

JSTH@ing.e.utexas.edu

S12, S13, Outlines

Basic Terminology OverviewModel type:

- If X contains binary indicators, ANOVA
- If X contains continuous covariates, Regression
- If X is a mix, 'linear model'.

• column space of X : $C(X)$, the span (set of all linear combos) of the cols of X .

• basis: a linearly indep. spanning set.

• If M is a ppm onto $C(X)$, and $v \in C(X)$, can write $Mv = v$.

If $v \in C(X)$, can write $v = Xb$; v is a linear combo of cols of X .

• vector space: $M \subseteq \mathbb{R}^n$ is a vec space if closed under addition & mult.

• $\forall x, y \in M, (x+y) \in M$

• $\forall x, y \in M, x'y \in M$

• $N \subseteq M$ is a subspace of M if N is also a vector space.

• linearly indep: $\{x_1, \dots, x_n\}$ are linearly dependent if at least one can be written as a linear combo of others.

• orthogonality: $x \perp y$ iff $x'y = 0$

• $\{x_1, \dots, x_p\}$ is an orthogonal basis if $x_i \perp x_j$, ie all elements orthog.

• orthonormal basis: An orthog basis, but all elements have length 1.

• Gram-Schmidt: Algorithm to get an orthonormal basis out of $\{x_1, \dots, x_n\}$ vecs.

• Key: Can decompose ANY vector x as $x = x_0 + x_1$, where $x_0 \in C(X)$, $x_1 \in C(X)^\perp$

This decomposition is unique.

Think of as an element in X plus a residual.

Perpendicular Projection Matrices

1) Formal definition:

M is a ppm onto K -dim subspace $C(x)$ if:

1. If $x \in C(x)$, then $Mx = x$ (leaves elems of $C(x)$ unchanged)
2. If $x \in C(x)^\perp$ then $Mx = 0$ (zeros out elems orthog to $C(x)$)

2) Constructive Definitions:

1. Let O be an orthonormal basis of $C(x)$. (Can always obtain O via Gram-Schmidt) Then $M = OO^T$.

2. $M = X(X^TX)^{-1}X'$ is a ppm to $C(x)$.
 $(I-M)$ is a ppm to $C(x)^\perp$.

3) Alt. criteria for being a ppm: M is a ppm to $C(x)$ iff

1. M is idempotent. $MM = M$
2. M is symmetric. $M = M'$

4) Rank of M : $\text{tr}(M) = r(M) = r(x)$.

Linear Algebra Review:

- A is non-singular if $\exists A^{-1}$ s.t. $AA^{-1} = A^{-1}A = I$.
- Rank: dimension of vec space spanned by cols of A :
 - # of linearly indep cols
 - $\text{tr}(A)$ if A idempotent.
 - A is full rank (rank = # cols) if non-singular. All ppm's are full rank.
- Null space: Set of all vectors x s.t. $Ax = 0$; $N(A)$
 - If A singular, $N(A) =$ just the zero vec
 - Rank of null space = # of redundant cols: If $r(A) = r$, $r(N(A)) = n - r$ for $A_{n \times n}$
- Eigenvals/vecs: $Ax = \lambda x \rightarrow x = \text{eigvec}, \lambda = \text{matching eigenvalue}$, for matrix A .
 - $(A - \lambda I)x = 0$ is another arrangement
 - multiplicity of an eigenvalue = # of times it occurs
 - Distinct eigenvalues' eigenvectors are orthogonal.
- Positive Semi-Definite: Cov matrices must be psd.
 - A is pos def (pd) if all eigenvals λ are > 0 .
 - A is pos semi def (psd) if all eigenvals λ are ≥ 0 .
- Or, A is psd if $\exists Q$ s.t. $A = QQ^T$. Useful for decomposing cov matrices.

Singular Value Decomp (SVD):

- $X = UDV^T$ where $D = \text{diag } (\lambda_j)$, where λ_j^2 are eigenvals of X^TX
 $(n \times p)$.
- $V_{p \times p} = [v_1 \dots v_p]$ with v_j = eigenvectors of X^TX , i.e. $X^TXv_j = D^2v_j$
- $U_{n \times p} = [u_1 \dots u_p]$ with u_j = eigenvectors of XX^T , i.e. $XX^Tu_j = D^2u_j$
- u_j 's and v_j 's are orthonormal

If X symmetric, reduces to $X = PDP^T$ where $D = \text{diag } (\lambda_j)$, λ_j = eigenvals of X and P = orthonormal eigenvectors of X .

Also called 'eigval decomp' or 'spectral decomp'.

Trace Properties:

- $\text{tr}(A) = \text{sum of diag elements}$
- For A symmetric, $\text{tr}(A) = \text{sum of eigenvals}$
- Cyclic movement ok inside trace: $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$
- $\text{tr}(M) = r(M)$ for ppm M .

For an idempotent matrix, only possible eigenvals are 0 and 1.

Chapter 1 Theorems:

1.1.1: For A, B constant $\sim X, Y$ random vectors, $\text{Cov}(AX + BY) = A \text{Cov}(X)A' + B \text{Cov}(Y)B' + A \text{Cov}(X, Y)B' + B \text{Cov}(Y, X)A'$

Quadratic Forms:

1.3.2: If $E(Y)=\mu$ and $\text{Cov}(Y)=\Sigma_i$, then $E(Y'Ay) = \text{tr}(A\Sigma_i) + \mu'A\mu$
(Expectation of Quadratic Forms)

1.3.3: If $Y \sim N(\mu, I)$ and M is a ppm, then $Y'MY \sim \chi^2(r(M), \mu'M\mu/2)$

1.3.6: If $Y \sim N(\mu, V)$ then $Y'AY \sim \chi^2(\text{tr}(AV), \mu'A\mu)$ under these conditions:

1. $VAVAV = VAV$
 2. $\mu'AV\mu = \mu'A\mu$
 3. $VAV\mu = V\mu$
- } where $AV = A$ in all these

1.3.7: If $Y \sim N(\mu, \sigma^2 I)$ and $BA=0$, then

1. $Y'AY + BY$ are indep (A symmetric)
2. $Y'AY$ and $Y'BY$ are indep (A, B symmetric)

2.1 - Identifiability + Estimability

Consider the linear model: $Y = X\beta + e$, $E(e) = 0$.

A parameterization for $E(Y)$ means writing $E(Y)$ as a fcn of some params β .

$$\cdot E(Y) = f(\beta)$$

A general linear model is a parameterization; $E(Y) = X\beta$

$$\text{why? } E(Y) = E(X\beta + e) = X\beta + E(e) = X\beta + 0 = X\beta$$

A parameterization is identifiable if knowing $E(Y)$ tells you β .

Identifiability:

- β is identifiable iff for any $\beta_1 \neq \beta_2$, $f(\beta_1) = f(\beta_2) \rightarrow \beta_1 = \beta_2$.
- Fcn $g(\beta)$ identifiable iff for any $\beta_1 \neq \beta_2$, $f(\beta_1) = f(\beta_2) \rightarrow g(\beta_1) = g(\beta_2)$.
- Ex: If $X\beta_1 = X\beta_2$, then $\beta_1 = (X'X)^{-1}X'X\beta_1 = (X'X)^{-1}X'X\beta_2 = \beta_2$.

Thm 2.1.2: A fcn $g(\beta)$ is identifiable iff $g(\beta)$ is a fcn of $f(\beta)$.

Estimability:

- An estimable fcn is an identifiable fcn which is linear.
- Means you can write as a fcn of design matrix X .

• A scalar-valued linear fcn of β : $\lambda'\beta = p'X\beta$ for vector p .

• A vector-valued linear fcn of β : $A'\beta = P'X\beta$ for matrix P .

• p and P are not unique, but MP is (and Mp), ie projection of p / P onto $C(X)$.

Linear Estimate: $f(Y)$ is a linear est of $\lambda'\beta$ if $f(Y) = a_0 + a_1'Y$ for scalar a_0 , vec a_1 .

• $\hat{\theta}$ is a linear est of θ if has form $\hat{\theta} = a_0 + a_1' \tilde{\theta}$

Biased Estimates:

• $\hat{\theta}$ unbiased if $E(\hat{\theta}) = \theta$

• $a_0 + a_1'Y$ is unbiased est of $\lambda'\beta$ iff $a_0 = 0$ and $a_1'X = \lambda'$

2.2 Least Squares Estimation

Model: $\mathbf{Y} = \mathbf{X}\beta + \mathbf{e}$, $E(\mathbf{e}) = \mathbf{0}$, $\text{Cov}(\mathbf{e}) = \sigma^2 \mathbf{I}$

$\hat{\beta}$ is a least sqs (LS) est if it minimizes $\min_{\beta} \{(Y - X\beta)^T(Y - X\beta)\}$

$\hat{\beta} = (X^T X)^{-1} X^T Y$, and $\hat{Y} = X\hat{\beta} = MY$ where M is ppm to $C(X)$: $M = X(X^T X)^{-1} X^T$

For estimable fctn $\lambda^T \beta = p^T X\beta$, $\hat{\lambda}^T \hat{\beta} = p^T M\hat{Y}$

If $\lambda^T = p^T X$, then $E(p^T M\hat{Y}) = \lambda^T \hat{\beta}$, i.e. $p^T M\hat{Y}$ unbiased.

Thrm 2.2.6: Let $r = r(X)$ and $\text{Cov}(\mathbf{e}) = \sigma^2 \mathbf{I}$. Then $\hat{\sigma}^2 = \frac{Y^T (I - M) Y}{n-r}$ is unbiased est. of σ^2 .

$\hat{\sigma}^2 = \frac{Y^T (I - M) Y}{n-r} = \frac{\text{RSS}}{n-r}$ RSS = squared length of residual vector.

RSS is also called SSE, sum of squared error

$\hat{\sigma}^2$ is also called MSE, mean squared error.

$(n-r) = \text{rank}(I - M) = \text{df error}$

2.3 Best Linear Unbiased Estimation (BLUE)

"Best" = smallest variance.

