

Week 2 QCB 408 / 508 Spring 2020

Continuous RV's (continued)

Normal (Gaussian) distribution

$$X \sim \text{Normal}(\mu, \sigma^2)$$

$$\mathcal{R} = (-\infty, \infty)$$

$$f(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

$$E[X] = \mu, \quad \text{Var}(X) = \sigma^2 > 0$$

$$dnorm(x, \text{mean}=2, \text{sd}=3)$$

Multivariate RV's

Two RV's X and Y

$$\text{Joint cdf } F(a, b) = \Pr(X \leq a, Y \leq b)$$

$$= \Pr(\{\omega : X(\omega) \leq a\} \cap \{\omega : Y(\omega) \leq b\})$$

Suppose continuous.

$$\frac{d^2}{dx dy} F(x, y) = f(x, y) \rightarrow \text{joint pdf}$$

Suppose discrete

joint pmf $f(x, y) = \Pr(X=x, Y=y)$

Marginal distributions

$$f(x) = \sum_{y \in \mathcal{R}_y} f(x, y)$$

$$f(x) = \int f(x, y) dy$$

See Law of Total Probability.

Independent rvs

If X and Y are independent

then:

$$f(x, y) = f(x)f(y) \quad \checkmark$$

$$\begin{aligned} F_{X, Y}(a, b) &= \Pr(X \leq a, Y \leq b) \\ &= \Pr(X \leq a) \Pr(Y \leq b) \\ &= F_X(a) F_Y(b) \end{aligned}$$

Conditional Distribution:

$$\rightarrow \Pr(X \leq a \mid Y \leq b) = \frac{\Pr(X \leq a, Y \leq b)}{\Pr(Y \leq b)}$$

$$\{\omega : X(\omega) \leq a\}$$

$$\{\omega : Y(\omega) \leq b\}$$

$$F_{X|Y}(a \mid Y \leq b) = \frac{F_{X,Y}(a,b)}{F_Y(b)}$$

$$f(x|y) = \frac{f(x,y)}{f(y)}$$

$$\text{Bayes Theorem } f(y|x) = \frac{f(x|y) f(y)}{f(x)}$$

All of the above extends to

vars X_1, X_2, \dots, X_n

$$f(x_1, x_2 \mid x_3, x_4, x_5) = \frac{f(x_1, x_2, x_3, x_4, x_5)}{f(x_3, x_4, x_5)}$$

$$f(x_3, x_4, x_5) = \iint f(x_1, x_2, x_3, x_4, x_5) dx_1 dx_2$$

If X_1, X_2, \dots, X_n are independent then

$$f(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i)$$

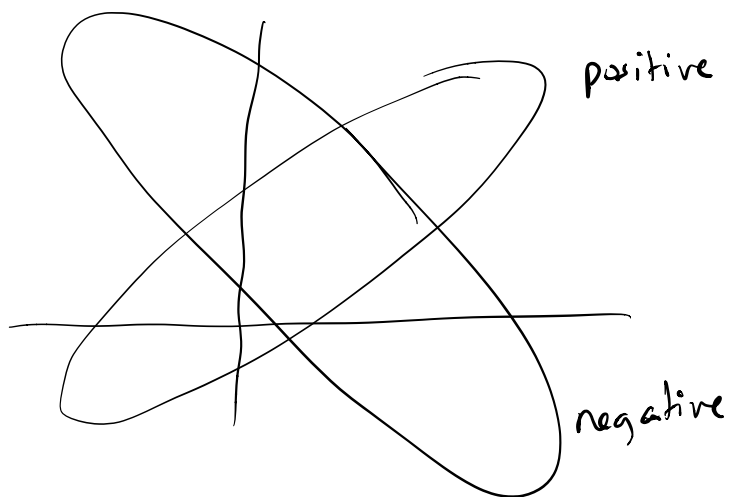
Moments of Joint Distributions

For a single rv, X :

$E[X^k]$ is the k^{th} moment

$$\begin{aligned} \text{Var}(X) &= E[(X - E[X])^2] \\ &= E[X^2] - E[X]^2 \end{aligned}$$

$$\text{Cov}(X, Y) = \underline{E[(X - E[X])(Y - E[Y])]} \checkmark$$



First set $E[X], E[Y]$

$$\text{Cov}(X, Y) = \iint (x - E[X])(y - E[Y]) f(x, y) dx dy$$

or

$$= \sum_x \sum_y (x - E[X])(y - E[Y]) \underbrace{f(x, y)}_{\downarrow}$$

$$\text{Pr}(X=x, Y=y)$$

Linear Transformations of RV's

X is a RV

a and b are constants

$$E[a + bX] = a + bE[X]$$

$$\text{Var}(a + bX) = b^2 \text{Var}(X)$$

$$E[\underline{a + bX}] = \int (a + bx) f(x) dx$$

$$= \int a f(x) dx + \int bx f(x) dx$$

$$= a \underbrace{\int f(x) dx}_1 + b \underbrace{\int x f(x) dx}_{E[X]}$$

$$= a + bE[X]$$

Exercise: $\text{Var}(a+bX) = b^2 \text{Var}(X)$

Law of Total Variance

Jointly distributed r.v's X and Y

→ $\text{Var}(X) = \text{Var}(E[X|Y]) + E[\text{Var}(X|Y)]$

$E[X|Y] \rightarrow ?$

$E[X|Y=y] = \underbrace{\int x f(x|y) dx}_{\text{function of } y}$

Have a pmf or pdf for y

$E[X|Y] = \underbrace{\int x f(x|Y) dx}_{\text{function of } Y}$

$E[E[X|Y]] = \int \int \underbrace{x f(x|y)}_{E[X|Y=y]} f(y) dy$
 $= \int \int \underline{x f(x|y) f(y)} dx dy$

