

Lec 1: March 28

- 753 Final exam notes

Q6: 2x2 table:

		Normal	Cancer
Smoker	55	40	
	45	60	

Poisson log-linear

$Z \sim \text{smoking} + \text{cancer}$

is smoking significantly associated w/ lung cancer?

· relationship between 2x2 table and log-linear

· independence model

→ each row has a row effect, each column has a column effect

→ test for association between rows and columns

by comparing indep. model to fully saturated model (residual deviance) $1 - \text{pchisq}$

found in R output

→ can't just look at z-values or p-values

· when can we use p-values?

b.) prospective probability that a smoker randomly sampled from the country has cancer.

→ just look at 2x2 table

→ write down log-odds from logistic regression

$$= \beta_0 + \beta_1 I(\text{smoking})$$

Non-smoker: $\log(40/55)$ } use this to recover log
smoker: $\log(60/45)$ } odds from prospective trial
} (b/c coefficients are same w/
retrospective except for
intercept)

can use i.'s in question to compute the intercept's offset

Q7: a.) residual deviance >> residual df
 \Rightarrow overdispersion

b.) to know we associated w/ disease risk?

Remember there's overdispersion!

$$\hat{\phi} \approx \frac{D}{n-p} = \frac{211.16}{85}$$

- need to use adjusted SE to construct t-statistic
- compare w/ t-dist w/ 85 df

c.) CI:

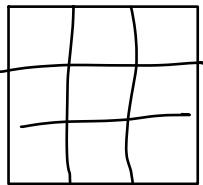
d.) intercept-only model vs. current model

- use null-model
- OVERDISPERSION !!
- construct F-statistic

$$\frac{(Resid - Null) / (n_1 - n_0)}{\hat{\phi}} \sim F$$

· Null deviance: deviance of intercept-only model |

Fixed Effects and Random Effects



· Effects: describe the impact of a level of the factor on the outcome/response (regression coefficients) → parameters

· Fixed effects: fixed constants

· Random effects: parameters are random

· Linear Mixed Model (LMM):

$$\begin{aligned}m_{ij} &= \mu + \alpha_i + \beta_j \\&\text{fixed} \quad \text{random} \\y_{ij} &= m_{ij} + \varepsilon_{ij}\end{aligned}$$

· GLM:

$$g(m_{ij}) = \mu + \alpha_i + \beta_j \quad \text{fixed constants}$$

· GLMM:

$$g(m_{ij}) = \mu + \alpha_i + \beta_j \quad \text{random}$$

· When to use fixed vs. random effects models?

EX:

(1) fixed effects

y_{ij} — outcome of patient j receiving drug i :
 $i=1$: placebo
 $i=2$: new drug

Evaluate whether new drug is better

Assume fixed effects model:

$$E(y_{ij}) = \mu_i = \mu + \alpha_i$$

drug-specific effect,
unknown fixed constant

(2) Random effects:

y_{ij} - outcome for patient j at hospital i

$$E(y_{ij}) = \mu + \alpha_i$$

predict
random variable, follows some dist.

Sometimes, main interest is to study variability across hospitals, not each individual hospital

i.e. $\alpha_i \sim N(0, \sigma^2)$

variability across hospitals

- also used for modelling correlated data

Inference

(1) Estimate

- Fixed effects: β
- random effects: β, σ^2

(a) Maximum likelihood approach

(b) REML - restricted maximum likelihood approach

- get rid of fixed effects parameter
e.g. integrate them out of likelihood

- try to estimate variance parameters

$$y = X\beta + Zu + \varepsilon$$

β fixed
 u RV
 ε noise

$$\kappa'y = \kappa'X\beta + \kappa'Zu + \kappa'\varepsilon$$

Find κ s.t. $\kappa X = 0$ $\sim N(0, \Sigma)$

$$\hat{y} = \kappa'y = \underbrace{\kappa'Zu}_{\text{fixed effects disappear}} + \kappa'\varepsilon \sim N(0, \sigma^2)$$

ex: $y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$$\text{ML: } \hat{\sigma}^2 = \frac{\sum (y_i - \bar{y})^2}{n}$$

$$\text{REML: } \hat{\sigma}^2 = \frac{\sum (y_i - \bar{y})^2}{n-1} \rightarrow \text{automatically accounts for df from fixed effects } \beta$$

(c) Bayesian approach

$$\theta = (\beta, \sigma^2)$$

Assume some prior distribution for the parameters.

$$\pi(\theta|y) = \frac{f(y|\theta)\pi(\theta)}{\int f(y|\theta)d\theta}$$

$$\int f(y|\theta)d\theta$$

inference based on posterior distribution

(d) estimating equations { Quasi-likelihood
Generalized estimating eqns

(e) Prediction

$$y_{ij} = \mu + a_i + b_{ij}$$

$$E(y_{ij}|a_i) = \mu + a_i \sim N(0, \sigma^2_a)$$

estimation

$$E(y_{ij}) \approx$$

Best predictor:

$$\beta = E(\text{an} | y)$$

Fixed Effects Model \rightarrow random:

One-way classification model

i. Normal and fixed effect

Consider m groups, each has n_i obs.

$$y_{ij} \quad i=1, \dots, m \\ j=1, \dots, n_i$$

$$y_{ij} = \mu_i + \varepsilon_{ij} \stackrel{iid}{\sim} N(0, \sigma^2)$$

$y_{ij} \stackrel{iid}{\sim} N(\mu_i, \sigma^2)$ ← each group has a group-specific mean that we want to estimate.

(i) MLE: $\hat{\mu}_i = \bar{y}_i \rightarrow$ unbiased

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i \sum_j (y_{ij} - \bar{y}_i)^2 \rightarrow$$
 biased

$$\downarrow \quad N = \sum_{i=1}^m n_i \quad E \hat{\sigma}^2 = \frac{N-m}{N} \sigma^2 \neq \sigma^2$$

(ii) REML:

$$S^2 = \frac{1}{N-m} \sum_i \sum_j (y_{ij} - \bar{y}_i)^2 \rightarrow$$
 unbiased

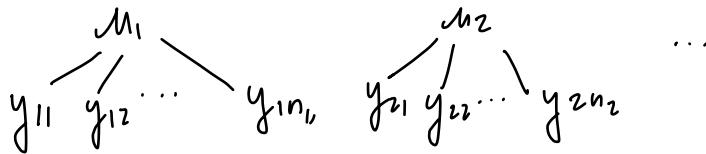
$$\frac{(N-m) S^2}{\sigma^2} \sim \chi^2_{N-m} \rightarrow \text{mean of } \chi^2_{N-m} = N-m$$

$\downarrow \text{Var} = 2/(N-m)$

$$\text{so } \text{Var}(S^2) = \frac{\sigma^2}{(N-m)^2} \cdot 2(N-m) = \frac{2\sigma^4}{N-m}$$

Sec 2: March 30

One-way classification (Normal)



1. Normal and fixed effects

$$y_{ij} \stackrel{\text{ind}}{\sim} N(\mu_i, \sigma^2)$$

(i) MLE:

$$\hat{\mu}_i = \bar{y}_i$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$$

REML:

$$\tilde{\sigma}^2 = \frac{1}{N-m} \sum_i \sum_j (y_{ij} - \bar{y}_i)^2$$

(ii) Likelihood-ratio test

$H_0: \mu_1 = \mu_2 = \dots = \mu_m$ MLE under null model

$LR \triangleq \Lambda = \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \right)^{-N/2}$ obtained by writing the likelihood of null model

MLE for variance

estimate under the

alternative model

and alternative model

within group
sum of squares

between group
sum of squares

$$\hat{\sigma}_0^2 = \frac{1}{N} \sum_i \sum_j (y_{ij} - \bar{y})^2 = \frac{1}{N} \left[\sum_i \sum_j (y_{ij} - \bar{y}_i)^2 + \sum_i \sum_j (\bar{y}_i - \bar{y})^2 \right]$$

$$\hat{\sigma}^2 = \frac{1}{N} \sum_i \sum_j (y_{ij} - \bar{y}_{ij})^2$$

within group sum of squares

$$\begin{aligned}
 -2 \log \Lambda &= N \log \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}^2} \right) \\
 &= N \log \left(1 + \frac{\sum_i \sum_j (\bar{y}_{ij} - \bar{y})^2}{\sum_i \sum_j (y_{ij} - \bar{y}_{ij})^2} \right) \\
 &= N \log \left(1 + \frac{(m-1) \sum_i \sum_j (\bar{y}_{ij} - \bar{y})^2 / (m-1)}{(N-m) \sum_i \sum_j (\bar{y}_{ij} - \bar{y}_i)^2 / (N-m)} \right) \\
 &= N \log \left(1 + \frac{m-1}{N-m} F \right)
 \end{aligned}$$

↑ F-statistic

⇒ LR test equivalent to ANOVA

(iii) Confidence intervals

CI for $\mu_i - \mu_j$

$$\bar{y}_i - \bar{y}_j \pm t_{N-m, \frac{\alpha}{2}} s \sqrt{\frac{1}{n_i} + \frac{1}{n_j}}$$

CI for a linear combination of μ_i 's: $\sum_i c_i \mu_i$

$$\sum_i c_i \bar{y}_i \pm t_{N-m, \frac{\alpha}{2}} s \sqrt{\sum_i \frac{c_i^2}{n_i}}$$

CI for σ^2

$$\frac{(N-m)s^2}{\sigma^2} \sim \chi^2_{N-m}$$
$$\Rightarrow \left(\frac{(N-m)s^2}{\chi^2_{N-m, 1-\alpha/2}}, \frac{(N-m)s^2}{\chi^2_{N-m, \alpha/2}} \right)$$

(iv) Test whether a linear combination of the group means is equal to a specific value

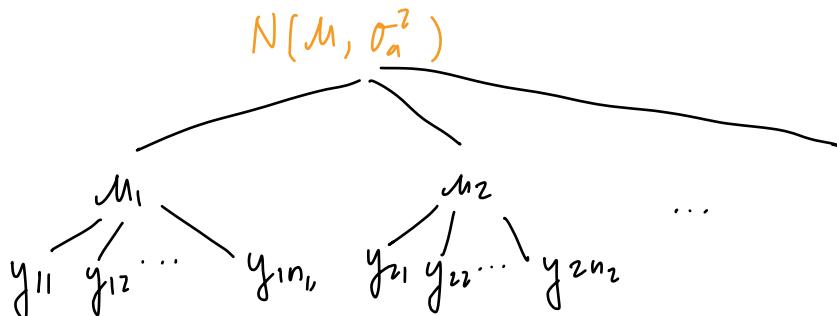
$$\text{Test } H_0: \sum_i c_i \mu_i = \mu_0$$

$$t = \frac{\sum_i c_i \bar{y}_i - \mu_0}{s \sqrt{\sum_i c_i^2 / n_i}} \sim t_{N-m}$$

determined by df of sample variance

2. Normal and Random Effects

- assume μ_i 's are drawn from a distribution



1. draw μ_i 's
2. draw y_{ij} 's from $N(\mu_i, \sigma^2_e)$

$$y_{ij} \mid \mu_i \stackrel{\text{ind}}{\sim} N(\mu_i, \sigma^2_e)$$

$$\mu_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2_a)$$

$$\mu_i = \mu + a_i$$

$$\begin{cases} y_{ij} \mid a_i \stackrel{\text{ind}}{\sim} N(\mu + a_i, \sigma^2_e) \\ a_i \stackrel{\text{iid}}{\sim} N(0, \sigma^2_a) \end{cases}$$