$a^T Y$ is BLUE for $\lambda^T \beta$ if $E(a^T Y) = \lambda^T \beta$

$\text{Var}(a^T Y) \leq \text{Var}(b^T Y)$ for any other unbiased $b^T Y$.

Gauss-Markov Theorem: If $\lambda^T \beta$ estimable, then $\hat{\lambda}^T \hat{\beta}_{\text{LS}}$ is BLUE.

If $\lambda^T \beta = p^T X\beta = X\beta$, $\hat{\beta}_{\text{LS}}$ is BLUE for β . \rightarrow Above is more general version.

Proof: Let M be the ppm to $C(X)$. $\lambda^T \beta = p^T X\beta$, so LS est of $\lambda^T \beta$ is $\lambda^T \hat{\beta}_{\text{LS}} = p^T M\hat{Y}$.

Let $a^T Y$ be another unbiased estimator of $\lambda^T \beta$.

Show $\text{Var}(a^T Y) \geq \text{Var}(p^T M\hat{Y})$.

$$\text{Var}(a^T Y) = \text{Var}(a^T Y - p^T M\hat{Y} + p^T M\hat{Y})$$

$$= \underbrace{\text{Var}(a^T Y)}_{\geq 0} + \text{Var}(p^T M\hat{Y}) - \underbrace{2 \text{Cov}(a^T Y - p^T M\hat{Y}, p^T M\hat{Y})}_{= 0}$$

Why does $\text{Cov}(a^T Y - p^T M\hat{Y}, p^T M\hat{Y}) = 0$?

$$\text{Cov}(a^T Y - p^T M\hat{Y}, p^T M\hat{Y}) = (a^T - p^T M)\text{Cov}(Y)M p$$

Then since $a^T Y$ unbiased, $\lambda^T \beta = E(a^T Y) = a^T X\beta$, so $p^T X\beta = a^T X\beta$, i.e. $p^T X = a^T X$.

Then $M = X(X^T X)^{-1} X^T$, so have:

$(a^T - p^T M)\text{Cov}(Y)M p = \sigma^2 (a^T - p^T M)M p = \sigma^2 (a^T M - p^T M)p$, and from above,

$$a^T M = p^T M, \text{ so } = 0.$$

Corollary 2.3.3: If $\sigma^2 > 0$, \exists a unique BLUE for any est. fctn $\lambda^T \beta$.

2.4 Max Lhood Estimation

Assume $Y \sim N(X\beta, \sigma^2 I)$. Then can estimate β MLE and σ^2 MLE by maximizing

$$L = (\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} (Y - X\beta)^T (Y - X\beta) \right]$$

For convenience, maximize log L. Gives same result as LS for $\hat{\beta} = (X^T X)^{-1} X^T Y$

$\hat{\sigma}^2_{\text{MLE}} = \frac{Y^T (I - M) Y}{n}$ is rarely used: MSE is used for $\hat{\sigma}^2$ bc unbiased.

2.6 Sampling Distributions of Estimators

$Y \sim N(X\beta, \sigma^2 I)$ $\hat{\beta} \sim N(\beta, \sigma^2 (X^T X)^{-1})$

For $\lambda^T \beta$, least sqs est is $\hat{\lambda}^T \hat{\beta} = p^T M\hat{Y}$. Can write sampling dist in 2 ways:

$$\hat{\lambda}^T \hat{\beta} \sim N(\lambda^T \beta, \sigma^2 \lambda^T (X^T X)^{-1} \lambda)$$

$$p^T M\hat{Y} \sim N(\lambda^T \beta, \sigma^2 p^T M p)$$

2.7 Generalized Least Squares

Model: $Y = X\beta + \mathbf{e}$, $E(\mathbf{e}) = \mathbf{0}$, $\text{Cov}(\mathbf{e}) = \sigma^2 V$ where V is psd, known. (1)

Write $V = QQ^{-1} \rightarrow Q^T V Q^{-T} = I$, so left-multiply and solve new model:

$$\begin{aligned} Q^{-1} Y &= Q^{-1} X\beta + Q^{-1} \mathbf{e} \\ Y^* &= X^* \beta + \mathbf{e}^* \end{aligned} \quad \text{with } E(\mathbf{e}^*) = \mathbf{0}, \text{Cov}(\mathbf{e}^*) = \sigma^2 I \quad (2)$$

Called Gen LS bc no longer minimizing squared dist b/w $Y = X\beta$, but minimizing a generalized squared distance specified by V^{-1} .

Thrm 2.7.1:(a) $\lambda^1\beta$ is estbl in model (1) iff estbl in model 2.(b) $\hat{\beta}$ is the gls est iff $\hat{\beta} = (X'V^{-1}X)^{-1}X'V^{-1}y$ · for any estimable $\gamma(\beta)$, \exists a unique LS soln.(c) $X'\hat{\beta}$ is the BLUE of $X'\beta$.(d) If $e \sim N(0, \sigma^2 V)$ then $X'\beta$ is the min. var unbiased estimator.(e) If $e \sim N(0, \sigma^2 V)$ then $X'\beta$ is also the MLE.Proof:(a) If $\lambda^1\beta$ estbl in (1), can write $\lambda^1 = p'X = (p'Q)Q^{-1}X$ so $X'\beta$ estbl in (2).If $X'\beta$ estbl in (2), can write $\lambda^1 = p'Q^{-1}X = (p'Q^{-1})X$ so $\lambda^1\beta$ estbl in (1).

(b) $(X'Q^{-1}Q^{-1}X)^{-1}X'Q^{-1}Q^{-1}Y = Q^{-1}X\beta \rightarrow (X'V^{-1}X)^{-1}X'V^{-1}Y = \hat{\beta}$

(c) $X'\hat{\beta}$ is BLUE for $\sigma^2 Y$ by Thrm 2.2.3. Any fit of Y can be obtained through $\sigma^2 Y$ bc Q^{-1} invertible. So the gls est is BLUE.2.8 - Normal Equations:Solving for β gives same $\hat{\beta}$ as LS.

· $X'X\beta = X'y$

· $X'V^{-1}X = X'V^{-1}y$ for gen LS

Regression in canonical form & ridge regression are two techniques to deal w. collinearity - when cols of X are nearly linearly dependent, ie very correlated.

15.2 - Regression in Canonical Form

· gives a convenient linear transformation so can read off LS soln.

· Model: $Y = X\beta + e$

· First, need SVD:

Thrm 15.1.2: SVD (Singular Value Decomp)· Let $X_{n \times p}$ be full rank. Can write $X = UDV^T$, where· $D_{p \times p} = \text{diag}(\lambda_j)$ where λ_j are sqrts of eigenvals of $X'X$ (λ_j^2 = eigenval)· $U_{n \times p}$ = eigenvecs of $X'X$, so $X'XU = UD^2$ · $V_{p \times p}$ = eigenvecs of $X'X$, so $X'XV = UD^2$ Canonical form:· Begin with $Y = X\beta + e$ and write $X = UDV^T$.· Write $U_* = [U_*, U_+]$ where cols of U_* are an orthonormal basis for \mathbb{R}^n .· Left-multiply by U_* to transform model. $U_*'Y = U_*'X\beta + U_*'e$ · Then model is $Y_* = U_*'X\beta + e_*$, $E(e_*) = 0$, $\text{Cov}(e_*) = \sigma^2 U_*'U_* = \sigma^2 I$ · Then use SVD again on $U_*'X\beta$ term:

· $U_*'X\beta = U_*'UDV^T\beta = \begin{bmatrix} U_* \\ U_+ \end{bmatrix} UDV^T\beta = \begin{bmatrix} L \\ 0 \end{bmatrix} V^T\beta$

· Reparameterize by letting $\gamma = V^T\beta$, to get canonical model:

· $Y_* = \begin{bmatrix} L \\ 0 \end{bmatrix} \gamma + e_*$, $E(e_*) = 0$, $\text{Cov}(e_*) = \sigma^2 I$

· Parameter estimation in this model:

· $\hat{\gamma} = D'U_*'y$

· $\text{Cov}(\hat{\gamma}) = \sigma^2 D^2$

· $\text{RSS} = Y_*' \begin{bmatrix} 0 & 0 \\ 0 & I_{n-p} \end{bmatrix} Y_*$

15.3 - Ridge Regression:

- Originally proposed to deal w. collinearity. Now viewed as form of penalized likelihood estimation. (Bayesian)
- ℓ_2 -norm penalty
- Shrinks coefficients, but none all the way to zero.
- Loss func: $\min_{\beta} \left\{ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \right\} = \min_{\beta} \left\{ \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$

MSE intuition: $MSE = E[(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)] = \text{tr}[\sigma^2 (X^T X)^{-1}]$ since $E(X^T Q X) = \text{tr}(Q \Sigma) + u^T Q u$

$$= \sigma^2 \sum_{j=1}^p \frac{1}{\lambda_j^2}$$
 where $\lambda_j^2 = \text{eigvals}$

$$\cdot E(\hat{\beta} - \beta) = u = 0$$

$$\cdot Q = I$$

$$\cdot \Sigma = \text{cov}(\hat{\beta})$$

- If 1+ eigenvals too close to 0, MSE blows up.
- For RR, mse becomes $\sigma^2 \sum_{j=1}^p \frac{1}{\lambda_j^2 + \lambda}$, preventing denoms from being too small.

Bayesian Interpretation:

- RR can be interpreted as a posterior mean w. a normal prior.
- Lhood: $p(y|\beta) \sim N(X\beta, \sigma^2 I)$
- Prior: $p(\beta) \sim N(0, \frac{\sigma^2}{\lambda} I)$

The normal-normal update yields posterior

$$p(\beta|y) \sim N(m, V) \text{ with } V = \left[\frac{1}{\sigma^2} X^T X + \frac{\lambda}{\sigma^2} I \right]^{-1}$$

$$m = V^{-1} X^T y = \left[\frac{1}{\sigma^2} (X^T X + \lambda I)^{-1} X^T y \right] = \frac{\hat{\beta}_{RR}}{\sigma^2}$$

LS solution: has closed form, $\hat{\beta}_{RR} = (X^T X + \lambda I)^{-1} X^T y$

How does RR handle collinearity? See SM1, HW2, Prob 5:

- Run RR on model with $X = [x_1 \dots x_p]$. The $\hat{\beta}_{RR,i} = a$ for x_1 . If we add $(m-1)$ copies of x_1 , how do the identical copy coeffs relate to a ?

