



ژنومیک محاسباتی

مطهری و فروغمند
پاییز ۱۴۰۰

امتحان پایان ترم

پاسخ سوال های امتحان پایان ترم
نگارنده: الهه بدلی - ۹۸۲۰۹۰۷۲

۱ پاسخ سوال ۱

شرایط گفته شده سوال، شرایط تعادل هاردی-وینبرگ است که جمعیت زیاد است. اگر دیگر شرایط این تعادل از جمله تصادفی بودن ازدواج و نبودن نیروهای انتخاب طبیعی و اینکه جهش نداشته باشیم و نسل ها قاطی نباشند یعنی ازدواج در هر نسل بین افراد همان نسل صورت گیرد. در این صورت طبق تعادل هاردی-وینبرگ می توان گفت بعد از یک نسل به تعادل می رسد و در همان تعادل باقی می ماند.

$$P[AA] = (f_1 + \frac{f_2}{4}) = f_A^2 = p^2$$

$$P[aa] = (f_3 + \frac{f_2}{4}) = f_a^2 = q^2$$

$$P[aA] = 1 - f_A^2 - f_a^2 = 2(f_1 + \frac{f_2}{4})(f_3 + \frac{f_2}{4}) = 2pq$$



۲ پاسخ سوال ۲

با در نظر V_i به عنوان نود i ام، V_i توزیع دو جمله ای دارد.

$$P[V_i = k] = \binom{2N}{k} \left(\frac{1}{2N}\right)^k \left(1 - \frac{1}{2N}\right)^{2N-k}$$

$$E[V_i | X_0 = 1] = 2Np = 2N * \frac{1}{2N} = 1$$

یعنی امید داریم در نسل بعد ۱ فرزند ببینیم.

$$Var[V_i] = 2Np(1-p) = 2N\left(\frac{1}{2N}\right)\left(1 - \frac{1}{2N}\right) = 1 - \frac{1}{2N}$$

$$Corr(V_i, V_j) = \frac{Cov(V_i, V_j)}{\sqrt{Var[V_i]Var[V_j]}}$$

$$Cov(V_i, V_j) = E[V_i, V_j] - E[V_i]E[V_j]$$

$$= E\left[\sum_{i=1}^{2N} V_i \sum_{i=1}^{2N} V_j\right] - 1 * 1$$

$$= \sum_{i \neq j} E[V_i]E[V_j] + \sum_{i=j} E[V_j, V_j] - 1$$

i و j را دو فرد متفاوت در نظر می گیریم. بنابراین جمع دوم صفر است.

$$= \frac{(2N)^2 - 2N}{(2N)^2} - 1 = -\frac{1}{2N}$$

بنابراین برای ضریب هم بستگی داریم:

$$Corr(V_i, V_j) = \frac{-\frac{1}{2N}}{\sqrt{\left(1 - \frac{1}{2N}\right)^2}} = -\frac{1}{2N-1}$$

قسمت دوم:

در این صورت والدها از بین N نفر انتخاب می شوند. بنابراین احتمال والد مشترک برابر است با: $\frac{1}{N}$. بنابراین داریم:

$$P[V_i = k] = \binom{N}{k} \left(\frac{1}{N}\right)^k \left(1 - \frac{1}{N}\right)^{N-k}$$

$$E[V_i] = N * p = 1 \quad Var[V_i] = Np(1-p) = \frac{N}{N} \left(1 - \frac{1}{N}\right) = 1 - \frac{1}{N}$$

بنابراین برای هم بستگی داریم:

$$Corr(V_i, V_j) = \frac{Cov(V_i, V_j)}{\sqrt{\left(1 - \frac{1}{N}\right)^2}} = \frac{-\frac{1}{N}}{1 - \frac{1}{N}} = -\frac{1}{N-1}$$

۳ پاسخ سوال ۳

قسمت اول: k برگ داریم: بدون داشتن هیچ شرطی محاسبات را انجام می دهیم. هر یال درخت $coalsence$ یک توزیع نمایی با پارامتر متناسب با تعداد افراد آن مرحله است. بنابراین:

$$E[] = \frac{1}{\binom{k}{1}} + \frac{1}{\binom{k-1}{1}} + \dots + \frac{1}{\binom{2}{1}} = \frac{2}{k(k-1)} + \frac{2}{(k-1)(k-2)} + \dots + 1 = \frac{2}{k-1} - \frac{2}{k} + \frac{2}{k-2} - \frac{2}{k-1} + \dots + \frac{2}{2} = 2 - \frac{2}{k} = \frac{2k-2}{k}$$



که به صورت حدی برابر ۲ است. یعنی ۲ برابر جمعیت. این عدد برای $k = 2$ این عدد برابر ۱ است.