- this assumption introduces marginal correlation between data points in the same group

(i) covariances

$$y_{ij} = \mu + a_i + \varepsilon_{ij} \stackrel{\sim}{\sim} N(0, \sigma^2_a)$$

$$\begin{aligned} a_i &\perp \varepsilon_{ij} & \text{const} \\ \text{cov}(y_{ij}, y_{in}) &= \text{cov}(\mu + a_i + \varepsilon_{ij}, \mu + a_i + \varepsilon_{in}) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{cov}(y_{ij}, y_{in}) &= \text{cov}(\mu + a_i + \varepsilon_{ij}, \mu + a_i + \varepsilon_{in}) \\ &= \text{cov}(a_i, a_i) + \cancel{\text{cov}(a_i, \varepsilon_{in})} + \cancel{\text{cov}(\varepsilon_{ij}, a_i)} \\ &\quad + \cancel{\text{cov}(\varepsilon_{ij}, \varepsilon_{in})} & 0 & 0 \\ &= \text{var}(a_i) & 0 \\ &= \sigma^2_a \neq 0 \end{aligned}$$

- don't have independent obs. in your data
- intraclass correlations.

$$\text{Var}(y_{ij}) = \text{Var}(\mu + \alpha_i + \epsilon_{ij}) = \sigma_a^2 + \sigma^2$$

Intra-class correlation:

correlation for two points
in the same group

$$\text{Cor}(y_{ij}, y_{ik}) = \frac{\sigma_a^2}{\sqrt{(\sigma_a^2 + \sigma^2)(\sigma_a^2 + \sigma^2)}} = \frac{\sigma_a^2}{\sigma_a^2 + \sigma^2}$$

variance for
random effects

$$f(y, a | \mu, \sigma_a^2, \sigma^2)$$

(variance for
measurement error
unobserved, integrate it out)

$$\int f(y, a | \mu, \sigma_a^2, \sigma^2) da$$

$$= f(y | \mu, \sigma_a^2, \sigma^2)$$

↳ this marginal distribution leads to the non-zero covariances

σ_a^2 characterizes between-group variances

(ii) Likelihood

$$\begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ y_{22} \\ \vdots \\ y_{2n_2} \end{bmatrix}$$

divide into blocks by groups.

$$y_i \triangleq \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{in_i} \end{bmatrix}$$

$$y_i \sim N(\mu \mathbf{1}, V_i)$$

$$V_i = \begin{bmatrix} \sigma_a^2 + \sigma^2 & \sigma_a^2 & & \\ \sigma_a^2 & \ddots & \sigma_a^2 & \\ & & \ddots & \sigma_a^2 + \sigma^2 \end{bmatrix}$$

off diagonals are
covariances of
elements in the
same group

$$V_i = \sigma^2 I_{n_i} + \sigma_a^2 J_{n_i}$$

↓

$$\begin{bmatrix} 1 & 1 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix}_{n_i \times n_i} = \mathbf{1} \cdot \mathbf{1}^T$$

$$V_i^{-1} = \frac{1}{\sigma^2} I_{n_i} - \frac{\sigma_a^2}{\sigma^2(\sigma^2 + n_i \sigma_a^2)} J_{n_i}$$

$$|V_i| = (\sigma^2 + n_i \sigma_a^2) (\sigma^2)^{n_i-1}$$

- likelihood is the product of likelihoods for different groups.

$$L = \prod_{i=1}^m (2\pi)^{-n_i/2} |V_i|^{-1/2} \exp \left\{ -\frac{1}{2} (y_i - \mu \mathbf{1})^T V_i^{-1} (y_i - \mu \mathbf{1}) \right\}$$

$$\ell = \log L = -\frac{N}{2} \log(2\pi) - \frac{1}{2} \sum_i \log(\sigma^2 + n_i \sigma_a^2) - \frac{1}{2}(N-m) \log \sigma^2$$

$$- \frac{1}{2\sigma^2} \sum_i \sum_j (y_{ij} - \mu)^2 + \frac{n \sigma_a^2}{2\sigma^2(\sigma^2 + n \sigma_a^2)} \sum_i \frac{(y_i - n_i \mu)^2}{\sigma^2 + n_i \sigma_a^2}$$

- in general, no closed form solution to score equations (need iterative algorithms)

special case: all n_i 's are the same (same # obs. in each group), balanced data:

$$\ell = -\frac{N}{2} \log 2\pi - \frac{1}{2} m \log(\sigma^2 + n \sigma_a^2) - \frac{1}{2} m(n-1) \log \sigma^2$$

$$- \frac{1}{2\sigma^2} \sum_i \sum_j (y_{ij} - \mu)^2 + \frac{n \sigma_a^2 \sum_i (y_i - \mu)^2}{2\sigma^2(\sigma^2 + n \sigma_a^2)}$$

$$SSA \triangleq \sum_i n(\bar{y}_i - \bar{y})^2$$

$$SSE \triangleq \sum_i (y_{ij} - \bar{y}_i)^2$$

$$\lambda \triangleq \sigma^2 + n\sigma_a^2$$

$$l = -\frac{N}{2} \log 2\pi - \frac{1}{2} m \log \lambda - \frac{1}{2} m(n-1) \log \sigma^2$$

$$-\frac{SSE}{2\sigma^2} - \frac{SSA}{2\lambda} - \frac{mn(\bar{y}-m)}{2\lambda}$$

ML equations:

$$\begin{cases} \frac{\partial l}{\partial m} = \frac{\partial l}{\partial \bar{y}} = \frac{mn(\bar{y}-m)}{\lambda} = 0 \rightarrow \hat{m} = \bar{y} \\ \frac{\partial l}{\partial \sigma^2} = \frac{\partial l}{\partial \hat{\sigma}^2} = \frac{-m(n-1)}{\sigma^2} + \frac{SSE}{2\sigma^4} = 0 \rightarrow \hat{\sigma}^2 = \frac{SSE}{m(n-1)} < df \\ \frac{\partial l}{\partial \lambda} = \frac{\partial l}{\partial \hat{\lambda}} = \frac{-m}{2\lambda} + \frac{SSA}{2\lambda^2} + \frac{mn(\bar{y}-m)^2}{2\lambda^2} = 0 \rightarrow \hat{\lambda} = \frac{SSA}{m} \text{ when } m = \bar{y} \end{cases}$$

→ Not MLE's yet b/c

$$\sigma_a^2 = \frac{\lambda - \sigma^2}{n}, \quad \hat{\sigma}_a^2 = \frac{\lambda - \sigma^2}{n} = \frac{(1 - \frac{1}{m})MSA - MSE}{n}$$

↑ not necessarily MLE b/c $\hat{\sigma}_a^2$ needs to be non-negative (not guaranteed by RHS)

MLE

If $(1 - \frac{1}{m}) MSA \geq MSE$, then

$$\begin{cases} \hat{\mu} = \bar{y} \\ \hat{\sigma}^2 = SSE \\ \hat{\sigma}_n^2 = \frac{(1 - \frac{1}{m}) MSA - MSE}{n} \end{cases} \quad \text{are MLEs}$$

If $(1 - \frac{1}{m}) MSA < MSE$, then

$$\underbrace{\Downarrow}_{\hat{\sigma}_n^2 = 0}$$

$\hat{\sigma}_n^2 = 0$, lies on boundary for allowed space of $\hat{\sigma}_n^2$

$$\Rightarrow \mu_1 = \mu_2 = \mu_3 \dots$$

Rewrite model based on this condition

$$\begin{cases} \hat{\sigma}_n^2 = 0 \\ \hat{\sigma}^2 = \frac{\sum_i \sum_j (y_{ij} - \bar{y})^2}{mn} = \frac{SST}{mn} \\ \hat{\mu} = \bar{y} \end{cases}$$

Mean

$$E(\hat{\mu}) = \mu$$

$E\hat{\sigma}^2$ and $E\hat{\sigma}_n^2$ don't have closed forms.

$$E\hat{\sigma}^2 = E(MSE | \hat{\sigma}_n^2 > 0) P(\hat{\sigma}_n^2 > 0)$$

$$+ E\left(\frac{SST}{mn} | \hat{\sigma}_n^2 < 0\right) P(\hat{\sigma}_n^2 < 0) \quad \text{generally cannot be computed}$$

Variance

$$\text{var}(\hat{\mu}) = \text{var}(\bar{y}) = \frac{\hat{\sigma}_n^2 + \frac{\sigma^2}{n}}{m} = \frac{\sigma^2 + n\hat{\sigma}_n^2}{mn}$$

$\bar{y}_1 \quad \bar{y}_2 \quad \dots$

$$\begin{aligned} y_{ij} &= \mu + \alpha_i + \varepsilon_{ij} \\ \bar{y}_i &= \mu + \alpha_i + \tilde{\varepsilon}_i \end{aligned}$$

$\text{var}\left[\frac{\hat{\sigma}^2}{\hat{\sigma}_{\alpha^2}}\right] \rightarrow$ no closed form but can derive asymptotic variance b/c it's inverse of fisher information

$$\begin{aligned} & \left\{ \begin{array}{l} l\sigma^2 \\ l\lambda \end{array} \right. \\ & \left\{ \begin{array}{l} -E l\sigma^2 \\ -E l\sigma^2\lambda \\ -E l\lambda\lambda \\ +E l\lambda\sigma^2 \end{array} \right. \rightarrow \text{Fisher info} \end{aligned}$$

$$\text{var}\left(\frac{\hat{\sigma}^2}{\hat{\sigma}_{\alpha^2}}\right) \rightarrow (\text{Fisher Info})^{-1}$$

$$\Rightarrow \text{var}\left(\frac{\hat{\sigma}^2}{\hat{\sigma}_{\alpha^2}}\right)$$