Model 1 Loss: $\min_{\beta} \left\{ \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\} = \min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \beta_0 - \beta_1 x_{i1} - \sum_{j=2}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$

Model 2 Loss: Now $j=1$ to m are x_1 predictors, and $j+1$ to p are $x_2 \dots x_p$.

$$\min_{\beta} \left\{ \sum_{i=1}^n \left(y_i - \tilde{\beta}_0 - \sum_{j=1}^m \tilde{\beta}_j x_{i1} - \sum_{j=m+1}^p \tilde{\beta}_j x_{ij} \right)^2 + \lambda \left(\sum_{j=1}^m \tilde{\beta}_j^2 + \sum_{j=m+1}^p \tilde{\beta}_j^2 \right) \right\}$$

Want $\tilde{\beta}_1 x_{i1} = \sum_{j=1}^m \tilde{\beta}_j x_{i1} \rightarrow a x_{i1} = m \tilde{\beta}_1 x_{i1} \rightarrow \tilde{\beta}_1 = a/m$.

VARIABLE SELECTION METHODS

1. LASSO: Regression with a ℓ_1 -norm penalty.

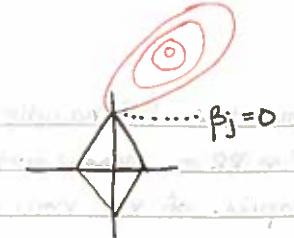
• Shrinks some coeffs all the way to 0, unlike ridge.

• Downside: shrinks even non-zeroed-out coeffs.

$$\hat{\beta}_L = \min_{\beta} \left\{ (y - X\beta)^T(y - X\beta) + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

• No closed-form soln; solve using SGD or other optimization techniques.

• Have to tune penalty parameter λ .



3. Horseshoe Prior:

$$L: p(y|\beta, \sigma^2) \sim N(X\beta, \sigma^2 I)$$

$$P: p(\beta_i | \lambda_i) \sim N(0, \lambda_i^2)$$

$$p(\lambda_i | \tau) \sim C^+(0, \tau)$$

$$p(\tau | \sigma^2) \sim C^+(0, \sigma^2)$$

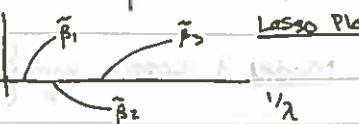
$$p(\sigma^2) \sim 1/\sigma^2 \text{ (Jeffreys)}$$

• Shrinks some params, leaves others as is.

• Less shrinkage on strongest predictors.

$C^+(0, \tau) = \text{half-cauchy w/ scale } \tau$

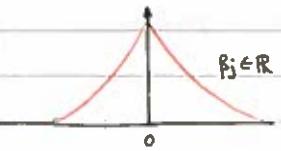
$$p(x) = \frac{1}{\pi \lambda [1 + x^2]}$$



2. Bayesian Lasso: We can do some reverse engineering to interpret $\hat{\beta}_L$ as the posterior mean of a Bayesian model.

• Lhood: $p(y|\beta, \sigma^2) \sim N(X\beta, \sigma^2 I)$

• Prior: $p(\beta_j) \sim \frac{\lambda}{2\sigma} \exp\left[-\frac{\lambda}{\sigma} |\beta_j|\right]$ The 'laplace', or double exponential, prior.



• Posterior: $\log p(\beta|y, \sigma^2) = c - \frac{1}{2\sigma^2} (y - X\beta)^T(y - X\beta) - \frac{\lambda}{\sigma} \sum_{j=1}^p |\beta_j|$

• Including the highlighted σ in $p(\beta_j)$ ensures posterior is unimodal

$$\text{Note: } p(\beta_j | \sigma) = \frac{1}{2} \cdot \frac{\lambda}{\sigma} \exp\left[-\frac{\lambda}{\sigma} |\beta_j|\right] = \int N(\beta_j | 0, \sigma^2 \gamma_j^2) \cdot \text{Exp}[\gamma_j^2 | \frac{\lambda^2}{2}] d\gamma_j^2$$

So we can replace the prior on β with:

$$p(\beta_j | \gamma_j^2, \sigma^2) \sim N(0, \sigma^2 \gamma_j^2)$$

$$p(\gamma_j) \sim \text{Exp}(\lambda^2/2)$$

$$\text{Proof: } \frac{\alpha}{2} e^{-\alpha|z|} = \int_0^\infty \frac{1}{\sqrt{2\pi s}} e^{-\frac{z^2}{2s}} \cdot \frac{\alpha^2}{2} e^{-\frac{\alpha^2 s}{2}} ds \quad \text{where } z = \beta_j / \sigma$$

$$\alpha = \lambda$$

$$s = \gamma_j^2$$

3.1 - Testing Background

- Model: $y = X\beta + e$, $E(e) = 0$, $\text{Cov}(e) = \sigma^2 I$
- The ppm M is the crucial item for model testing.
- $C(X)$ = estimation space, $C(X)^\perp$ = error space, so $(I-M)$ is ppm to error space.
- Any two models w. same estimation space are essentially same model.

3.2 - Testing Models

- Full model, which we assume correct: $y = X\beta + e$, $E(e) = 0$, $\text{Cov}(e) = \sigma^2 I$. (M1)
- Want to test if a reduced model is better: $y = X_0\gamma + \eta$, $E(\eta) = 0$, $\text{Cov}(\eta) = \sigma^2 I$. (M0)
- $C(X_0) \subseteq C(X)$
- Test statistic: $\frac{y'(M_1 - M_0)y / r(M_1 - M_0)}{y'(I - M_0)y / r(I - M_0)} \sim F(df_1, df_0, ncp)$
- $ncp: ncp = \beta' X' (M - M_0) X \beta / (2\sigma^2)$
- $ncp = 0$ under M_0 , since $(M - M_0) = M(I - M_0)$, and $I - M_0 \in C(X_0)^\perp \subseteq C(X)^\perp$, so $M(I - M_0) = 0$.
- Null Hypothesis: $H_0: M_0$ is sufficient, reject for large F stat.

Proof of distribution of test stat:

- $y'(M_1 - M_0)y \perp\!\!\!\perp y'(I - M_0)y$, since by Thrm 1.3.7, if $(M_1 - M_0)(I - M_0) = 0$, are indep. and $(M_1 - M_0)(I - M_0) = (M_1 - M_0)M_0(I - M_0) = 0$
- Numerator: $y'(M_1 - M_0)y \rightarrow$ First, $y \sim N(X\beta, \sigma^2 I)$ so $y/\sigma \sim N(X\beta, I)$. Then by Thrm 1.1.3, $y'(M_1 - M_0)y \sim \chi^2[r(M_1 - M_0), \frac{\beta' X' (M - M_0) X \beta}{2\sigma^2}]$
- Denominator: $y'(I - M_0)y \rightarrow$ By Thrm 1.1.3, $y'(I - M_0)y \sim \chi^2[r(I - M_0), \frac{\beta' X' (I - M_0) X \beta}{2\sigma^2}]$

We can set up a reduced model by specifying estimable constraints, $\lambda' \beta = 0$.

But sometimes this doesn't work, like for $H_0: \beta_1 = \beta_2 + 5$. Need a generalized testing procedure.

3.2.1 - Generalized Test Procedure:

$$M_1: Y = X\beta + e$$

$$M_0: Y = X_0\gamma_0 + \eta + Xb \text{ where } b = \text{constant}$$

$$\text{Procedure: } M_1: Y - Xb = (X\beta - Xb) + e \rightarrow M_1: Y = X\beta + e$$

$$\text{where } \beta^* = \beta - b$$

$$M_0: Y^* = X_0\gamma^* + \eta$$

Then F-test like usual.

$$\text{Example: } M_1: Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + e$$

$$\text{Want to test } H_0: \beta_1 = 0, \text{ or } \beta_2 = \beta_3 + 5$$

$$M_0: Y = \beta_0 + (\beta_3 + 5)X_2 + \beta_3 X_3 + e$$

$$Y = \beta_0 + \beta_3(X_2 + X_3) + 5X_2 + e$$

$$Y = \beta_0 + \beta_3(X_2 + X_3) + 5X_2 + e$$

$$(Y - 5X_2) = (\beta_0 - 5X_2) + (\beta_3 - 5X_2)(X_2 + X_3) + e$$

$$\text{Numerator for F-Test: } (Y - 5X_2)^T (M_1 - M_0)(Y - 5X_2)$$

3.3 - Testing Linear Parametric Facts:

$$\cdot M_1: Y = X\beta + e$$

$$\cdot M_0: Y = X\beta + e, \text{ subject to } \lambda' \beta = 0, \text{ where } \lambda' \beta \text{ testable, s.t. } \lambda' \beta = P' X \beta.$$

$$\cdot \text{ Pick a matrix } U \text{ so } C(U) = C(\lambda), \text{ i.e. } U \text{ s.t. } \lambda' U = 0.$$

$$\cdot \text{ Set } X_0 = XU$$

$$\cdot \text{ Then constraint is } \beta = U\gamma \rightarrow M_0 \text{ is now } Y = X_0\gamma + \eta$$

$$\cdot \text{ Number cols for } U: \dim(U) = \dim(X) - \dim(\lambda)$$

Example: 1-Way ANOVA

$$Y = \begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{21} \\ Y_{22} \\ Y_{31} \\ Y_{32} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \cdot \text{ So } \beta = [u, d_1, d_2, d_3]^T$$

$$\cdot \text{ Want to test } \lambda' = (0, 1, 0, -1)$$

$$\cdot \text{ Can pick } U = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

$$\text{so that } \lambda' U = (0, 1, 0, -1) \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = (0, 0, 0) \checkmark$$