قسمت دوم:

در صورتی که ۲ فرد از نصف والد‌ها حق انتخاب برای جد مشترک داشته باشند، آن گاه: احتمال جد مشترک: $\frac{1}{N}$ می‌شود.

بنابراین:

$$p = \frac{\binom{k}{2}}{N}, \quad a = \frac{p}{M} a = \frac{\binom{k}{2}}{NM}$$

در صورتی که فرض کنیم $M = \frac{1}{N}$ محاسبات مشابه قبل خواهد بود.

۴ سوال ۴

قسمت ۱:

در مدل WF اگر در نسل اول فرکانس A برابر $\frac{x}{2N}$ باشد، وقتی در زمان پیش می‌رویم $fixation$ اتفاق می‌افتد.

اثبات: با فرض $X_0 = \frac{x}{2N}$ آن‌گاه فرکانس الیل A در نسل بعد از توزیع دو جمله‌ای است:

$$P[X_1 = k | X_0 = \frac{x}{2N}] = \binom{2N}{k} \left(\frac{x}{2N}\right)^k \left(1 - \frac{x}{2N}\right)^{2N-k}$$

$$E[X_1 | X_0 = \frac{x}{2N}] = \frac{x}{2N}$$

یعنی انتظار داریم همین نسبت در مراحل بعدی حفظ شود. همچنین این تعریف، یک مارتن‌گل است. $E[X_n | X_0 = \frac{x}{2N}] = \frac{x}{2N}$. حال $fixation$ یا قبل از لحظه n ام اتفاق افتاده است یا بعد از آن لحظه. بنابراین امید را تفکیک می‌کنیم.

$$E[X_n | X_0 = \frac{x}{2N}, fixation\ before\ n] P[fixation\ before\ n] + E[X_n | X_0 = \frac{x}{2N}, fixation\ after\ n] P[fixation\ after\ n]$$

که اگر n زیاد باشد، $fixation$ قبل از n رخ داده است. بنابراین:

$$\frac{x}{2N} = E[X_n | X_0 = \frac{x}{2N}, fixation\ before\ n]$$

حال $fixation$ یا برای a اتفاق افتاده است که امید آن صفر است و یا برای A اتفاق افتاده است که امید آن برابر ۱ است. بنابراین:

$$P[fixation\ is\ A] = \frac{x}{2N}$$

بنابراین احتمال $fixation$ برای هر الیل برابر فرکانس آن در نسل اول است.

قسمت ۲:

فرض کنیم U و V دو متغیر تصادفی با توزیع نمایی به ترتیب با پارامترهای λ_1 و λ_2 است. بنابراین:

$$P[U > t] = e^{-\lambda_1 t}$$

و

$$P[U > t] = e^{-\lambda_1 t}$$

حال:

$$Z = \min(U, V)$$

$$P[\min(U, V) > t] = P[U > t, V > t]$$



U و V مستقل اند. بنابراین:

$$= P[U > t]P[V > t] = e^{(\lambda_1 + \lambda_2)t}$$

قسمت ۳:

از منظر ۲: با فرض اینکه جهش با نرخ θ رخ بدهد، آن گاه تعداد جهش نمایی با پارامتر $k\theta$ است. در این صورت داریم:

$$P[S_n = s] = P[S_n = s | \text{mutation first}]P[\text{mutation first}] + P[S_n = s | \text{coalescence first}]P[\text{coalescence first}] =$$

$$P[S_n = s - 1] \cdot \frac{k\theta}{\binom{k}{2} + k\theta} + P[S_{n-1} = s] \cdot \frac{\binom{k}{2}}{\binom{k}{2} + k\theta}$$

البته در منابع کلاس نرخ جهش، $\frac{\theta}{4}$ در نظر گرفته شده است.

قسمت ۴:

معمولا مدل هایی مثل مدل WF فرض هایی دارند که واقعی نیستند. مثلا تعداد جمعیت همیشه ثابت است یا ازدواج تصادفی است یا مثلا تعداد مرد و زن یکسان است. اما در واقعیت اینگونه نیست. در واقعیت جمعیت موثر مطرح می شود که از جمعیت اصلی کوچک تر است و با تقریب خوبی مدل هایی که قبلا داشتیم را می توان استفاده کرد. برای مثال می توان از روش یافتن جمعیت موثر بر پایه واریانس استفاده کرد. در مدل WF ، با فرض X تعداد A ها در این نسل، داریم:

$$P[X = k] = \binom{2N}{k} p^k (1-p)^{2N-k}$$

$$E[X] = 2Np$$

$$Var[X] = 2Np(1-p)$$

$$p' = \frac{X}{2N}$$

$$E[p'] = \frac{E[X]}{2N} = \frac{2Np}{2N} = p$$

$$Var[p'] = \frac{Var[X]}{(2N)^2} = \frac{2Np(1-p)}{(2N)^2} = \frac{p(1-p)}{2N}$$

بنابراین:

$$2N_{eff} = \frac{p(1-p)}{V[p']}$$

بنابراین با داشتن تخمینی از واریانس می توان N_{eff} را محاسبه نمود.