) variable transformation

One-way classification

$$\text{Var}(\hat{\mu}) = \text{Var}(\bar{y}) = \frac{\sigma^2 + n\sigma^2_a}{mn}$$

$$\hat{\sigma}_\sigma^2, \hat{\sigma}_a^2$$

$$\begin{cases} \ell_{\sigma^2} = \frac{\partial \ell}{\partial \sigma^2} = \frac{-m(n-1)}{2\sigma^2} + \frac{SSE}{2\sigma^4} \\ \ell_\lambda = \frac{\partial \ell}{\partial \lambda} = \frac{-m}{2\lambda} + \frac{SSA}{2\lambda^2} + \frac{mn(\bar{y}-\mu)^2}{2\lambda^2} \end{cases}$$

$$\begin{cases} \ell_{\sigma^2 \sigma^2} = \frac{\partial^2 \ell}{\partial \sigma^2 \partial \sigma^2} = \frac{m(n-1)}{2\sigma^4} - \frac{2SSE}{2\sigma^6} \\ \ell_{\sigma^2 \lambda} = \frac{\partial^2 \ell}{\partial \sigma^2 \partial \lambda} = 0 \\ \ell_{\lambda \lambda} = \frac{\partial^2 \ell}{\partial \lambda \partial \lambda} = \frac{m}{2\lambda^2} - \frac{2SSA}{2\lambda^3} - \frac{mn(\bar{y}-\mu)^2}{2\lambda^5} \end{cases}$$

$$\begin{cases} -E \ell_{\sigma^2 \sigma^2} = \frac{m(n-1)}{2\sigma^4} \\ -E \ell_{\lambda \lambda} = \frac{m}{2\lambda^2} \end{cases}$$

inverse of Fisher information

$$\Rightarrow \text{Var} \begin{bmatrix} \hat{\sigma}^2 \\ \hat{\lambda} \end{bmatrix} \approx \begin{bmatrix} \frac{2\sigma^4}{m(n-1)} & 0 \\ 0 & \frac{2\lambda^2}{m} \end{bmatrix}$$

$$\hat{\sigma}_a^2 = \frac{\hat{\lambda} - \hat{\sigma}^2}{n}$$

$$\Rightarrow \text{var} \left[\frac{\hat{\sigma}^2}{\hat{\sigma}_n^2} \right] \underset{n \rightarrow \infty}{\cong} 2\sigma^4 \begin{bmatrix} \frac{1}{m(n-1)} & \text{cov}(-\frac{\sigma^2}{n}, \sigma^2) \\ \text{cov}(\frac{\sigma^2}{n}, \sigma^2) & \text{var}(\sigma_n^2) \end{bmatrix}$$

$$= 2\sigma^4 \begin{bmatrix} \frac{1}{m(n-1)} & \frac{-1}{mn(n-1)} \\ \frac{-1}{mn(n-1)} & \frac{1}{n^2} \left(\frac{\lambda^2/\sigma^4}{m} + \frac{1}{m(m-1)} \right) \end{bmatrix}$$

$$\left\{ \begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1m} \\ y_{21} \\ \vdots \\ y_{2n} \\ \vdots \\ y_{m1} \\ \vdots \\ y_{mn} \end{bmatrix} \right. = \left. \begin{bmatrix} 1 & \dots & m \\ 0 & 0 & \dots & 0 \\ \vdots & & & \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} + \begin{bmatrix} e_{11} \\ e_{12} \\ \vdots \\ e_{mn} \end{bmatrix} \right.$$

$x \quad \beta \quad + \quad \varepsilon$

fixed effects model.

$$u_i = \mu + \alpha_i \sim N(0, \sigma^2_\alpha)$$

$$\left\{ \begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{1m} \\ y_{21} \\ \vdots \\ y_{2n} \\ \vdots \\ y_{m1} \\ \vdots \\ y_{mn} \end{bmatrix} \right. = \left. \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} M + \begin{bmatrix} a_1 \\ \vdots \\ a_m \end{bmatrix} + \begin{bmatrix} e_{11} \\ e_{12} \\ \vdots \\ e_{mn} \end{bmatrix} \right.$$

e

$\epsilon \sim \Sigma$, non-diagonal variance matrix

$$\Sigma : \begin{bmatrix} \sigma^2 + \alpha^2 & \alpha^2 & & & \\ \alpha^2 & \ddots & & & \\ & & \ddots & & \\ & & & \sigma^2 + \alpha^2 & \\ \hline & & & & \vdots \\ 0 & & & & 0 \\ 0 & & & & 0 \\ 0 & & & & 0 \\ 0 & & & & 0 \end{bmatrix}$$

within group variances

- have correlation within the same group.
- trying to model a set of correlated data w/ random effects model
 - ↳ common usage of ranef models
- introduces a more parsimonious model for correlated data.

REML $y = X\beta + Zu + \varepsilon$

- first try to get rid of fixed effects.
- Balanced data: $n_1 = n_2 = \dots = n_m = n$.

$$L(\mu, \sigma^2, \alpha^2; y) \propto \exp\left(-\frac{1}{2} \left[\frac{\text{SSE}}{\sigma^2} + \frac{\text{SSA}}{\lambda} + \frac{(\bar{y} - \mu)^2}{\lambda/mn} \right] \right)$$

can treat as a normal density function
 ⇒ can integrate out

$$= \mathcal{L}(m, y) \mathcal{L}(\sigma^2, \sigma_a^2; SSA, SSE)$$

$$\frac{1}{(2\pi)^{1/2} \left(\frac{\lambda}{mn}\right)^{1/2}} \exp\left(\frac{-(\bar{y}-m)^2}{2\lambda/mn}\right)$$

$$\mathcal{L}(\sigma^2, \sigma_a^2; y) = \int \mathcal{L}(m, \sigma^2, \sigma_a^2; y) dm$$

SSA, SSE
sufficient statistics

REML log-likelihood:

$$\begin{aligned} l_R &= -\frac{1}{2} (mn-1) \log 2\pi - \frac{1}{2} \log mn \\ &\quad - \frac{1}{2} m(n-1) \log \sigma^2 - \frac{1}{2} (m-1) \log \lambda \\ &\quad - \frac{SSE}{2\sigma^2} - \frac{SSA}{2\lambda} \end{aligned}$$

taking derivative and solving:

$$\left\{ \begin{array}{l} \hat{\sigma}^2 = \frac{SSE}{m(n-1)} = MSE - \text{sum of MCE} \\ \hat{\lambda} = \frac{SSA}{n-1} = MSA - \text{mean between group sum of squares} \end{array} \right.$$

$$\text{MCE: } \frac{SSA}{m}$$

$$\hat{\sigma}_a^2 = \hat{\lambda} - \hat{\sigma}^2 = \frac{1}{n} (MSA - MSE)$$

Do not necessarily give values in allowed parameter space.

If $MSA \geq MSE$: $\hat{\sigma}_\alpha^2 \geq 0$

$$\begin{cases} \hat{\sigma}^2 = MSE \\ \hat{\sigma}_\alpha^2 = \frac{1}{n}(MSA - MSE) \end{cases}$$

If $MSA < MSE$: $\hat{\sigma}_\alpha^2 < 0$

$$\begin{cases} \hat{\sigma}^2 = \frac{SST}{mn-1} \\ \hat{\sigma}_\alpha^2 = 0 \end{cases} \Rightarrow \text{no within group variance}$$

Derivations are more complicated for unbalanced data.

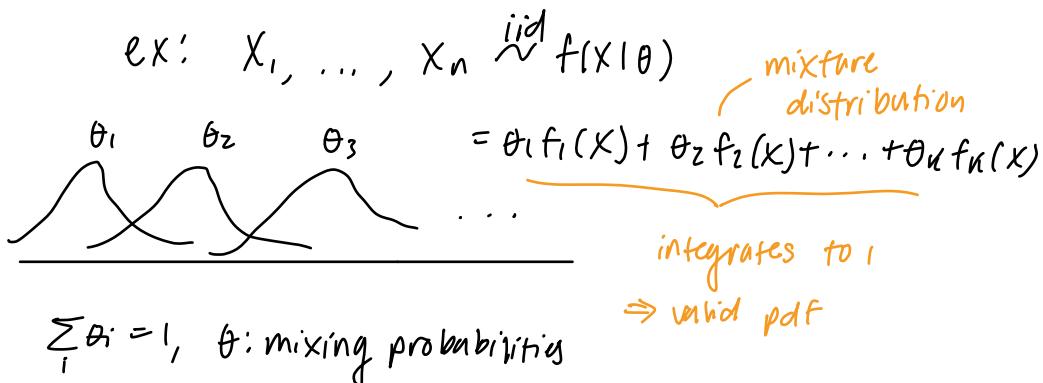
Expectation-Maximization (EM)

pdf $f(x, y | \theta) \rightarrow$ complete data likelihood
(\downarrow observed \downarrow missing) \rightarrow complete data

$f(x | \theta)$ can be derived by integrating y out.
 $= \int f(x, y | \theta) dy \rightarrow$ marginal prob. for observed data.

Goal: find MLE for observed data likelihood
(find θ to maximize $\int f(x, y | \theta) dy$)

- problem: not easy to work w/ observed data likelihood.
easier to work w/ complete data likelihood.



Assume f_i 's are pre-specified. Only unknowns are θ 's,
MLE approach: write down likelihood for obs. data

$$l(\theta) = \log L(\theta) = \sum_{j=1}^n \log [\theta_1 f_1(x_j) + \dots + \theta_n f_n(x_j)]$$

take deriv. + solve for each θ_i .

$$\left\{ \begin{array}{l} \frac{\partial}{\partial \theta_k} = \sum \frac{f_k(x)}{\theta_1 f_1(x) + \dots + \theta_n f_n(x)} = 0 \\ \text{subject to } \sum_i \theta_i = 1 \end{array} \right. \quad \text{difficult to solve.}$$

Instead, we can use the complete data likelihood.

introduce a set of indicators to indicate which of
the normal components was used to generate the data.

Step 1: randomly pick a category $1, \dots, K$ based on multinomial prob's $\theta_1, \dots, \theta_K$

Step 2: sample X from f_{θ_k} , your selected category.

introduce indicators to augment the system:

indicators: U_1, \dots, U_n

$$U_j = k \Rightarrow X \sim f_k(\cdot)$$

a priori.

$$\begin{cases} P(U_j = k) = \theta_k \\ f(X_j | U_j = k) = f_k(x_j) \end{cases}$$

two sets of RVs: U 's, X 's. only X 's are observed.
missing

$(U, X) \rightarrow$ complete data. can now describe
complete data likelihood.

$$f(X, U | \theta) = \prod_{j=1}^n \left\{ \prod_{i=1}^K [\theta_i f_i(x_j)] I(U_j = i) \right\}$$

multinomial likelihood

joint probability for U and X

viewing as a function of θ , it's the complete data likelihood.