• Can pick any $U \rightarrow Y$ not so easy to see, use Gram-Schmidt.

$$\cdot \text{ Now } M_0: Y = X_0\gamma + \eta \text{ where } X_0 = XU$$

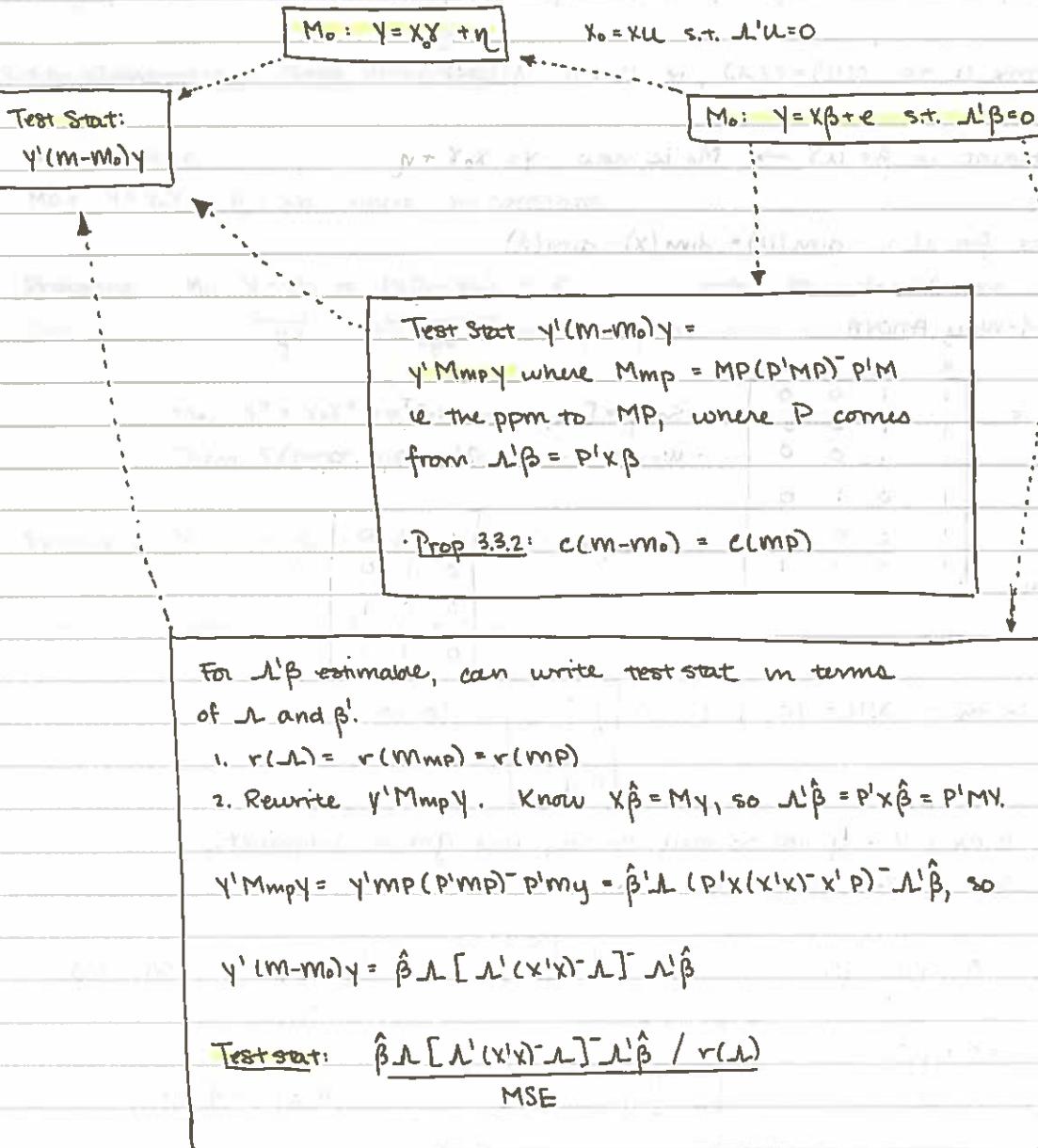
$$\cdot \text{ Example 2: Consider joint constraint } \lambda' = \begin{bmatrix} \lambda'_1 \\ \lambda'_2 \end{bmatrix} \text{ where } \lambda'_1 = (0, 1, 0, -1) \text{ s.t. } \lambda' \beta = 0$$

$$\cdot \text{ Want } \lambda' U = 0: \begin{bmatrix} 0 & 1 & 0 & -1 \\ 0 & 1 & 1 & 1 \end{bmatrix} \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & -2 \\ 0 & 1 \end{bmatrix}}_U = 0$$

cols $U = 4 - 2 = 2$

$$\cdot X_0 = XU, \text{ and } M_0: Y = X_0\gamma + \eta$$

How to get back & forth, if have full model $y = X\beta + e$?



3.3 - What happens if you pick a non-estimable constraint?

- Say pick $M_0: y = X\beta + e, L'\beta = 0$ for a non-estimable constraint.
- Call $\lambda_0'\beta$ the estimable part; then $\lambda_0'\beta$ and $L'\beta$ give the same reduced model, since the non-estimable part doesn't induce any constraint on $C(X)$.
- Thm 3.3.6: If $C(\lambda) \cap C(X) = C(\lambda_0)$, and $C(u_0) = C(\lambda_0)^\perp$, then $C(Xu) = C(Xu_0)$. Thus $L_0'\beta$ and $L'\beta$ induce the same model.
- Prop 3.3.7: If $L'\beta = 0$ estab, and $\lambda \neq 0$, then $\lambda'\beta = 0 \rightarrow C(Xu) \neq C(Xu_0)$. (for an estab constraint, there is always a restraint on $C(X)$.)
- Cor. 3.3.8: $C(\lambda) \cap C(X) = \{0\}$ iff $C(Xu) = C(X)$, ie if $L'\beta$ has no estab part, the constraint doesn't affect the model.

3.4 - Decomposition of Sums of Squares

Total sum of squares (SS): $SS_{\text{tot}} = SS = \mathbf{y}'\mathbf{y}$

$\cdot SSR(x) = \text{sum of sqs for regression on } x: SSR(x) = \mathbf{y}'M\mathbf{y}$

$\cdot SSE = \text{sum of sqs for error: } SSE = \mathbf{y}'(I-M)\mathbf{y}$

$$SS = \mathbf{y}'M\mathbf{y} + \mathbf{y}'(I-M)\mathbf{y} \quad (\text{by adding/subtracting } M\mathbf{y} \text{ to } \mathbf{y}'\mathbf{y})$$

\cdot Say M_0 is a reduced model: $x = [x_0; x_1]$ with M_0 as ppm for reduced model.

\cdot Then we can decompose $\mathbf{y}'M\mathbf{y}$:

$$\mathbf{y}'M\mathbf{y} = \mathbf{y}'(M-M_0) + \mathbf{y}'M_0\mathbf{y}$$

$$SSR(x) = SSR(x_1|x_0) + SSR(x_0)$$

\cdot Other notation: Can write in terms of parameter names, say $x = [\alpha; \eta]$

$$\text{Then } SSR(x) = R(\alpha|\alpha) + R(\eta)$$

\cdot Conceptual note: $M-M_0=Mx$, so M is ppm to $C(M-M_0)$.

\cdot Say $R = [r_1 \dots r_p]$ is an orthonormal basis of $C(x)$.

$$\text{Then } M = RR^T = \sum_{s=1}^p R_s R_s^T$$

↓

\cdot So we can write any space as an additive decomp of orthog. projection matrices. We only want to do this when $M^*|A$ have meaning, like M_0 .

When order added to model matters?

\cdot When is $R(\alpha|\eta, \mu) = R(\alpha|\mu)$? This means α, η not confounded; order doesn't matter.

Prop 3.6.3: $R(\alpha|\eta, \mu) = R(\alpha|\mu)$ iff $C(M_i - M_j) \perp C(M_0 - M_j)$

\cdot Equivalently, $\underbrace{(M_i - M_j)}_{\text{Tests dropping } n} = \underbrace{(M_0 - M_j)}_{\text{Tests dropping } n}$ iff $(M_i - M_j)(M_0 - M_j) = 0$

Tests dropping
n
from model
with
 μ, η

$\left\{ \begin{array}{l} M_i = \text{model with } \mu, \eta \\ M_0 = \text{model with } \mu, \alpha \\ M = \text{model with } \mu, \alpha, \eta \end{array} \right.$

4.1 - One way ANOVA Model

\cdot Model: $y_{ij} = \mu + \alpha_i + \epsilon_{ij}$ for $i=1\dots t, j=1\dots n_i$ (t indexes trts, j indexed obs in ea. trt)

$\cdot \mu = \text{grand mean}$

$\cdot \alpha_i = \text{effect of } i\text{th treatment}$

\cdot Common reduced model: $y_{ij} = \mu + \alpha_{ij}; H_0: \alpha_1 = \dots = \alpha_t = 0$

\cdot Can also test various contrasts.

$$\cdot \mathbf{y} = \mathbf{x}\beta + \mathbf{e}, E(\mathbf{e}) = 0, \text{cov}(\mathbf{e}) = \sigma^2 I$$

$$\cdot \mathbf{x} = \begin{bmatrix} 1 & x_1 & x_2 & x_3 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \end{bmatrix} \quad \text{and } \beta = \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} \quad (\text{size } n \times n \text{ matrix of } 1's)$$

$$\begin{bmatrix} \frac{1}{n_1} & \frac{1}{n_1} & \frac{1}{n_1} \\ \frac{1}{n_1} & \frac{1}{n_1} & \frac{1}{n_1} \\ \frac{1}{n_2} & \frac{1}{n_2} & \frac{1}{n_2} \\ 0 & \frac{1}{n_2} & \frac{1}{n_2} \\ 0 & 0 & \frac{1}{n_3} \end{bmatrix}$$

$$\cdot M = \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}' = \text{block diagonal s.t. } \text{diag}[\frac{1}{n_1} J_{n_1}^{n_1}] =$$

$$(n_1=3, n_2=2, n_3=1)$$

\cdot can't just invert $(\mathbf{x}'\mathbf{x})$ since not full rank - used ginv.