سوال ۵:

سوال ۱:

تخمین واترسون: برای تخمین زدن نرخ جهش استفاده می شود. با حل از منظر ۱: تعداد واقعی جهش ها روی درخت به ازای n نفر را S_n در نظر می گیریم.

$$P[S_n = s] = \frac{n-1}{\theta} \sum_{i=1}^{n-1} (-1)^{(i-1)} \binom{n-2}{i-1} \left(\frac{\theta}{i+\theta}\right)^{s+1}$$

با به دست آوردن $E[S_n]$ می توان تخمین واترسون را به دست آورد.

در صورتی که تعداد جهش ها در شاخه ی n ام را X_n در نظر بگیریم.

$$S_n = X_n + X_{n-1} + \dots + X_1$$

تعداد جهش ها در هر شاخه از توزیع پواسون هستند.

$$E[S_n] = E[X_n] + E[X_{n-1}] + \dots + E[X_1]$$

و داریم:

$$E[S_k] = E_t[E[S_k|t]]$$

با توجه به اینکه kt تا فاصله روی درخت داریم و تعداد جهش ها دارای توزیع پواسون هستند، بنابراین:

$$E[S_k|t] = \sum_{i=0}^{\infty} i \binom{k}{i} e^{-t} t^i$$

بنابراین:

$$E[S_k] = \sum_{i=0}^{\infty} i \binom{k}{i} E[t] e^{-t} t^i = \frac{\sum_{i=0}^{\infty} i \binom{k}{i} t^i}{\sum_{i=0}^{\infty} \binom{k}{i} t^i} = \frac{k}{k-1}$$

زیرا با فرض $\frac{\theta}{k} = \sum_{i=0}^{\infty} i \binom{k}{i} e^{-t} t^i$ و با فرض $x = 1$ داریم: $\theta = \sum_{i=0}^{\infty} i \binom{k}{i} e^{-t} t^i$. بنابراین:

$$\frac{\theta}{k-1} = \frac{\theta}{k-1}$$

حال با جمع زدن برای n های مختلف داریم:

$$\frac{\theta}{n-1} + \frac{\theta}{n-2} + \dots + \frac{\theta}{1}$$

که برابر است با:

$$E[S_k] = \theta \sum_{i=1}^{n-1} \frac{1}{i} = \theta h_n$$

که در آن: $h_n = \log n$

سوال ۲:

در روش $TajimaD$ ابتدا فاصله جفت جفت رشته ها را به دست می آوریم. فرض می کنیم d_{ij} فاصله ۲ رشته است. این فاصله با

تعداد جهش های روی ۲ شاخه برابر است.

حال با فرض وجود k نفر، داریم:

احتمال اینکه اول جهش داشته باشیم:

$$\frac{\frac{\theta k}{2}}{\frac{\theta k}{2} + \binom{k}{2}}$$

برای دو نفر داریم: احتمال اینکه اول جهش داشته باشیم: $\frac{\theta}{\theta+1}$ و احتمال اینکه اول $coalsence$ برابر $\frac{1}{\theta+1}$ است.

حال فاصله این دو نفر را این گونه در نظر می گیریم که چند جهش تا قبل از $coalsence$ دارند.

$$P[d_{ij} = k] = \left(\frac{\theta}{\theta+1}\right)^k \frac{1}{1+\theta}$$

یعنی k بار جهش رخ داده است و در نهایت $coalsence$ انجام شده است.

$$E[d_{ij}] = \frac{1-p}{p}, \quad k = 0, 1, 2, \dots$$

$$E[d_{ij}] = \frac{1 - \frac{1}{1+\theta}}{\frac{1}{1+\theta}} = \frac{\theta}{1+\theta} = \theta$$

و برای تخمین $TajimaD$ نیز داریم:

$$E[\hat{\theta}_T] = \frac{\sum_{i,j} E[d_{ij}]}{\binom{n}{2}} = \theta$$

سوال ۳:

استفاده از $TajimaD$ به این صورت است که در حالتی که جهش‌ها خنثی باشند باید نرخ جهش تخمین زده شده توسط $Waterson$ و $TajimaD$ یکی باشند. بنابراین باید اختلافشان صفر باشد و D در این نام هم به همین معنا است.