- can count # points in each category if U_j 's are known \Rightarrow turns $f(u, x | \theta)$ into multinomial likelihood.
- But U 's (membership) are not known a priori.
- But if you know θ and x , you can infer U using Bayes Thm.

iterate between these two steps (EM algorithm)

EM algorithm

- Recall the goal: Find MLE for $L(\theta) = f(x | \theta)$

Define $Q(\theta | \theta_n) = E_{Y|X, \theta_n} \log f(x, y | \theta)$ using complete data likelihood

$($ y is missing, try to get its posterior then integrate out y from complete data likelihood.

$$= \int [\log f(x, y | \theta)] f(y | X, \theta_n) dy$$

can be computed from joint and marginal of x
(assume joint dist. is known.)

Start w/ initial estimate θ_0 .

\downarrow
E-step: compute $Q(\theta | \theta_n)$

\downarrow
M-step: try to find θ that maximizes $Q(\theta | \theta_n)$

$$\theta_{n+1} = \operatorname{arg\,max}_{\theta} Q(\theta | \theta_n)$$

iteratively run E-and M-step until convergence.

In this way, the observed data likelihood will not decrease

$$l(\theta_{n+1}) \geq l(\theta_n) \quad (\text{hill climbing algorithm})$$



algorithm will converge to a local maximum.

lec 4: April 6

Expectation Maximization

- observed data likelihood:

$f(x|\theta) \leftarrow$ not easy to maximize

introduce: $f(x, y|\theta)$ (augmented system)

complete data likelihood

$$Q(\theta | \theta^n) = E_{y|x} \log f(x, y|\theta)$$

↑
current parameter estimate

- Why does this algorithm work?

① Jensen's inequality

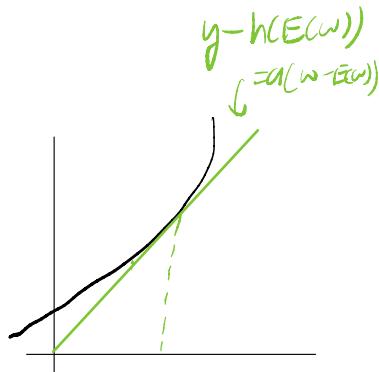
$w : RV$

$h(x) : \text{convex function}$

$$\text{Then, } E(h(w)) \geq h(E(w))$$

equation holds

iff $w = EW$ a.s.



$$h(w) \geq y = h(Ew) + a(w-Ew)$$

$$E(h(w)) \geq E(y) = h(Ew)$$

② Information Inequality

Consider: $f(x), g(x)$

$$E_f \log f(x) \geq E_f \log g(x)$$

$$E_f \log f(x) - E_f \log g(x)$$

$$= E_f \log \frac{f}{g} \quad \int \left[\log \frac{f(x)}{g(x)} \right] f(x) dx$$

$$= E_f \left[-\log \frac{g}{f} \right]$$

by Jensen's inequality,

$$E_f \left[-\log \frac{g}{f} \right] \geq -\log \left(E_f \left(\frac{g}{f} \right) \right)$$



$$= -\log \left(\int \frac{g(x)}{f(x)} f(x) dx \right) \quad \text{b/c } g \text{ is a density}$$

$$= -\log(1) = 0$$

so it integrates to 1

$$\Rightarrow E_f \log f(x) \geq E_f \log g(x)$$

③ The ascent property

$$\log f(x, y | \theta) = \log f(x | \theta) + \log f(y | x, \theta)$$

part we are trying to maximize

$$\log f(x | \theta) = \log f(x, y | \theta) - \log f(y | x, \theta)$$

taking expectation:

$$E_{Y|X,\theta^n} \log f(x|\theta) = \underbrace{E_{Y|X,\theta^n} \log f(x, Y|\theta)}_Q$$

doesn't contain y $-E_{Y|X,\theta^n} \log f(Y|X,\theta)$

$$\log f(x|\theta) = Q(\theta|\theta^n) - E_{Y|X,\theta^n} \log f(Y|X,\theta)$$

in EM, we compute Q and find θ^{n+1} that maximizes Q .

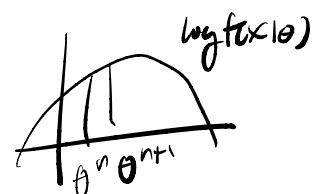
$$(1) Q(\theta^{n+1} | \theta^n) \geq Q(\theta^n | \theta^n)$$

$$(2) -E_{Y|X,\theta^n} \log f(Y|X, \theta^{n+1}) \geq -E_{Y|X,\theta^n} \log f(Y|X, \theta^n)$$

not the same distributions

old estimate: $\log f(x|\theta^n)$

new estimate: $\log f(x|\theta^{n+1})$



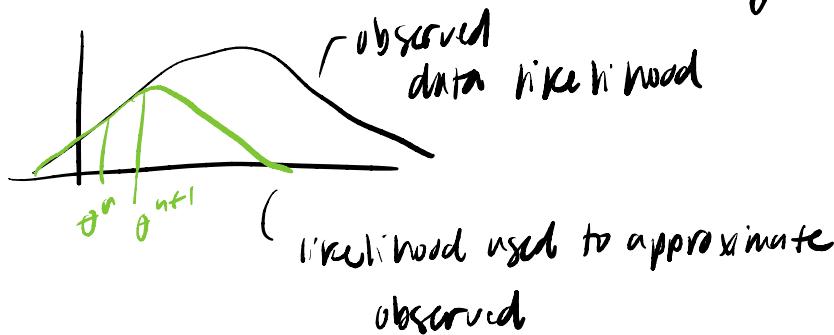
$\log f(x|\theta^n) \leq \log f(x|\theta^{n+1})$ by the above two inequalities

E-step: Evaluate $Q(\theta | \theta^n)$

M-step: Try to find θ to maximize Q

$$\theta^{n+1} = \operatorname{argmax}_{\theta} Q(\theta | \theta^n)$$

observed data
log lik



Example 1: Mixture model

x_1, \dots, x_n independently drawn from

$$f(x) = \theta_1 f_1(x) + \theta_2 f_2(x) + \dots + \theta_K f_K(x)$$

introduce indicators: u_1, \dots, u_n

view θ 's as multinomial probabilities

$$f(x, u | \theta) = \prod_{j=1}^n \prod_{i=1}^K [\theta_i f_i(x_j)] I(u_j=i)$$

$$\log f(x, u | \theta) = \sum_{j=1}^n \sum_{i=1}^K \left\{ \log \theta_i + \log f_i(x_j) \right\} I(u_j=i)$$

Assume given
constant

Goal: find θ that maximizes the loglik

E-step: $u|x, \theta^t$

$$E_{u|x, \theta^t} I(u_j = i)$$

$$= P(u_j = i | x, \theta^n)$$

$$= \frac{\theta_i^t f_i(x_j)}{\sum_i \theta_i^t f_i(x_j)} \triangleq \hat{u}_{j,i}$$

$$\text{then } E \log f(x, u|\theta) = \sum_{j=1}^n \sum_{i=1}^k \hat{u}_{j,i} \{ \log \theta_i + \log f_i(x_j) \}$$

$$Q(\theta|\theta^n)$$

psuedocounts

$$\text{M-step: } Q(\theta|\theta^t) = \sum_{j=1}^n \sum_{i=1}^k \hat{u}_{j,i} \log \theta_i + \text{const}$$

multinomial loglik

$$= \sum_{i=1}^k \left\{ \left[\sum_{j=1}^n \hat{u}_{j,i} \right] \log \theta_i \right\}$$

$$\hat{\theta}_i^{(t+1)} = \frac{\sum_{j=1}^n \tilde{u}_{ji}}{\sum_{i=1}^n \sum_{j=1}^n \tilde{u}_{ji}}$$

$$= \frac{\sum_{j=1}^n \tilde{u}_{ji}}{n}$$

Example 2: Mixed effects model

$$\begin{array}{c}
 \alpha_i \sim N(0, \sigma_\alpha^2) \\
 \text{---} \\
 M_i = \mu + \alpha_i \\
 \text{---} \\
 y_{ij} \sim N(M_i, \sigma^2)
 \end{array}$$

assume each group has n_i datapoints (not the same across groups)

$$\begin{array}{c}
 P(Y, M_i | \mu, \sigma^2, \sigma_\alpha^2) \\
 \text{---} \\
 \text{obs} \quad \text{missing}
 \end{array}$$

- $P(Y | \mu, \sigma^2, \sigma_a^2)$ is too hard to work with
- use complete data likelihood instead.

$$P(Y, M | \mu, \sigma^2, \sigma_a^2)$$

$$= \prod_{i=1}^m \left\{ \left[\prod_{j=1}^{n_i} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y_{ij} - \mu_i)^2}{2\sigma^2}\right\} \right] \right.$$

$$\left. \frac{1}{\sqrt{2\pi\sigma_a^2}} \exp\left\{-\frac{(\mu_i - \mu)^2}{2\sigma_a^2}\right\} \right\}$$

generate μ_i

$$= \prod_{i=1}^m (2\pi\sigma^2)^{-n_i/2} (2\pi\sigma_a^2)^{-1/2}$$

$$\exp\left\{-\frac{\sum (y_{ij} - \mu_i)^2}{2\sigma^2} - \frac{(\mu_i - \mu)^2}{2\sigma_a^2}\right\}$$

follows
normal dist.
conditional
on y and
other priors

Posterior dist. for μ_i : prior + data
means weighted by amount
of info from both.

$$\mu_i | Y, M, \sigma^2, \sigma_a^2$$

$$\sim N\left(\frac{\frac{\mu_i}{\sigma^2} \bar{y}_i + \frac{1}{\sigma_a^2} \mu}{\frac{n_i}{\sigma^2} + \frac{1}{\sigma_a^2}}, \frac{1}{\frac{n_i}{\sigma^2} + \frac{1}{\sigma_a^2}}\right)$$

derived
for HW

$$\begin{aligned}
 l_c &= \log P(Y, M_i | \mu, \sigma^2, \sigma_a^2) \\
 &= -\frac{1}{2}(N+m) \log 2\pi - \frac{N}{2} \log \sigma^2 - \frac{m}{2} \log \sigma_a^2 \\
 &\quad - \sum_{i=1}^m \sum_{j=1}^{n_i} \frac{(y_{ij} - m_i)^2}{2\sigma^2} - \sum_{i=1}^m \frac{(m_i - \mu)^2}{2\sigma_a^2} \\
 &= -\frac{1}{2}(N+m) \log 2\pi - \frac{N}{2} \log \sigma^2 - \frac{m}{2} \log \sigma_a^2 \\
 &\quad - \frac{1}{2\sigma^2} \sum_{i=1}^m \sum_{j=1}^{n_i} (y_{ij} - 2m_i y_{ij} + m_i^2) \\
 &\quad - \frac{1}{2\sigma_a^2} \sum_{i=1}^m (m_i^2 - 2m_i \mu + \mu^2)
 \end{aligned}$$

terms involving missing data

E-step: compute $Q(\theta | \theta^t) = E_{M|Y} l_c$

m_i 's are the missing data

Need to evaluate $E(m_i | Y, \mu, \sigma^2, \sigma_a^2)$

$E(m_i^2 | Y, \mu, \sigma^2, \sigma_a^2)$

We already know the distribution of m_i ,
so these expectations are easy to evaluate.

$$E(m_i | y, m^t, \sigma^2, \alpha_a^2) = \frac{\frac{m_i}{\sigma^2} \bar{y}_i + \frac{1}{\sigma_a^2} \alpha_a^t}{\frac{m_i}{\sigma^2} + \frac{1}{\sigma_a^2}} \triangleq \tilde{m}_i$$

$$E(m_i^2 | y, m^t, \sigma^2, \alpha_a^2) = \frac{1}{\frac{m_i}{\sigma^2} + \frac{1}{\sigma_a^2}} + \tilde{m}_i^2 \triangleq \tilde{n}_i$$

in LC, replace m w/ \tilde{m} , m^2 w/ \tilde{n} to get
Q function.

