$, we will proceed to do further experiments. Define the positive predictive value $PPV(c, \mu, \pi)$ to be the fraction, among the genes we get with this cutoff, that are truly differentially expressed ($Z_i = 0$), and define the negative predictive value $NPV(c, \mu, \pi)$ to be the fraction, among the genes we leave out of this cutoff, that are equally expressed ($Z_i = 1$). Use the MLEs for the parameters μ, π . Among all pairs of $(PPV(c), NPV(c))$ with a $PPV > 90\%$, what should c be to give the highest NPV ? Hence state which genes from the data set you would pick based on that c .

Proof. Recall

$$\begin{aligned} PPV(c) &= P(Z_i = 0 | Y_i > c) \\ &= \frac{P(Y_i > c | Z_i = 0)P(Z_i = 0)}{P(Y_i > c | Z_i = 0)P(Z_i = 0) + P(Y_i > c | Z_i = 1)P(Z_i = 1)} \\ &= \frac{\Phi(\mu - c)(1 - \pi)}{\Phi(\mu - c)(1 - \pi) + \Phi(-c)\pi} \end{aligned}$$

and

$$\begin{aligned} NPV(c) &= P(Z_i = 1 | Y_i \leq c) \\ &= \frac{P(Y_i \leq c | Z_i = 1)P(Z_i = 1)}{P(Y_i \leq c | Z_i = 1)P(Z_i = 1) + P(Y_i \leq c | Z_i = 0)P(Z_i = 0)} \\ &= \frac{\Phi(c)\pi}{\Phi(c)\pi + \Phi(c - \mu)(1 - \pi)} \end{aligned}$$

It is clear that PPV is increasing in c and NPV is decreasing in c . Thus, the minimum c such that $PPV \geq 0.9$ gives the highest NPV among the c such that $PPV > 0.9$. That value is

$$PPV(c^*) =$$

$$NPV(c^*) =$$

□