\cdot Can also derive M using $Z = [\alpha_1 \alpha_2 \alpha_3]$ b/c $C(x) = C(Z)$. So $M = Z(Z'Z)^{-1}Z'$

$$\cdot$$
 Fitted values: $\hat{\mathbf{y}} = M\mathbf{y} = [\bar{y}_1 \bar{y}_2 \dots \bar{y}_t]^T$

\cdot logical bc model assumes same means for all obs under ea. trt.

$$\cdot \hat{\mathbf{y}} \in C(x)$$

$$\cdot \mathbf{y} = \hat{\mathbf{y}} + \hat{\mathbf{e}}, \text{ so residuals } \hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - M\mathbf{y} = (I-M)\mathbf{y}, \text{ so } \hat{\mathbf{e}} \in C(x)^\perp$$

Projection matrices:

$$\cdot M_{\alpha} = \mathbf{x}(\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}' \text{ for } \mathbf{x} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} = \frac{1}{n} J_n^n, \text{ size } (n \times n) \text{ matrix of } 1/n.$$

$$\cdot M_{\alpha\mu} = M - M_{\mu}$$

$$\cdot \text{Then } M_{\alpha\mu}y = (M - M_{\mu})y = (\bar{y}_1 - \bar{y}, \dots, \bar{y}_t - \bar{y})$$

\cdot Testing for Contrasts: Use same F-test as w/ other $M_0: \mathbf{y} = \mathbf{x}\beta + \mathbf{e}$ s.t. $\mathbf{L}'\beta = 0$ for estimable $\mathbf{L}'\beta$ fctns.

(cont'd with bottom of next pg)

Decomp of SS:

$$\mathbf{y}'\mathbf{y} = \mathbf{y}'\mathbf{M}\mathbf{y} + \mathbf{y}'(\mathbf{I}-\mathbf{M})\mathbf{y} \quad \text{where } \mathbf{y}'\mathbf{M}\mathbf{y} = \mathbf{y}'\mathbf{M}\mathbf{u}\mathbf{y} + \mathbf{y}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{y}$$

For ANOVA Table, need Expected Mean Square, ie $E(\mathbf{y}'\mathbf{A}\mathbf{y})$ for each SS term.

In general, we know $E(\mathbf{y}) = \mathbf{x}\beta$, and $\text{Var}(\mathbf{y}) = \sigma^2\mathbf{I}$, and $E(\mathbf{y}'\mathbf{A}\mathbf{y}) = \text{tr}(\mathbf{A}\text{Var}(\mathbf{y})) + \mathbf{u}'\mathbf{A}\mathbf{u}$

Source	Expected Mean Square: $E(\text{SS}/df)$	df
$\mathbf{y}'\mathbf{M}\mathbf{u}\mathbf{y}$	$E(\mathbf{y}'\mathbf{M}\mathbf{u}\mathbf{y}/1) = \sigma^2 + \beta'\mathbf{x}'\mathbf{M}\mathbf{u}\mathbf{x}\beta$	$r(\mathbf{M}\mathbf{u})=1$
$\mathbf{y}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{y}$	$E(\mathbf{y}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{y}/t-1) = \sigma^2 + \beta'\mathbf{x}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{x}\beta/t-1$	$r(\mathbf{M}-\mathbf{M}\mathbf{u})=t-1$
$\mathbf{y}'(\mathbf{I}-\mathbf{M})\mathbf{y}$	$E(\mathbf{y}'(\mathbf{I}-\mathbf{M})\mathbf{y}/n-t) = \sigma^2 + \beta'\mathbf{x}'(\mathbf{I}-\mathbf{M})\mathbf{x}\beta/n-t = \sigma^2$ $= 0 \text{ since } (\mathbf{I}-\mathbf{M})\mathbf{x}=0$	$r(\mathbf{I}-\mathbf{M})=n-t$

ANOVA Table:

Source	df	SS	$E(\text{MS}) = E(\text{SS}/df)$
Grand Mean	1	$\mathbf{y}'\mathbf{M}\mathbf{u}\mathbf{y}$	$\sigma^2 + \beta'\mathbf{x}'\mathbf{M}\mathbf{u}\mathbf{x}\beta$
Treatments	$t-1$	$\mathbf{y}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{y}$	$\sigma^2 + \beta'\mathbf{x}'(\mathbf{M}-\mathbf{M}\mathbf{u})\mathbf{x}\beta/t-1$
Error	$n-t$	$\mathbf{y}'(\mathbf{I}-\mathbf{M})\mathbf{y}$	σ^2
Total	n	$\mathbf{y}'\mathbf{y}$	

4.2 - Estimation & Contrasts for ANOVA

$\lambda'\beta$ where \mathbf{u} is not involved: $\lambda' = (0, \lambda_1, \dots, \lambda_t)$

Need to find a contrast λ so can write as $\mathbf{p}'\mathbf{x}\beta$, where $\mathbf{p}'\mathbf{J}=0$ so 1st element of $\lambda=0$.

$$\mathbf{x}'\hat{\beta} = \sum_{i=1}^t \lambda_i \bar{y}_i = \mathbf{p}'\mathbf{M}\mathbf{y}$$

Use same F-test as for any other estimable fctn constraint:

$$H_0: \mathbf{y} = \mathbf{x}\beta + \epsilon \text{ s.t. } \lambda'\beta = 0$$

A regression model is any linear model $\mathbf{y} = \mathbf{x}\beta + \epsilon$ where $\mathbf{x}'\mathbf{x}$ is non-singular (rank p).

6.1 - Simple Linear Regression:

Everything we did before applies; now just deriving non-matrix forms.

Model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$, or $\mathbf{y} = \mathbf{x}\beta + \epsilon$ where $\mathbf{x} = [\mathbf{1} \ \mathbf{x}]$, $\beta = [\beta_0 \ \beta_1]^T$

$$\begin{aligned} \text{LS estimates: } \hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} & \hat{\beta}_1 &= \frac{SS_{xy}}{SS_x} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \end{aligned}$$

Derivation with mean-centered covariates:

$$\mathbf{x} = [\mathbf{J} \ \mathbf{x}] \rightarrow \mathbf{x}_* = [\mathbf{J} \ \mathbf{x} - \bar{x}]$$

note that $\mathbf{J} \perp \mathbf{L}(x_i - \bar{x})$, so two cols are orthog.

$$\text{Can write } \mathbf{x}_* = \mathbf{x}\mathbf{U}, \text{ where } \mathbf{U} = \begin{bmatrix} 1 & -\bar{x} \\ 0 & 1 \end{bmatrix} \text{ so } [\mathbf{J} \ \mathbf{x}_*] \begin{bmatrix} 1 & -\bar{x} \end{bmatrix} = [\mathbf{J} \ \mathbf{x}_* - \bar{x}]$$

Model spec. is unchanged: $E(\mathbf{y}) = \mathbf{x}\beta = \mathbf{x}_*\beta$

$$\text{Solve for } \mathbf{y}: \mathbf{x}\beta = \mathbf{x}\mathbf{U}\beta, \text{ so } \mathbf{y} = \mathbf{U}'\beta$$

$$\hat{\beta} = (\mathbf{x}_*^T \mathbf{x}_*)^{-1} \mathbf{x}_*^T \mathbf{y}$$

$$(\mathbf{x}_*^T \mathbf{x}_*) = \begin{bmatrix} n & 0 \\ 0 & \sum_{i=1}^n (x_i - \bar{x})^2 \end{bmatrix} \rightarrow (\mathbf{x}_*^T \mathbf{x}_*)^{-1} = \begin{bmatrix} 1/n & 0 \\ 0 & 1/\sum_{i=1}^n (x_i - \bar{x})^2 \end{bmatrix}$$

$$\hat{\beta} = \begin{bmatrix} 1/n & 0 \\ 0 & 1/\sum_{i=1}^n (x_i - \bar{x}) \end{bmatrix} \begin{bmatrix} -\bar{x} \\ -\bar{y} - \bar{x}\bar{\beta} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \bar{y} \\ SS_{xy}/SS_x \end{bmatrix}$$

↑ since $\bar{y}=0$, $\sum (x_i - \bar{x})y_i = SS_{xy}$

6.2 - Multiple Regression

- Model: $y = X\beta + e$ where $e \sim N(0, \sigma^2 I)$ $\rightarrow y \sim N(X\beta, \sigma^2 I)$
- $\hat{\beta} = (X'X)^{-1}X'y$
- $\text{Cov}(\hat{\beta}) = \sigma^2 (X'X)^{-1}$
- $\text{SSR}(x) = y'My = \hat{\beta}'X'X\hat{\beta}$
- $\text{SSE} = y'(I-M)y$
- $\text{dFE} = r(I-M) = n-p$
- $\hat{\sigma}^2 = \text{MSE} = y'(I-M)y / r(I-M)$
- $M = X(X'X)^{-1}X'$

- Any $\lambda'\beta$ is estimable, since X is full rank.
- $\lambda'\hat{\beta} \sim N(\lambda'\beta, \sigma^2 \lambda'(X'X)^{-1}\lambda)$

- Tests - CI based on $\lambda'\beta$ are:

$$TS = \frac{\lambda'\hat{\beta} - \lambda'\beta}{\sqrt{\text{MSE} \cdot \lambda'(X'X)^{-1}\lambda}} \sim t_{dF}$$

- SS Decomp:

$$y'y = y'My + y'(I-M)y$$

$$\begin{aligned} SS &= \text{SSR}(x) + \text{SSE} & \text{SSR}(x) &= y'(M - M_S)y + y'M_Sy \\ && \text{SScovariates} & \text{SSint} \end{aligned}$$

- Alt. Mean-Centred Model:

$M_S = \frac{1}{n} J_n^T$, the ppm for the intercept-only model.