M-step:

$$m^{t+1} = \frac{1}{m} \sum_{i=1}^m \tilde{m}_i \quad (\text{maximized by taking deriv. and solving})$$

$$\sigma^2^{t+1} = \frac{1}{N} \sum_{i=1}^m \sum_{j=1}^n (y_{ij} - 2y_j \tilde{m}_i + \tilde{n}_i)$$

similar to sum of squared residuals.

$$\alpha_a^{2,t+1} = \frac{1}{m} \left\{ \sum_{i=1}^m \tilde{n}_i - \left(\frac{\sum \tilde{m}_i}{m} \right)^2 \right\}$$

$$= \frac{1}{m} \sum_{i=1}^m \tilde{\eta}_i - [m^{t+1}]^2$$

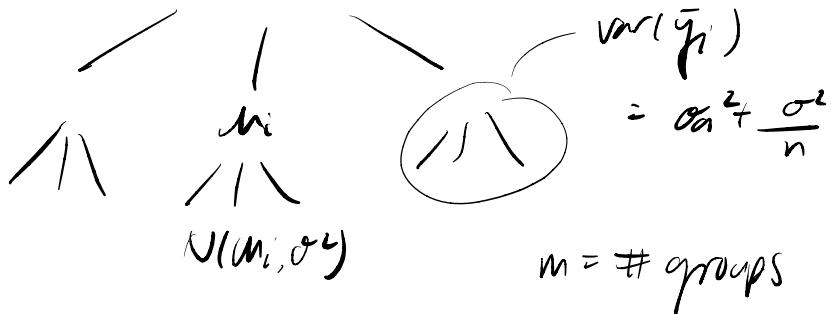
Lee 5: April 11

Random effects model:

$$\mu, \sigma^2, \sigma_a^2$$

$\sigma_a^2 = 0 \Rightarrow$ all means are equal

$$N(\mu, \sigma^2)$$



Balanced data

$$\hat{\mu}, \bar{y} \sim N\left(\mu, \frac{\sigma_a^2 + \frac{\sigma^2}{n}}{m}\right) \quad S^2 \sim \frac{\sigma^2}{N-m} \chi^2_{N-m}$$

$$\hat{\sigma}^2: \text{MSE} = \frac{\sum \sum (y_{ij} - \bar{y}_{..})^2}{N-m} \triangleq S^2$$

$$\hat{\sigma}_a^2: H_0: \sigma_a^2 = 0$$

$$\lambda = \sigma^2 / n \sigma_a^2$$

$$\Rightarrow \lambda = \sigma^2$$

used to estimate
 σ_a^2

estimated using MSA

(between group variance divided by corresponding df)

$$\frac{\text{MSA}}{\text{MSE}} = \frac{\text{SSA}/(m-1)}{\text{SSE}/(N-m)} \sim F_{m-1, N-m}$$

used to evaluate whether
 $\sigma^2 = 0$

unbalanced data:

- take deriv. of obs. data likelihood, can also take 2nd deriv.

$$\frac{\partial \log f(x|\theta)}{\partial \theta}$$

plug in $\hat{\theta}$ and evaluate (easy)

$$\frac{\partial^2 \log f(x|\theta)}{\partial \theta^2}$$

\Rightarrow can evaluate Fisher information \rightarrow take inverse to get asymptotic variance.

- Since these are MLEs, we can use asymptotic properties of MLEs to do inference.

- In mixed effect models, sometimes we are interested in estimating μ_i 's (group specific means)

Predicting random effects:

X, Y are RVs

Y is unobserved, X is observed.

find $g(x)$ that can predict Y .

The best predictor is the one with minimum squared error.

$$E(Y - g(x))^2 \quad \text{if } g(x) \text{ should minimize this.}$$

$$g(x) = E(Y | X) \text{ minimizes this.}$$

so best predictor of μ_i 's is

$$E(\mu_i | y, n, \sigma^2, \sigma_a^2) \quad \underbrace{\text{last time: showed}}_{\mu_i \text{ has a normal dist.}}$$

$$= \frac{n_i}{\sigma^2} \quad \text{estimated using EM}$$

$$= \frac{\frac{n_i}{\sigma^2} + \frac{1}{\sigma_a^2}}{\frac{n_i}{\sigma^2} + \frac{1}{\sigma_a^2}} \bar{y}_i + \frac{\frac{1}{\sigma_a^2}}{\frac{n_i}{\sigma^2} + \frac{1}{\sigma_a^2}} \mu$$

$$\underbrace{\quad}_{\text{shrinkage factor}}$$

combination of sample group mean and prior mean.

• plug in EM MLE's of m , σ^2 , σ_{α}^2

$$m_i = m + a_i \sim N(0, \sigma_{\alpha}^2)$$

$$E(a_i | y, m, \sigma^2, \sigma_{\alpha}^2)$$

b/c $a_i = m_i - m$

$$= \frac{m}{\sigma^2} \left(\bar{y}_i - m \right) + \frac{\frac{1}{\sigma_{\alpha}^2}}{\frac{m}{\sigma^2} + \frac{1}{\sigma_{\alpha}^2}} 0$$

Note: we are trying to pull the sample mean (\bar{y}_i) towards its prior (m) in $E(a_i | y, \sigma^2, \sigma_{\alpha}^2, m)$

And trying to pull a_i estimate to 0.

Both uses: trying to shrink the data-driven estimate towards the prior.

Posterior mean also called **shrinkage estimator**

let $\beta = \text{shrinkage factor}$

$$= \frac{\frac{1}{\sigma_{\alpha}^2}}{\frac{m}{\sigma^2} + \frac{1}{\sigma_{\alpha}^2}}$$

σ_{α}^2 large \rightarrow info contributed by prior is small, rely on data more.

$$E(m|y, m, \sigma^2, \alpha^2) = (1-\beta)y_i + \beta m$$

α^2 small \rightarrow more info contributed by prior.

- groups very different \rightarrow mostly use data
- groups similar \rightarrow mostly use prior

Empirical Bayes

- used data from all groups to estimate overall mean
- for group specific mean: combined overall mean w/ group-specific empirical mean (borrowed info from other groups)

Genomics example:

20,000 genes

	Patient			Control			interested in DE genes
	1	2	3	1	2	3	
1	5.1	4.8	5.3	1.6	2.1	1.7	
2							
3							
4	0.9	1.0	1.1	1.5	1.4	1.5	very small sample variance
5							
6							
7							
8							
9							
10							

- can run a t-test for each row.

$H_0:$ patient mean = control mean
 gene expr gene expr

- lots of false positives (variance can't be stably estimated w/ small sample size)

$$\frac{\bar{x} - \bar{y}}{\sqrt{\left(\frac{1}{n_1} + \frac{1}{n_2}\right)S^2}}$$

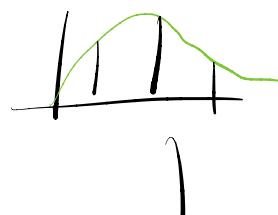
assume the
two groups have same variance

$$S^2 = \frac{\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2}{n_1 + n_2 - 2} \quad \left. \begin{array}{l} \text{pooled} \\ \text{sample variance} \end{array} \right\}$$

Suppose data is normally distributed:

$$X \sim N(\mu_X, \sigma^2), \quad Y \sim N(\mu_Y, \sigma^2)$$

$$\frac{df S^2}{\sigma^2} \sim \chi_{df}^2$$



↙

probability of getting a very small sample variance is low, but we are considering $\sim 20,000$ genes.

Scenario: No DE gene.

- But, mean difference is non-negligible, and s^2 is small. A few hundred out of 20,000 will have this \rightarrow huge t -statistic (false positive)
- very low statistical power to detect true DE genes.



using shrinkage to solve this problem:

$$x_{ijk} | \mu_{ij}, \sigma_i^2 \sim N(\mu_{ij}, \sigma_i^2)$$

group 1: patient, 2: control

gene sample replicate
1, 2, 3

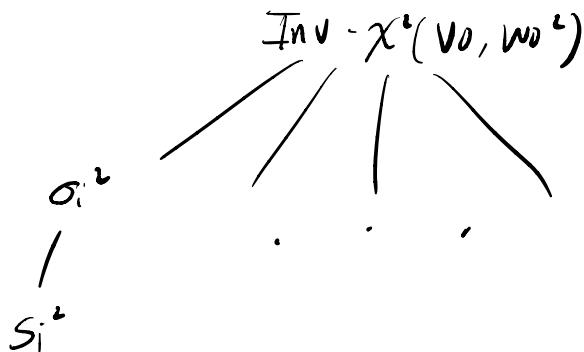
$$S_i^2 = \frac{\sum (x_{ij} - \bar{x}_{ij})^2}{n_1 + n_2 - 2}$$

$$S_i^2 | \sigma_i^2 \sim \frac{\sigma_i^2 \chi_v^2}{v} \text{ under normality assumption}$$

$v = n_1 + n_2 - 2$

If we assume σ_i^2 are random effects:

$$\sigma_i^2 \sim \text{Inv-}\chi^2(v_0, w_0^2)$$



$$X \sim \chi^2_v, \quad Y = \frac{v_0 w_0^2}{X} \sim \text{Inv-}\chi^2(v_0, w_0^2)$$

conjugate prior for χ^2

PDF of Inv- $\chi^2(v, w^2)$: \Rightarrow posterior will be

$$f(\theta) = \frac{(\frac{v}{2})^{\frac{v}{2}}}{\Gamma(\frac{v}{2})} (\omega^2)^{\frac{v}{2}} \theta^{-(\frac{v}{2}+1)} e^{-\frac{v w^2}{2\theta}}$$

$$E(\theta) = \frac{V}{V-2} w^2$$

$$\text{var}(\theta) = \frac{2V^2}{(V-2)^2(V-4)} w^4$$

$$\text{mod}(\theta) = \frac{V}{V+2} w^2$$

$$E s_i^2 = E \sigma_i^2 = \frac{V_0}{V_0 - 2} w_0^2$$

$$\hat{\sigma}_i^2 = E(\sigma_i^2 | s_i^2, V_0, w_0^2)$$

[data-based estimate]
[prior mean for σ_i^2]

$$= (1 - B) s_i^2 + B \left(\frac{V_0}{V_0 - 2} w_0^2 \right)$$

Shrinkage estimator

$$B = \frac{E(\text{var}(s_i | \sigma_i^2))}{\text{var}(s_i^2)}$$

integrate missing data out
to get marginal dist of s_i^2

} ratio between
within group var
and between
group var.