$$\begin{aligned}
&= \iint x f(x,y) dx dy \\
&= \int \int x f(x,y) dy dx \\
&= \int x \underbrace{\int f(x,y) dy}_{f(x)} dx \\
&= \int x f(x) dx \\
&= E[X]
\end{aligned}$$

$$f(x) = \int f(x,y) dy$$

$$\text{Var}(X|Y) = E[(X - E[X|Y])^2 | Y]$$

Exercise : Prove Law of Total Variance
 ↓
 or
 read

Hardy-Weinberg Equilibrium

SNP CC, CT, TT
 ↓ ↓ ↓
 0 1 2

Genotype $X \in \{0, 1, 2\}$

$$\Pr(X=k) = r_k \quad k=0, 1, 2$$

$$E[X] = r_1 + 2r_2 \quad \underline{E\left[\frac{X}{2}\right] = \frac{r_1}{2} + r_2 \equiv p}$$

freq. of T

$Y \in \{0, 1\}$ is the transmitted

$$\Pr(Y=1 | X=k) = \begin{cases} 0 & k=0 \\ 1/2 & k=1 \\ 1 & k=2 \end{cases}$$

$$\Pr(Y=1) = \sum_{k=0}^2 \Pr(Y=1 | X=k) \Pr(X=k)$$

$$= 0 \cdot r_0 + \frac{1}{2} r_1 + r_2$$

$$= p$$

$$\Rightarrow Y \sim \text{Bernoulli}(p)$$

$Z \in \{0, 1, 2\}$ is next generation genotype

$Y_1, Y_2 \stackrel{iid}{\sim} \text{Bernoulli}(p)$

iid = independent, identically distributed

$$Z = Y_1 + Y_2$$

$$\Rightarrow Z \sim \text{Binomial}(2, p)$$

$$E[Z] = 2p = E[X] \Rightarrow \text{transmission probs for } Z \text{ are identical to } X$$

\Rightarrow all ^{future} generations' genotypes are drawn from Binomial(2, p) \Rightarrow equilibrium

$$Pr(Z=0) = (1-p)^2$$

$$Pr(Z=1) = 2p(1-p) \leftarrow$$

$$Pr(Z=2) = p^2$$

Typically $\Pr(Z=1) < 2p(1-p)$ when HWE is violated.

Inbreeding

IBD = identical by descent

Two alleles are IBD if they are copies from a common ancestor

Let I be a rv such that it indicates whether two randomly drawn alleles are IBD or not.

Y_1 and Y_2 :

$$\begin{aligned} \longrightarrow I=1 &\Rightarrow Y_1 = Y_2 \sim \text{Bernoulli}(p) \\ I=0 &\Rightarrow Y_1, Y_2 \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p) \end{aligned}$$

$$I \sim \text{Bernoulli}(f)$$

f is the inbreeding coefficient (usually F)

$$Z | I=0 \sim \text{Binomial}(2, p)$$

$$Z/2 | I=1 \sim \text{Bernoulli}(p)$$

$Y_1, Y_2, Y_3 \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

$Z | I=0 \sim \text{Binomial}(2, p)$

$Z | I=1 \sim \text{Bernoulli}(p)$

$$Z = (Y_1 + Y_2)(1 - I) + 2Y_3 I$$

HNÉ →

$$\Pr(Z = k | I = 0) = \begin{cases} (1-p)^2 & k=0 \\ 2p(1-p) & k=1 \\ p^2 & k=2 \end{cases}$$

$$\Pr(Z = k | I = 1) = \begin{cases} 1-p & k=0 \\ 0 & k=1 \\ p & k=2 \end{cases}$$

i.e., $Z | I=1 \sim 2\text{-Bernoulli}(p)$

$$\Pr(Z=0) = \Pr(Z=0 | I=0) \Pr(I=0) + \Pr(Z=0 | I=1) \Pr(I=1)$$

$$= (1-p)^2 (1-f) + (1-p) f$$

$$= (1-p)^2 + p(1-p) f$$

$$\Pr(Z=1) = 2p(1-p)(1-f) + 0 \cdot f$$

$$= 2p(1-p)(1-f) < 2p(1-p)$$

$$\Pr(Z=2) = p^2(1-f) + p f$$

$$= p^2 + p(1-p) f$$

$$E[Z] = E[E[Z|I]] =$$

$$\text{Note: } 2p(1-p) = 2p(1-p)(1-f) +$$

$$2p(1-p) f$$

$$= 2p(1-p)(1-f) +$$

$$\rightarrow p(1-p) f +$$

$$\rightarrow p(1-p) f$$

$$\text{Var}(Z) = \underline{E[\text{Var}(Z|I)]} + \text{Var}(E[Z|I])$$

$$\text{Var}(Z | I=0) = 2p(1-p) \quad \{\text{Binomial}(2, p)\}$$

$$\begin{aligned}\text{Var}(Z | I=1) &= \text{Var}(2Y_3) \\ &= 4\text{Var}(Y_3) \\ &= 4p(1-p)\end{aligned}$$