$$y = [j: (1-M_S)Z] \begin{bmatrix} Y_0 \\ Y_* \end{bmatrix} + e, \text{ ie } y = X_*y + e. \text{ Note } C(X_*) = C(X).$$

y to β mapping: $y_i = \bar{y}_0 + \sum_{j=1}^p (x_{ij} - \bar{x}_{ij})\beta_j + \eta_i$ Then distribute:

$$y_i = \bar{y}_0 - \sum_{j=1}^p (\bar{x}_{ij}Y_j) + \sum_{j=1}^p x_{ij}Y_j + \eta_i$$

$\underbrace{\bar{y}_0}_{\beta_0}$ $\underbrace{\sum_{j=1}^p x_{ij}Y_j}_{\text{for } j=1 \text{ to } p, Y_j = \beta_j}$

- Normal Eqns:

$$\begin{bmatrix} X_*'X_* \\ n & 0 \\ 0 & Z'(1-M_S)Z \end{bmatrix} \begin{bmatrix} Y \\ Y_0 \\ \dots \\ Y_* \end{bmatrix} = \begin{bmatrix} X_*'y \\ \bar{y}_0 \\ \dots \\ Z'(1-M_S)y \end{bmatrix}$$

Coeff of Determination + Multiple Corr Coef:

$$R^2 = \frac{\text{SSReg}}{\text{SSTO}} = \frac{y'(M - M_S)y}{y'y - y'M_Sy}$$

"fraction of non-intercept SS explained by model"

- proportion of variability explained by covariates, out of total variability excl. intercept.
- The sample equivalent of the mult corr coeff r^2 .
- $r^2 = \max_{\lambda \in \mathbb{R}} \{ \text{corr}(y, \lambda + \beta'x) \}$ maximum corr with y achievable by a linear predictor.
- Achieved by the Best Linear Predictor (BLP), so $r^2 = \text{corr}(y, \bar{y}_0 + \beta'_* (x - M_S x))$
- Can also write as $r^2 = \frac{\beta'_* V_{xx} \beta_*}{\sigma_y^2}$

6.3.3 - Best Prediction

- Now we will treat both x & y as random variables - no observed data yet! This is purely theoretical.

- x and y have some joint distribution (x, y)

- $x = (x_1, \dots, x_{p-1})^T$, $y = (y_1, \dots, y_p)^T$ (y can be univariate or multivariate)

- Conditional Expectation: Let $m(x) = E(y|x)$. The condit. expectation $m(x)$ has the smallest MSE of any predictor of y , ie:

$$(*) E[(y - m(x))^2] = E[(y - f(x))^2] \text{ for univariate } y. \quad f(x) \text{ is any other predictor of } y.$$

$$(**) E[(y - m(x))^T(y - m(x))] = E[(y - f(x))^T(y - f(x))] \text{ for multivariate } y.$$

Proof (univariate case):

$$E[(y - f(x))^2] = E[(y - m(x) + m(x) - f(x))^2] = E[\{(y - m(x)) + (m(x) - f(x))\}^2] = (y - m(x))^2 + (m(x) - f(x))^2$$

$$= E[(y - m(x))^2] + \underbrace{E[(m(x) - f(x))^2]}_{\text{MSE for } m(x)} + 2E[(y - m(x))(m(x) - f(x))]$$

≥ 0 since squared

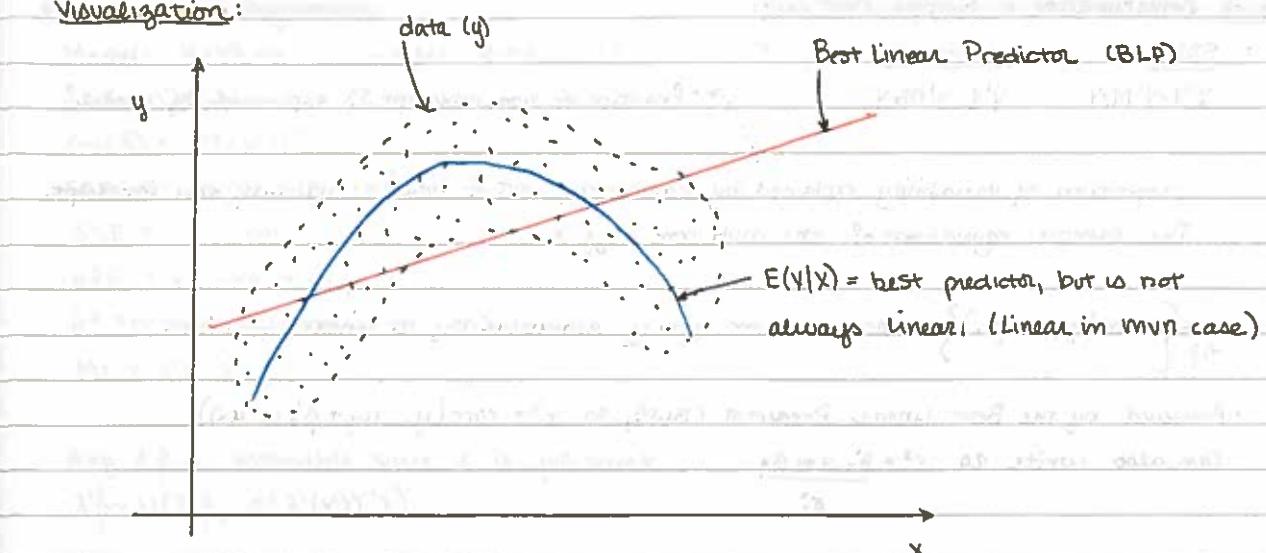
$$= 0, \text{ since } E[y - m(x)] = E(y) - E(m(x))$$

$$= E(y) - E(E(y|x))$$

$$= E(y) - E(y)$$

$$= 0$$

Visualization:



6.3.4 - Best Linear Predictor (BLP):

- Any linear predictor of y has form $f(x) = \alpha + \beta'x$.
- If y is univariate, α is scalar, $\beta + x$ $(p-1)$ length vectors.
- If y multivariate (length q), α is $(q \times 1)$, β is $(p-1 \times q)$ matrix

$$\begin{bmatrix} 1 \\ y \\ 1 \end{bmatrix}_{q \times 1} = \begin{bmatrix} 1 \\ x \\ 1 \end{bmatrix}_{q \times 1} + \begin{bmatrix} -\beta_1 \\ \vdots \\ -\beta_q \end{bmatrix}_{(q \times p-1)} \begin{bmatrix} 1 \\ x \\ 1 \end{bmatrix}_{(p-1)} = \begin{bmatrix} 1 \\ \alpha \\ 1 \end{bmatrix}_{q \times 1} + \begin{bmatrix} \beta'x \\ \vdots \\ \beta'x \end{bmatrix}_{q \times 1}$$

- $\hat{m}(x) = \hat{E}(y|x) = \mu_y + \beta'_x(x - \mu_x)$ is the best linear predictor (smallest MSE).

$$\mu_y = E(y)$$

$$\mu_x = E(x)$$

$$V_{xx} = \text{Cov}(x)$$

$$V_{xy} = \text{Cov}(x, y) = V_{yx}$$

$$\beta^* = \text{soln to } V_{xx}\beta = V_{xy}$$

↓

$$\text{Same as } \beta^* \text{ is soln to } \beta' V_{xx} = V_{yx}$$

- $\hat{E}(y|x)$ is also called 'the linear expectation'.

- If $y \sim \text{MVN}$, the BLP is the best predictor, ie $E(y|x) = \hat{E}(y|x)$.

• Proof that $\hat{E}(y|x)$ is best linear predictor:

• Goal: Show $E[(y - \hat{E}(y|x))^2] \leq E[(y - f(x))^2]$.

• Define an arbitrary linear predictor, $f(x) = \eta + \beta'(x - \mu_x)$.

$$E[(y - f(x))^2] = E[(y - \hat{E}(y|x)) + (\hat{E}(y|x) - f(x))]^2$$

$$= E[(y - \hat{E}(y|x))^2] + E[(\hat{E}(y|x) - f(x))^2] + 2E[(y - \hat{E}(y|x))(\hat{E}(y|x) - f(x))]$$

MSE for BLP ≈ 0 bc squared ⊕ = 0 as shown below

$$\oplus = E[(y - \hat{E}(y|x))(\hat{E}(y|x) - f(x))] = E[(y - \mu_y) - \beta'_x(x - \mu_x)] \{ \mu_y + \beta'_x(x - \mu_x) - \eta - \beta'(x - \mu_x) \}$$

$$= E[(y - \mu_y) - \beta'_x(x - \mu_x)] \{ (\mu_y - \eta) + (\beta_x - \beta)'(x - \mu_x) \}$$

$$= E[(y - \mu_y)(\mu_y - \eta)] - E[\beta'_x(x - \mu_x)(\mu_y - \eta)] + E[(\beta_x - \beta)'(x - \mu_x)(y - \mu_y)] - E[\beta'_x(x - \mu_x)(x - \mu_x)'(\beta_x - \beta)]$$

① ② ③ ④

$$\textcircled{1} = 0 \text{ since } E(y - \mu_y) = \mu_y - \mu_y = 0$$

$$\textcircled{2} = 0 \text{ since } E(x - \mu_x) = \mu_x - \mu_x = 0 \text{ (expectation is wrt both } x \text{ & } y)$$

$$\textcircled{3} = (\beta_x - \beta)' E((x - \mu_x)(y - \mu_y)) = (\beta_x - \beta)' V_{xy}$$

$$\textcircled{4} = \beta'_x E((x - \mu_x)(x - \mu_x)'(\beta_x - \beta)) = \beta'_x V_{xx}(\beta_x - \beta)$$

$$\text{Then } \beta_x \text{ is soln to } V_{xx}\beta = V_{xy}, \text{ so } \beta'_x V_{xx} = V_{yx}, \text{ so } \textcircled{4} = V_{yx}(\beta_x - \beta) = (\beta_x - \beta)' V_{xy}$$

Then $\textcircled{3}$ and $\textcircled{4}$ are identical, so $\oplus = 0$.