- borrowing information from other groups
- estimate group i's true variance

In the normal model:

$$\beta_3 = \frac{\frac{1}{\sigma_a^2}}{\frac{n_i}{\sigma_i^2} + \frac{1}{\sigma_a^2}} \cdot \frac{\frac{\sigma_a^2 \sigma^2}{n_i}}{\frac{\sigma_a^2 \cdot \sigma^2}{n_i}}$$

$$= \frac{\frac{\sigma^2}{n_i}}{\underbrace{\frac{\sigma^2}{n_i} + \sigma_a^2}_{\text{total variance}}} \quad \begin{array}{l} \text{within group} \\ \text{variance} \end{array}$$

Lee 6: April 13

HW 2: due April 27
Proj. S: random + mixed effects
districts of Baltimore

Shrinkage:

(1-B) Data + B Prior

$$\hat{\beta} = \frac{\text{within-group variance}}{\text{total variance}}$$

$$= \frac{E[\text{var}(s_i^2 | \sigma_i^2)]}{\text{var } s_i^2}$$

$$= \frac{v_0 - 2}{v_0 + v_i - 2}$$

prior has higher df
(\Rightarrow contains more info, give it more weight)

$$s_i^2 | \sigma_i^2 \sim \frac{\sigma_i^2 \chi_{v_i}^2}{v_i}$$

df contributed by the prior

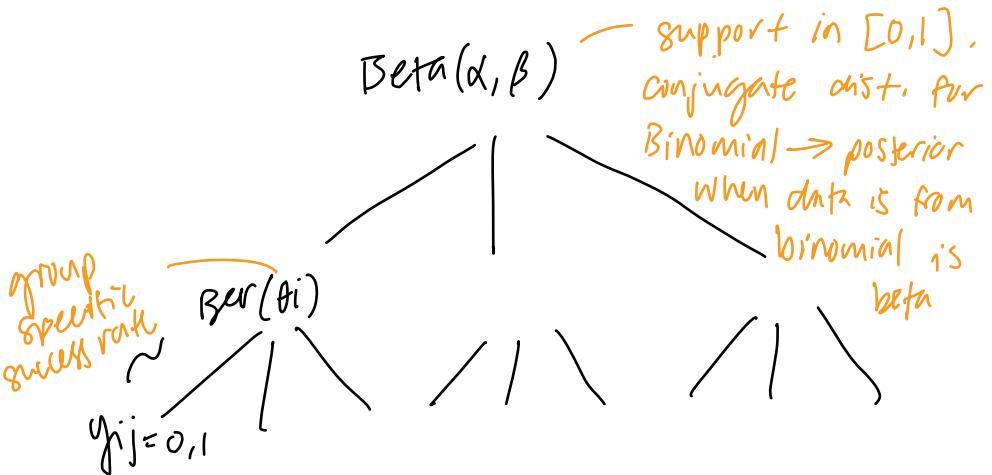
df contributed by gene i

$$\sigma_i^2 \sim \frac{v_0 w_0^2}{\chi_{v_0}^2}$$

- small within-group variance \rightarrow more difference between groups, β should be close to 0.

How to handle non-normal models?

- Fixed effects model: compute π from data.
- Mixed effects: π are random



$P(y_i, \theta | \alpha, \beta)$

m groups.
indep.

$$= \prod_{i=1}^m \left[\frac{\theta_i^{\alpha-1} (1-\theta_i)^{\beta-1}}{\text{Beta}(\alpha, \beta)} \prod_{j=1}^{n_i} \left[\theta_i^{y_{ij}} (1-\theta_i)^{1-y_{ij}} \right] \right]$$

group
 θ

$\propto \prod_{i=1}^m \theta_i^{\alpha+y_{it}-1} (1-\theta_i)^{\beta+n_i-y_{it}-1}$

sum of all y_{ij} 's in
the i th group
(total # success)
total # failures

joint dist. of θ, y .

has form of Beta dist.

if we treat y as known,
 $P(\theta | y, \alpha, \beta) \sim \text{Beta}(\alpha + y_{it}, \beta + n - y_{it})$

best predictor of θ_i
is posterior mean

$$E(y_{ij}) = E[E(y_{ij} | \theta_i)]$$

$$= E\theta_i$$

$$= \frac{\alpha}{\alpha + \beta}$$

View α, β as pseudocounts
 α : prior count of first category
 β : prior count of second category

$$\text{var}(y_{ij}) = \text{var}[E(y_{ij} | \theta_i)] + E[\text{var}(y_{ij} | \theta_i)]$$

$$= \frac{\alpha \beta}{(\alpha + \beta)^2} = E(y_{ij})[1 - E(y_{ij})]$$

$\text{cov}(y_{ij}, y_{ik}) = 0$ if $i \neq k$ b/c their θ 's are independently generated.

$$\text{cov}(y_{ij}, y_{ik}) = \text{cov}[E(y_{ij} | \theta_i), E(y_{ik} | \theta_i)]$$

$$+ E[\text{cov}(y_{ij}, y_{ik} | \theta_i)]$$

independently generated given θ_i

$$= \text{cov}(\theta_i, \theta_i) + E(0)$$

$$= \text{var}(\theta_i) = \frac{\alpha \beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}$$

How to estimate α , β , and predict θ_i ?

$$P(y | \alpha, \beta) = \int P(y, \theta | \alpha, \beta) d\theta$$
$$= \prod_{i=1}^m \frac{\text{Beta}(\alpha + y_{it}, \beta + n_i - y_{it})}{\text{Beta}(\alpha, \beta)}$$

$$\log P(y | \alpha, \beta) = \sum_{i=1}^m \log \left(\frac{\text{Beta}(\alpha + y_{it}, \beta + n_i - y_{it})}{\text{Beta}(\alpha, \beta)} \right)$$

(Beta functions can be expressed as gamma functions, log-gamma can be computed in R.)

posterior mean is
the best predictor
for θ_i

$$E(\theta_i | y, \alpha, \beta) = \frac{\alpha + y_{it}}{\alpha + y_{it} + \beta + n_i - y_{it}}$$

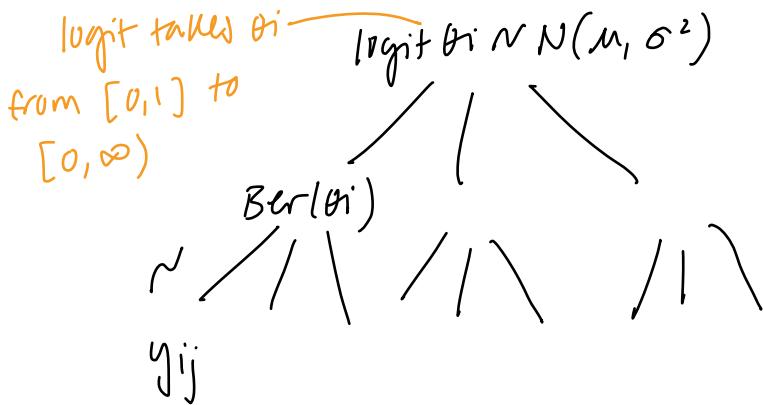
$$= \frac{\alpha + y_{it}}{\alpha + \beta + n_i}$$

$$= \frac{\alpha + \beta}{\alpha + \beta + n_i} \left(\frac{\alpha}{\alpha + \beta} \right) + \frac{n_i}{\alpha + \beta + n_i} \left(\frac{y_{it}}{n_i} \right)$$

reduces MSE
b/c variance is
reduced.

only use this term
in the fixed effects
model

If we use a different prior:



$$P(y_i | \theta | \mu, \sigma^2)$$

$$= \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\text{logit}\theta_i - \mu)^2}{2\sigma^2}} \prod_{j=1}^{n_i} \theta_i^{y_{ij}} (1-\theta_i)^{1-y_{ij}}$$

generate θ_i from
Normal

- $P(y_i | \theta | \mu, \sigma^2)$ is hard to integrate wrt θ
- EM also does not help b/c can't compute expectation
- use Markov chain Monte Carlo (MCMC)

Markov Chains

1. A MC is a stochastic process $\{X_i\}$, such that

$$P(X_i | X_0, \dots, X_{i-1}) = P(X_i | X_{i-1})$$

the current state of the chain depends only on
the immediate previous state

State space : the collection of values of X_i

* $i = 0, 1, 2, \dots$ discrete time

t - continuous time

{
· finite state space
· countable state space
* }

2. One-step transition probability:

$$P_{ij}^{\text{time}} = \Pr\{X_{n+1} = j | X_n = i\}$$

If $P_{ij}^{\text{time}} = P_{ij}$ (doesn't depend on time), it is a

stationary transition probability.

$$\Pr(X_{n+i} = j \mid X_n = i) = \Pr(X_i = j \mid X_0 = i)$$

$$\sum_j p_{ij} = 1$$

3. n-step transition probability

$$P_{ij}^{(n)} = \Pr\{X_n = j \mid X_0 = i\}$$

$$P^{(n)} = \begin{bmatrix} & \cdots & & s \\ 1 & & & \\ \vdots & & P_{ij}^{(n)} & \\ s & & & \end{bmatrix}_{s \times s}$$

If the MC has a stationary one-step transition probability, then

$$P^{(n)} = P^n, \quad P = \left[\begin{array}{c} p_{ij} \end{array} \right]_{s \times s}$$

4. Right Multiplication

$$E[f(X_n) \mid X_{n-1} = i]$$

$$= \sum_j p_{ij} f(j)$$

5. Left multiplication

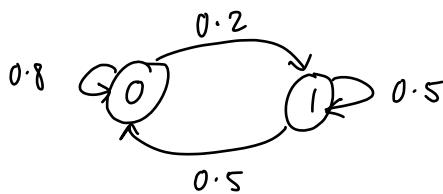
$$\pi = (\pi_1, \dots, \pi_s), \quad \sum_{i=1}^s \pi_i = 1$$

$$P(X_i=j | X_0=i) = p_{ij}$$

$$P(X_i=j) = \sum_i P(X_i=j | X_0=i) P(X_0=i)$$

$$= \sum_i \pi_i p_{ij}$$

Markov chains



$$P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

} one-step transition probability matrix

Assuming stationary transition probabilities:

$$P^{(n)} = P^n$$

Suppose $\Pr(X_0 = i) = \pi_i$, then

$$\Pr(X_1 = j | X_0 = i) = p_{ij}$$

$$\Pr(X_1 = j) = \sum_i \pi_i p_{ij}$$

} marginalize over X_0

If $\Pr(X_1 = j) = \pi_j$, $\pi P = \pi$ then π is called the stationary distribution of the Markov chain.