$$\text{Var}(Z | I) = 2p(1-p)(1-I) + 4p(1-p)I$$

$$\begin{aligned}E[\text{Var}(Z|I)] &= 2p(1-p)(1-f) + 4p(1-p)f \\ &= 2p(1-p)(1+f)\end{aligned}$$

$$Z | I=0 = Y_1 + Y_2$$

$$Z | I=1 = 2Y_3$$

$$E[Z | I=0] = 2p$$

$$E[Z | I=1] = 2E[Y_3] = 2p$$

$$E[Z | I] = 2p$$

$$\text{Var}(E[Z|I]) = 0$$

$$\Rightarrow \text{Var}(Z) = 2p(1-p)(1+f)$$

$$f = \frac{\text{Var}(Z) - \text{Var}(Z|I=0)}{\text{Var}(Z|I=0)}$$

$$f = 1 - \frac{\text{Pr}(Z=1)}{\text{Pr}(Z=1|I=0)}$$

Drift

Make allele frequency random.

Then allow for HWE style mating.

$$Z|Q \sim \text{Binomial}(2, Q)$$

Q is a random variable

$$Z|Q=q \sim \text{Binomial}(2, q) \quad \checkmark$$

Q changes according to a distribution

p = ancestral allele frequency

f = fixation index, "inbreeding"

$$Q \sim \text{Beta} \left(\underset{\alpha}{\frac{1-f}{f} p}, \underset{\beta}{\frac{1-f}{f} (1-p)} \right)$$

$$= \text{BN}(p, f)$$

Balding-Nichols

$$Q \sim \text{BN}(p, f)$$

$$E[Q] = p, \quad \text{Var}(Q) = p(1-p)f \quad \checkmark$$

$$E[Z] = E[E[Z|Q]] = E[2Q] = 2p$$

$$\text{Pr}(Z=2) = \int \text{Pr}(Z=2|Q=q) f(q) dq$$

$$= \int q^2 f(q) dq$$

$$= E[Q^2]$$

$$= \text{Var}(Q) + E[Q]^2$$

$$= p(1-p)f + p^2$$

$$\begin{aligned} \text{Var}(Q) &= E[Q^2] - E[Q]^2 \end{aligned}$$

$$\text{Verify: } P_r(Z=0) = p(1-p)f + (1-p)^2$$

$$P_r(Z=1) = 2p(1-p)(1-f)$$

$$\rightarrow \text{Var}(Z) = E[\text{Var}(Z|Q)] + \text{Var}(E[Z|Q])$$

$$= E[2Q(1-Q)] + \text{Var}(2Q)$$

$$= E[2Q] - E[2Q^2] + 4p(1-p)f$$

$$= 2p - 2E[Q^2] + 4p(1-p)f$$

$$= 2p - 2[\text{Var}(Q) + E(Q)^2] + 4p(1-p)f$$

$$= 2p - 2p(1-p)f - 2p^2 + 4p(1-p)f$$

$$= 2p(1-p) + 2p(1-p)f$$

$$= 2[p(1-p) + p(1-p)f]$$

$$= 2p(1-p)(1+f)$$

Model of RNA-seq Data

Assume $i = 1, 2, \dots, m$ genes

$j = 1, 2, \dots, n$ observations

Consider a single biological condition or population.

Observe RNA-seq read counts Y_{ij} for gene i in observation j .

Target: True proportion of gene i expression, call it a_i

The proportion of mRNA for gene i in this population is a_i .

$$\sum_{i=1}^m a_i = 1, \quad \text{most } a_i \text{ are small.}$$

Instructive, idealized example to follow.

Step 1. Sample cells and mRNA molecules

Step 2. Sequence mRNA molecules and obtain counts.

Step 1:

M_j is the number of mRNA molecules sampled for observation j

X_{ij} is the number of mRNA molecules sampled for gene i in observation j

Assume completely random sampling of mRNA molecules.

$$\Rightarrow X_{ij} | M_j \sim \text{Binomial}(M_j, a_i)$$

M_j is large, a_i is small

$\Rightarrow X_{ij} \dot{\sim} \text{Poisson}(M_j a_i)$

$$\text{Var}(X_{ij} | M_j) = M_j a_i \underbrace{(1 - a_i)}_{\approx 1}$$

$$\pi_{ij} = \frac{X_{ij}}{M_j} \quad \text{random proportion of gene } i \text{ mRNA in observation } j$$

$E[\pi_{ij}] =$ to be continued next week (week 3)...