Therefore, $\hat{E}(y|x)$ has smallest MSE, so is best linear predictor. ■

Properties of Linear Expectations:

$$1. \hat{E}(Ay + bx) = A\hat{E}(y|x) + b$$

$$2. \hat{E}(x_i|x) = x_i \quad (\text{one of the coordinates of } x)$$

$$3. \hat{E}(x^T \beta | x) = x^T \beta$$

$$4. \hat{E}(y|Ax+b) = \hat{E}(y|x)$$

$$5. \hat{E}(y|x) \text{ if } \text{Cov}(x,y)=0 \text{ is my}$$

Principal Components Analysis

- PCA is the problem of finding linear combos of y , $a^T y$, that are Best Linear Predictors of the original y .
- Solution: eigenvectors of $\text{Cov}(y)$, ordered by largest eigen.

Two ways to think of PCA:

- As finding BLPA for y .
 - As finding directions a_i to project y s.t. projection onto (ca_i) has smallest MSE.
- In both cases, the principal components are the first r eigenvectors of $\Sigma = \text{Cov}(y)$, a_1, \dots, a_r , ordered by eigenvalues ϕ_1, \dots, ϕ_r s.t. $\phi_1 \geq \dots \geq \phi_r$.
 - Eigenvalues ϕ tell what % of variance is explained by each PC.
 - We center the data for PCA.

$$\text{Projection onto } ca_i \text{ for PC } a_i: \text{proj} = \frac{a_i^T y a_i}{\|a_i\|} = a_i^T y a_i$$

if $\|a_i\|=1$

Proof that eigenvals + eigenvcts are soln:

$$\min_a \left\{ E[(y - a^T y a)^T (y - a^T y a)] \right\} = \min_a \left\{ E[y^T y - 2y^T a^T y a + a^T y a^T y a] \right\} \quad (a^T a = 1)$$

$$= \min_a \left\{ E[y^T y - 2(a^T y)^2 + (a^T y)^2] \right\} = \min_a \left\{ \underbrace{E[y^T y]}_{\textcircled{1}} - \underbrace{E[(a^T y)^2]}_{\textcircled{2}} \right\}$$

$$\textcircled{1} \text{ Doesn't depend on } a, \text{ can ignore. } y^T y = \text{tr}(y^T y) = \text{tr}(yy^T) = \text{tr}(\Sigma) = \sum \phi_j \text{ for } \Sigma \text{ (sum of eigenvals)}$$

$$\textcircled{2} E[(a^T y)(a^T y)] = E(a^T y y^T a) = E(a^T \Sigma a) = a^T \Sigma a$$

Then minimizing $\text{tr}(\Sigma) - a^T \Sigma a$ is equiv to maximizing $a^T \Sigma a$.

$\max_a \{ a^T \Sigma a \}$ subj. to $\|a\| = a^T a = 1$, otherwise could just make a very large.

Write Lagrangian: $L(a) = a^T \Sigma a - \lambda(a^T a - 1)$

$$\frac{\partial L(a)}{\partial a} = 0 \rightarrow 2\Sigma a - 2\lambda a = 0 \rightarrow \Sigma a = \lambda a, \text{ so } a = \text{eigvec}, \lambda = \text{eigval}.$$

$\text{Var}(a^T y) = a^T \Sigma a = a^T \lambda a = \lambda(a^T a) = \lambda$, so eigvals explain variance.

GRUSS-MARKOV THEOREM

Example PCA Problem: $(H_{11}, \#2)$

- Let $a'_1 \dots a'_r$ be a set of r principal components for y .
- Show that $\text{tr}\{\text{Cov}[y - \hat{E}(y|a'_1 y, a'_2 y)]\} = \sum_{j=r+1}^n \phi_j$, where ϕ_j are eigenvalues.

$$\cdot \text{Cov}(y - \hat{E}(y|x)) = V_y - V_{yx}V_x^{-1}V_{xy}$$

where $x = [a'_1 y, \dots, a'_r y]^T$
($r \times 1$)

$$\cdot V_y = \text{Cov}(y) = \Sigma$$

$$\cdot V_x:$$

$$V_x = \{\text{Cov}(a'_i y, a'_j y)\} = \{\alpha_i' \Sigma \alpha_j\} = \begin{bmatrix} \alpha_1' \Sigma \alpha_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \alpha_r' \Sigma \alpha_r \end{bmatrix}$$

since $\alpha_i' \alpha_j = 0$ for $i \neq j$
bc eigenvectors orthonormal

$$\cdot V_{yx} = [\text{cov}(y, a'_1 y) \dots \text{cov}(y, a'_r y)] = [\alpha_1' \Sigma \dots \alpha_r' \Sigma]$$

Then put it together, without trace yet: $V_y - V_{yx}V_x^{-1}V_{xy}$

$$= \Sigma - [\alpha_1' \Sigma \dots \alpha_r' \Sigma] \begin{bmatrix} (\alpha_1' \Sigma \alpha_1)^{-1} & & \\ & \ddots & \\ & & (\alpha_r' \Sigma \alpha_r)^{-1} \end{bmatrix} \begin{bmatrix} \Sigma \alpha_1 \\ \vdots \\ \Sigma \alpha_r \end{bmatrix}$$

$$= \Sigma - \left(\frac{\alpha_1' \Sigma \Sigma \alpha_1}{\alpha_1' \Sigma \alpha_1} + \dots + \frac{\alpha_r' \Sigma \Sigma \alpha_r}{\alpha_r' \Sigma \alpha_r} \right)$$

Then α_i is eigenvalue for Σ , so $\Sigma \alpha_i = \phi_i \alpha_i$

$$= \Sigma - \left(\frac{\phi_1^2 \alpha_1' \alpha_1}{\phi_1 \alpha_1' \alpha_1} + \dots + \frac{\phi_r^2 \alpha_r' \alpha_r}{\phi_r \alpha_r' \alpha_r} \right) = \Sigma - (\phi_1 + \dots + \phi_r)$$

since $\alpha_i' \alpha_i = 1$
(orthonormal eigenvectors)

Now add trace:

$$\cdot \text{tr}(\Sigma - (\phi_1 + \dots + \phi_r)) = \text{tr}(\Sigma) - \sum_{j=1}^r \phi_j$$

Then $\text{tr}(\Sigma) = \text{sum of its eigenvalues}$; $\sum_{j=1}^r \phi_j$

$$\cdot \sum_{j=1}^r \phi_j - \sum_{j=r+1}^n \phi_j$$

$$\cdot \sum_{j=r+1}^n \phi_j$$

Thm: $\hat{\beta}$ is the BLUE among all linear predictors, with best = smallest variance.

Proof: Let $a'y$ be an arbitrary unbiased linear predictor.

$$\cdot E(y) = X\beta, \text{Var}(y) = \sigma^2 I, \hat{\beta} \sim N(X\beta, \sigma^2(X'X)^{-1}). \text{ Also, } \hat{\beta} = p'My \text{ for } \lambda'\beta = p'\hat{\beta}$$

$$\cdot \text{Var}(a'y) = \text{Var}(a'y - \hat{\beta} + \hat{\beta})$$

$$= \text{Var}(a'y - p'My + p'My)$$

$$= \text{Var}(p'My) + \underbrace{\text{Var}(a'y - p'My)}_{> 0 \text{ bc is a variance}} + \underbrace{2\text{Cov}(a'y - p'My, p'My)}_{= 0}$$

$$\text{Cov}(a'y - p'My, p'My) = (a' - p'M) \text{Cov}(y) M p$$

$$= \sigma^2 (a' p'M) M p$$

$$= \sigma^2 (a' M - p'M) p = 0$$

Then $E(a'y) = \lambda'\beta = p'X\beta \geq p' = a'$, so
And $E(a'y) = a'E(y) = a'X\beta \geq a' = a'$

Therefore, $\hat{\beta}$ has smallest variance among all unbiased linear estimators.

Exponential Families

1. Natural parameterization: $f(x|\theta) = h(x) \cdot c(\theta) \cdot \exp\left[\sum_{i=1}^k w_i(\theta) t_i(x)\right]$
where θ is a vec with k params.

• 1-param families are just $f(x|\theta) = h(x) \cdot c(\theta) \cdot \exp[w(\theta) t_1(x)]$

2. Natural Parameterization: $p(x|\theta) = h(x) \cdot \exp[\theta x - \psi(\theta)]$ where $E(X) = \psi'(\theta)$.

Ex 1) Is $\text{pareto}(a, b)$ with a fixed, b unknown an exponential family?

$$p(x|a, b) = \frac{ba^b}{x^{b+1}}, \quad a < x < \infty, \quad a > 0, b > 0$$

$$= ba^b x^{-(b+1)} I(x) \cdot I(a) \cdot I(b)$$

$$= ba^b I(x) \cdot I(a) \cdot I(b) e^{\log(x^{-(b+1)})} = ba^b I(x) I(a) I(b) e^{-(b+1) \log x}$$

Is an exp family, with $h(x) = I(x)$

$$w(\theta) = b+1$$

$$c(\theta) = ba^b I(a) I(b) \quad t(x) = -\log x$$

Ex 2) Is $\text{pareto}(a, b)$ with a unknown, b fixed an exponential family?

$$p(x|a, b) = ba^b x^{-(b+1)} \cdot I(x) \cdot I(a) \cdot I(b)$$

No, bc of the $I(x)$ term. Cannot separate the a and x here.