NOT the same as

stationary one-step transition probability

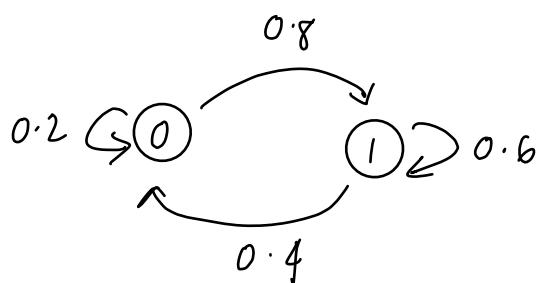
Stationary distribution:

$$\pi = \pi P = \pi P \cdot P = \pi P^n$$

So $\Pr(X_n = j) = \pi_j$

Limiting distribution:

$$P = \begin{bmatrix} 0.2 & 0.8 \\ 0.4 & 0.6 \end{bmatrix}$$



$$P^2 = \begin{bmatrix} 0.36 & 0.64 \\ 0.32 & 0.68 \end{bmatrix}$$

$$P^4 = \begin{bmatrix} 0.33 & 0.67 \\ 0.33 & 0.67 \end{bmatrix} \quad \text{each row becomes the same}$$

Suppose you start your chain with state 0, $X_0 = 0$

$$[1 \ 0] P^n = [0.33 \ 0.67]$$

probability that
 $X_n = 0$

probability that
 $X_n = 1$

Starting from state 1 instead:

$$\begin{bmatrix} 0 & 1 \end{bmatrix} P^n = \begin{bmatrix} 0.33 & 0.67 \end{bmatrix}$$

regardless of where you start your chain,
the probability of visiting the two states
at step n will be the same.

For a finite state space MC:

- P is regular if $\exists k$ such that P^k has all its elements > 0 .
- A MC with regular P has limiting distribution $\pi = (\pi_1, \dots, \pi_n)$ that satisfy

$$\lim_{n \rightarrow \infty} P_{ij}^{(n)} = \pi_j \leftarrow \text{does not depend on initial state } i$$

$$\sum_j \pi_j = 1, \quad \pi_j > 0$$

$$\pi \lim_{n \rightarrow \infty} P^{(n)} = \pi$$

$$\lim_{n \rightarrow \infty} \pi P^{n+1} = \pi P^n$$

$$= (\pi P) P^n = \pi P^n \Rightarrow \pi P = \pi$$

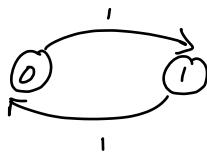
\curvearrowright

solve this equation to get
the limiting distribution
(sol'n is unique)

- Stationary distribution \neq limiting distribution
- limiting \Rightarrow stationary, stationary $\not\Rightarrow$ limiting

• ex:

$$P = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$



$$\pi = [0.5 \ 0.5]$$

$$\pi P = [0.5 \ 0.5] = \pi \text{ (stationary)}$$

$$X_0 = 0 \quad P_{01}^{(1)} = 1, \quad P_{01}^{(2)} = 0 \quad P_{01}^{(3)} = 1$$

$$P_{01}^{(n)} = 1, 0, 1, 0, 1, 0 \dots$$

$$P_{00}^{(n)} = 0, 1, 0, 1, 0, 1 \dots$$

$$P_{10}^{(n)} = 1, 0, 1, 0, 1, 0 \dots$$

$$P_{11}^{(n)} = 0, 1, 0, 1, 0, 1 \dots$$

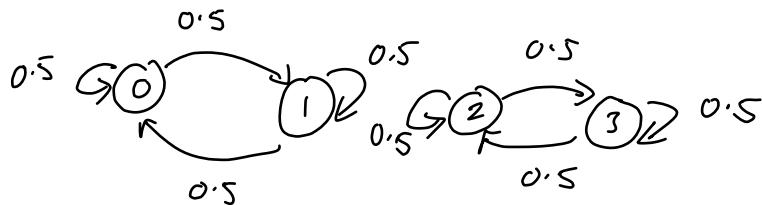
don't have limits
as $n \rightarrow \infty$

\Rightarrow no limiting distribution

- can't find a n such that P^n has all entries > 0

$$P^2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad P^3 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \dots$$

ex:



$$P = \begin{bmatrix} 0 & 1 & 2 & 3 \\ 0.5 & 0.5 & 0 & 0 \\ 0.5 & 0.5 & 0 & 0 \\ 0 & 0 & 0.5 & 0.5 \\ 0 & 0 & 0.5 & 0.5 \end{bmatrix}$$

- this MC has a stationary distribution.

$$\pi = [0.25 \ 0.25 \ 0.25 \ 0.25]$$

$$\pi P = \pi$$

$$X_0 = 0, \quad P_{00}^{(1)} = 0.5 \quad P_{00}^{(2)} = 0.5$$

$X_0 = 2 \quad P_{20}^{(1)} = 0 \quad P_{20}^{(2)} = 0$

have different limits
(depends on initial state i)

thus, we don't have a limiting distribution.

- If you have a stationary distribution, your limiting distribution will be unique.
- Reasons for not having a limit:
 - periodic behavior (first example)
 - two different groups of states (second example)

Classification of states of a MC

① Accessible: $i \rightarrow j$ if $\exists k$ such that $P_{ij}^{(n)} > 0$

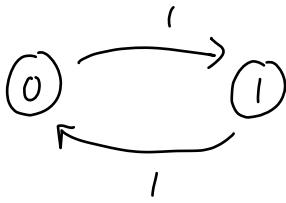
initial state is i , after n -step transitions, there is non-zero probability of transitioning to j .

② Communicate: $i \leftrightarrow j \Leftrightarrow \begin{matrix} i \rightarrow j \\ j \rightarrow i \end{matrix}$

③ Irreducible: $\nexists i, j, i \leftrightarrow j$ or $P_{ij}^{(n)} > 0$ for some n .

④ Period of a state i , $d(i)$:

the greatest common divisor (GCD) of all integers $n \geq 1$ for which $P_{ii}^{(n)} > 0$



$$d(0) = ?$$

$$P_{00}^{(1)} = 0$$

$$P_{00}^{(2)} = 1 \rightarrow n = 2$$

$$P_{00}^{(3)} = 0$$

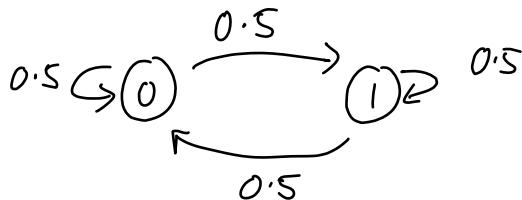
$$P_{00}^{(4)} = 1 \rightarrow n = 4$$

⋮

sequence of n 's : 2, 4, 6, 8, ...

$$G(0) = 2 = d(0)$$

$d(1) = 2$ as well.



$$P = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

$$P_{00}^{(1)} = 0.5$$

$$P_{00}^{(2)} > 0$$

> 0

1

2

3 ..

$d(0) = 1 \Rightarrow 0$ is an aperiodic state

⑤ Aperiodic: period = 1

- If $i \leftrightarrow j$, then $d(i) = d(j)$
- In an irreducible MC, all states have the same period.
 - Period is a class property for irreducible MCs.

⑥ $f_{ii}^{(n)} = \Pr \{ X_n = i, X_v \neq i \text{ for } v=1, \dots, n-1 \mid X_0 = i \}$

The probability that starting from i , the first return to i occurs at the n^{th} transition.

$$f_{ii}^0 = 0, f_{ii}^1 = p_{ii},$$

$$P_{ii}^{(n)} = \sum_{k=0}^n f_{ii}^k p_{ii}^{(n-k)}$$

$$t=0 \quad 1 \quad 2 \quad n=3$$

state: ② 0 0 0

start ① 0 0 0 ← probability you're
still in state 1 after n steps? $p_{i,i}^n$

② 0 0 0

③ 0 0 0

$$\{X_0 = i, X_n = i\} = \underbrace{\{1^{\text{st}} \text{ return}\}_{n=1}}_{\text{disjoint events}} \cup \{2\} \cup \{3\} \dots$$

Adding them up shows

$$P_{ii}^{(n)} = \sum_{k=0}^n f_{ii}^{(k)} P_{ii}^{(n-k)}$$

⑦ recurrent:

$$\text{State } i \text{ is recurrent} \iff \sum_{n=1}^{\infty} f_{ii}^{(n)} = 1$$

i.e. the probability to return to state i after some finite length of time is equal to 1.

• Class property

⑧ transient: a non-recurrent state

$$\text{i.e. } \sum_{n=1}^{\infty} f_{ii}^{(n)} < 1$$

i is recurrent if and only if $\sum_{n=1}^{\infty} P_{ii}^{(n)} = \infty$

$$\sum_{n=1}^{\infty} E[I(X_n=i) \mid X_0=i] = \sum_{n=1}^{\infty} P_{ii}^{(n)} = \infty$$



count # times
you've visited state i

⑨ Positive recurrent (strongly ergodic):

$$\sum_{n=0}^{\infty} n f_{ii}^n < \infty$$

Null recurrent (weakly ergodic):

$$\sum_{n=0}^{\infty} n f_{ii}^n = \infty$$

expected length of
time for first occurrence

class properties $i \leftrightarrow j$, i positive recurrent $\Rightarrow j$ positive recurrent.

Markov chains

- An irreducible MC with finite state space is positive recurrent.

Basic Limit Theorem of MCs:

- ① A recurrent, irreducible, aperiodic MC:

$$\lim_{n \rightarrow \infty} P_{ii}^{(n)} = \frac{1}{\sum_{n=0}^{\infty} n f_{ii}^{(n)}} \quad \text{probability that first return occurs at time } n$$

$$\stackrel{\Delta}{=} \frac{1}{m_i} = 0 \quad \text{if MC is null-recurrent}$$

As well,

$$\lim_{n \rightarrow \infty} P_{ij}^{(n)} = \lim_{n \rightarrow \infty} P_{ii}^{(n)} \quad \begin{matrix} \text{converge to a number that} \\ \text{doesn't depend on } i \end{matrix}$$

- ② For positive recurrent, irreducible, aperiodic MCs:

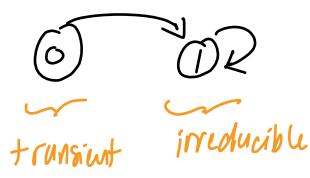
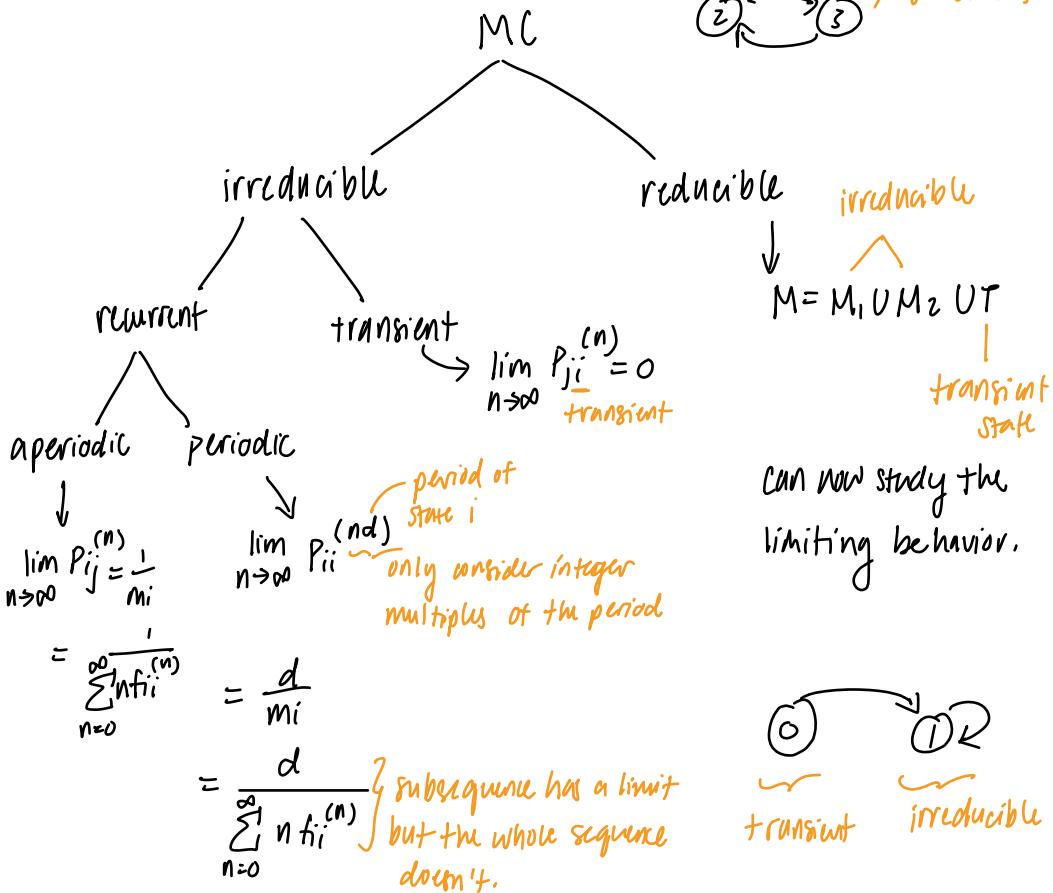
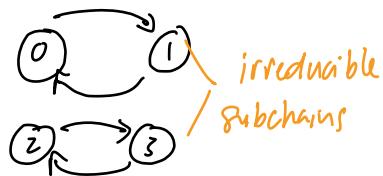
$$\lim_{n \rightarrow \infty} P_{jj}^{(n)} = \pi_j = \sum_{i=0}^{\infty} \pi_i P_{ij}$$

$$\sum_{i=0}^{\infty} \pi_i = 1$$

i.e. limiting distribution exists and it is equal to the stationary distribution. π is unique and (limiting distribution)

uniquely determined by:

$$\begin{cases} \pi = \pi P \\ \sum \pi_i = 1 \end{cases}$$



$$\lim_{n \rightarrow \infty} a_n = c \Rightarrow \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n = c$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n P_{ji}^{(n)} = \lim_{n \rightarrow \infty} P_{ji}^{(n)} = \frac{1}{m_i}$$

$$\begin{aligned} \frac{1}{n} \sum_{k=1}^n p_{ji}^{(k)} &= \frac{1}{n} \sum_{k=1}^n E \left[\underbrace{\mathbb{I}(X_n=i) | X_0=j}_{\text{count # times state } i \text{ has been visited from first transition to } n^{\text{th}}} \right] \\ &= \frac{1}{n} E \left[\underbrace{\sum_{k=1}^n \mathbb{I}(X_n=i) | X_0=j}_{\text{empirical frequency of visiting state } i \text{ has same limit as transition probability}} \right] \end{aligned}$$

empirical frequency of visiting state i has same limit as transition probability

longitudinal behavior of the MC

- empirical longitudinal behavior can be used to study the cross-sectional behavior because they converge to the same limiting distribution.

Reversible Markov Chains

- A reversible MC is a MC such that:

$$\exists \pi, \pi_i p_{ij} = \pi_j p_{ji} \quad (*)$$

detailed balance condition

$$\pi_i = \sum_j \pi_i p_{ij} = \sum_j \pi_j p_{ji}$$

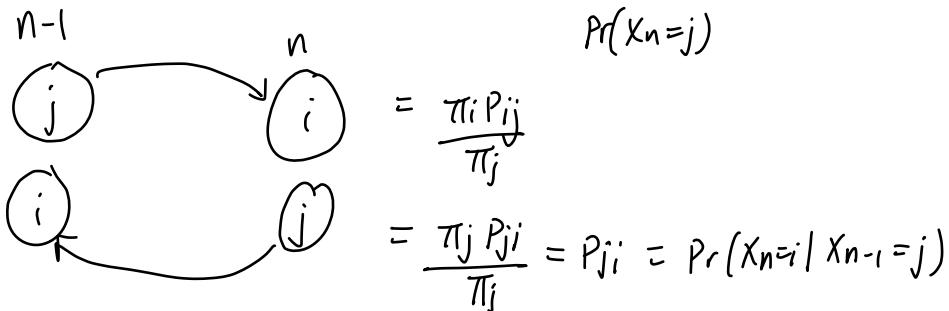
$$\Rightarrow \pi = \pi P$$

• Why "reversible"?

$$\Pr(X_{n-1} = i \mid X_n = j) = \frac{\Pr(X_{n-1} = i, X_n = j)}{\Pr(X_n = j)}$$

reverse + transition probability

$$= \frac{\Pr(X_{n-1} = i) \Pr(X_n = j \mid X_{n-1} = i)}{\Pr(X_n = j)}$$



Markov Chain Monte Carlo

Metropolis - Hastings

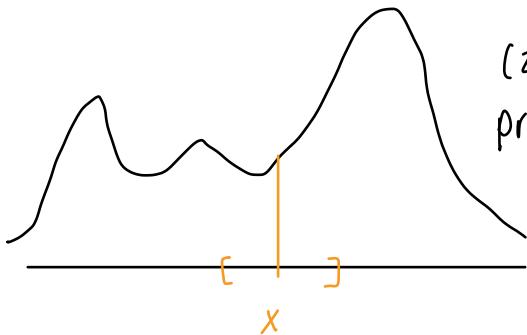
Goal: Draw samples from a target distribution $\pi(\cdot)$

initial value: X_0

$$X_{n-1} = x$$

(1) simulate a candidate from the jumping distribution

$$y \sim q(x, \cdot)$$



(2) Accept the candidate with probability

$$\alpha(x, y) = \min \left[\frac{\pi(y) q(y, x)}{\pi(x) q(x, y)}, 1 \right]$$

conditional that

current state is x ,
what's the probability
of transitioning to y

- In other words, draw u from $\text{unif}(0, 1)$ and accept y if

$u \leq \alpha(x, y)$. If y is accepted, $x_n = y$. Otherwise,

$$X_n = X_{n-1} = x.$$

i.e. $X_n = \begin{cases} y & \text{if } u \leq \alpha(x, y) \\ x & \text{if } u > \alpha(x, y) \end{cases}$

- Sequence of X 's is a Markov chain.

Why does MH work?

(1) creates a MC

(2) chain is reversible w/ $\pi(\cdot)$.

Pf:/ Let $\kappa(x, y)$ denote the transition density of the MC.

$$\kappa(x, y) = \alpha(x, y) q(x, y) \quad \text{if } x \neq y$$

$$+ \delta_x(y) \left[\sum_t [1 - \kappa(x, t)] q(x, t) \right] \quad \text{if } x = y$$

indicator for $x = y$

propose a new t , reject
it and keep x

$$\text{WTS: } \pi(x) \kappa(x, y) = \pi(y) \kappa(y, x)$$

$$\pi(x) \kappa(x, y) = \pi(x) \alpha(x, y) q(x, y)$$

$$+ \pi(x) \delta_x(y) \left[\sum_t [1 - \kappa(x, t)] q(x, t) \right]$$

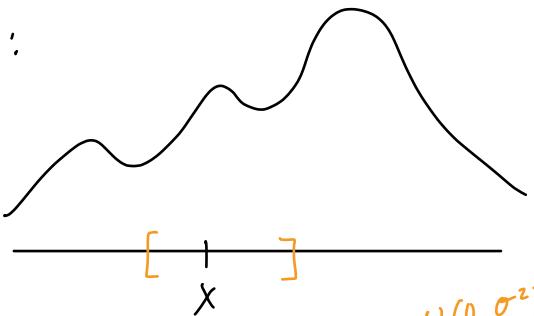
$$\begin{aligned} \pi(x) \kappa(x, y) q(x, y) &= \pi(x) q(x, y) \min \left(1, \frac{\pi(y) q(y, x)}{\pi(x) q(x, y)} \right) \\ &= \min(q(y, x) \pi(y), \pi(x) q(x, y)) \end{aligned}$$

$$\begin{aligned} \pi(x) \kappa(x, y) &= \min(q(y, x) \pi(y), \pi(x) q(x, y)) \\ &+ \pi(x) \delta_x(y) \left[\sum_t [1 - \kappa(x, t)] q(x, t) \right] \end{aligned}$$

$$\begin{aligned} \pi(y) \kappa(y, x) &= \pi(y) \alpha(y, x) q(y, x) \\ &+ \pi(y) \delta_y(x) \left[\sum_t [1 - \kappa(y, t)] q(y, t) \right] \\ &= \min(\pi(x) q(x, y), \pi(y) q(y, x)) \\ &+ \pi(y) \delta_y(x) \left[\sum_t [1 - \kappa(y, t)] q(y, t) \right] \end{aligned}$$

$$\text{so } \pi(x) \kappa(x, y) = \pi(y) \kappa(y, x)$$

Ex:



$$\text{propose: } y = x + \epsilon \sim N(0, \sigma^2) \sim \text{Unif}[-a, a]$$

$$q(y|x) = N(x, \sigma^2) \text{ if we use } N(0, \sigma^2)$$