Ex 3) Is $\text{ga}(a, b)$ a 2-param exponential family?

$$f(x|a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}, \quad a > 0, b > 0$$

$$= \frac{b^a}{\Gamma(a)} \cdot I(a) \cdot I(b) \cdot e^{\log(x^{a-1} e^{-bx})}$$

$$= \frac{b^a}{\Gamma(a)} I(a) I(b) e^{(a-1) \log x + (-bx)}$$

Is a 2-param exp family, with $h(x) = 1$

$$c(\theta) = \frac{b^a}{\Gamma(a)} I(a) I(b)$$

$$w_1(\theta) = a-1$$

$$w_2(\theta) = -b$$

$$t_1(x) = \log x$$

$$t_2(x) = x$$

Ex 4) Is $X \sim \text{Binom}(n, p)$ where $n=20$ known a 1-param exp family?

$$f(x|n, p) = \binom{n}{x} p^x (1-p)^{n-x}, \quad p \in [0, 1]$$

$$\cdot \binom{n}{x} I(p) p^x (1-p)^{n-x}$$

$$\cdot \binom{n}{x} I(p) (1-p)^n e^{\log(p^x (1-p)^{-x})}$$

$$= \binom{n}{x} I(p) (1-p)^n e^{\log(\frac{p}{1-p})^x}$$

$$= \binom{n}{x} I(p) (1-p)^n e^{x \log(\frac{p}{1-p})}$$

Is 1-param exp family, with $h(x) = \binom{n}{x}$

$$c(\theta) = I(p) (1-p)^n$$

$$w(\theta) = \log(\frac{p}{1-p})$$

$$t(x) = x$$

Ex 5 (FY 15 Prelim #3):

$p(x|\theta) = h(x) \exp[\theta x - \psi(\theta)]$ is natural parameterization.

A] Let X be sample space of RV $X \sim p(x|\theta)$, from exp family w param θ .

Use identities below to show $E(x) = \psi'(\theta)$.

$$\cdot \int_x p(x|\theta) dx = 1 \quad \text{and} \quad \frac{d}{d\theta} \int_x p(x|\theta) dx = \int_x \frac{d}{d\theta} p(x|\theta) dx$$

Proof: $E(x) = \int_x x p(x|\theta) dx = \int_x x h(x) \exp[\theta x - \psi(\theta)] dx$

$$= e^{-\psi(\theta)} \int_x x \cdot h(x) e^{\theta x} dx \quad \text{Then } xe^{\theta x} = \frac{d}{d\theta} e^{\theta x}, \text{ so}$$

$$= e^{-\psi(\theta)} \int_x h(x) \cdot \frac{d}{d\theta} e^{\theta x} dx = e^{-\psi(\theta)} \frac{d}{d\theta} \int_x h(x) e^{\theta x} dx \quad \text{by 2nd identity.}$$

Then $\int_x h(x) e^{\theta x} dx$ integrates to $e^{\psi(\theta)}$ bc of 1st identity.

$$= e^{-\psi(\theta)} \frac{d}{d\theta} (e^{\psi(\theta)}) = \psi'(\theta) e^{\psi(\theta)} e^{-\psi(\theta)} = \psi'(\theta). \blacksquare$$

B] Express $N(\mu, \sigma^2)$ with unknown μ , known σ^2 in natural exponential form

$$p(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x-\mu)^2\right] = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[\frac{1}{2\sigma^2}(x^2 - 2\mu x + \mu^2)\right]$$

$$= \underbrace{\frac{e^{-\frac{x^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma}}_{h(x)} \cdot \exp\left[\frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right]$$

$\uparrow \quad \uparrow$
 $x = \mu \quad \theta = \mu$

C] Consider a Bayesian model where $x \sim N(\theta, 1)$ and $\theta \sim \pi(\theta)$. The marginal

$$m(x) = \int_{\mathbb{R}} \phi(x|\theta, 1) \pi(\theta) d\theta \quad \text{where } \phi \sim \text{normal density.}$$

Prove the posterior mean can be written as $E(\theta|x) = x + \frac{d}{dx} \log m(x)$.

$$\text{Posterior: } \frac{\phi(x|\theta, 1) \pi(\theta)}{m(x)} \cdot \frac{1}{\sqrt{2\pi}} \cdot \pi(\theta) \cdot \exp\left[-\frac{1}{2}(x-\theta)^2\right] m(x)^{-1}$$

$$= \frac{\pi(\theta)}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x-\theta)^2 - \log m(x)\right]$$

Now recall since working with posterior $p(\theta|x)$, x is like θ in formulas

$$= \underbrace{\frac{\pi(\theta)}{\sqrt{2\pi}} \exp\left[-\frac{\theta^2}{2}\right]}_{h(\theta)} \exp\left[x\theta - \underbrace{\left(\frac{x^2}{2} + \log m(x)\right)}_{\psi(x)}\right]$$

From part A, know that $E(x) = \psi'(\theta)$, so:

$$E(\theta|x) = \psi'(\theta) = \frac{2x}{2} + \frac{d}{dx} \log m(x) = x - \frac{d \log m(x)}{dx}$$

Monte-Carlo Integration Basics:

• Setting: want to evaluate integrals of the form, $I = \int g(x) f(x) dx$
where $f(x)$ is a density, $g(x)$ is some func.

• Approach: sample $x_i | A \stackrel{iid}{\sim} f(x)$, then plug $x_i | A$ into $\hat{I}_n = \frac{1}{N} \sum_{i=1}^n g(x_i)$

2 Key Results:

$$1) E(\hat{I}_n) = I, \text{ since } E(g(x)) = \int g(x) f(x) dx$$

$$2) \text{Var}(\hat{I}_n) = \frac{1}{N^2} (\text{Var}(g(x_1)) + \dots + \text{Var}(g(x_N))) = \frac{N \text{Var}(g(x))}{N^2} = \frac{\text{Var}(g(x))}{N}$$

• Convergence: $\hat{I}_n \xrightarrow{P} I$ which can be shown by Chebyshev's Inequality.

$$\cdot P(|\hat{I}_n - E(\hat{I}_n)| \geq t \cdot \text{Var}(\hat{I}_n)) \leq \frac{1}{t^2} \equiv P(|\hat{I}_n - E(\hat{I}_n)| \geq t) \leq \frac{\text{Var}(\hat{I}_n)}{t^2}$$

$$\cdot P(|\hat{I}_n - E(\hat{I}_n)| \geq t) \leq \frac{\text{Var}(\hat{I}_n)}{t^2}$$

$$= P(|\hat{I}_n - I| \geq t) \leq \frac{\text{Var}(\hat{I}_n)}{t^2}$$

$$= P(|\hat{I}_n - I| \geq t) \leq \frac{\text{Var}(g(x))}{nt^2} \quad \text{and} \quad \frac{\text{Var}(g(x))}{nt^2} \rightarrow 0 \text{ as } n \rightarrow \infty.$$

Therefore, $\hat{I}_n \xrightarrow{P} I$.

Rejection Sampling:

• Goal: Sample from some density $f(x)$.

• Steps:

i) Write $f(x) \propto l(x) \cdot h(x)$ where $l(x)$ is a non-neg fn with upper bound M .

• $h(x)$ is a density we can sample from.

ii) Check $l(x)$ is bounded by taking $\frac{d}{dx} l(x) = 0$ and solve for root(s) x^* and $M = l(x^*)$.

iii) To get $f(x)$ in form $l(x) \cdot h(x)$, can split up $f(x)$, or can multiply & divide by a density.

iv) Generate $U \sim U(0,1)$.

v) Generate h from $h(x)$.

vi) Calculate l_x as $l(h)$.

vii) Acceptance test: if $U \leq \frac{l_x}{M}$ then accept l_x as a sample from $f(x)$. Else discard & restart.

• Probability of acceptance:

If $f(x) = l(x) \cdot h(x)$, $p(\text{accept}) = 1/M$

If $f(x) \propto l(x) \cdot h(x)$, $p(\text{accept}) \approx \frac{\# \text{accepts}}{\# \text{draws}}$ (cannot calculate directly)

• Strategy: Works best if can find $h(x)$ similar in shape to $f(x)$. \rightarrow higher $p(\text{accept})$.

• Downside to approach: don't learn from previous samples; always starting over. ARS (adaptive rej sampling) addresses this.

Adaptive Rejection Sampling:

• Goal: Sample from some density $f(x)$. Must be log-concave: $\frac{d^2 \log f(x)}{dx^2} \leq 0$

• Approach: Construct a piecewise linear envelope from $\log f(x)$, and sample from envelope by exponentiating the linear parts (bc know how to sample from exponentials.) Done via rejection sampling, and each time a point is rejected, update the envelope to include an intersection of tangent lines at the sampled point.

• Steps:

i) Let $f(x)$ be the target fn. Verify it is log-concave.

• Let $h(x) = \log f(x)$

• Let $h'(x) = \frac{d}{dx} \log f(x)$

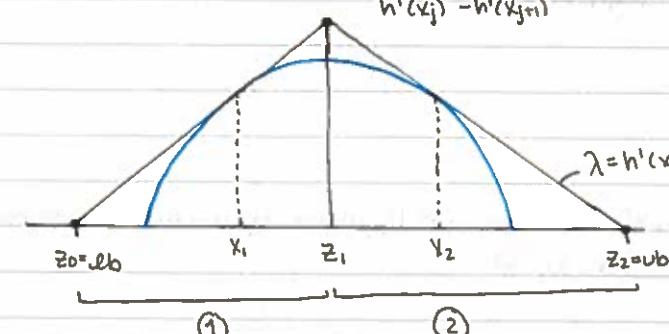
ii) Begin w. a vector of at least two x vals to try. Calculate the envelope at each point, ie the slope of the tangent and the intercept which ensures that the tan line intersects the point x_i .

• slopes: $h'(x) = m$

• intercepts: $b = h(x_j) - x_j \cdot h'(x_j)$ (bc $y = mx + b \rightarrow h(x) = h'(x) \cdot x + b$)

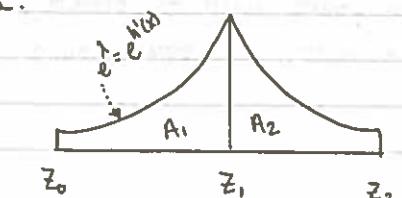
• points z_i , where tan lines meet, to divide envelope into segments:

$$z_j = \frac{h(x_{j+1}) - h(x_j) + x_{j+1} \cdot h'(x_{j+1}) - x_j \cdot h'(x_j)}{h'(x_j) - h'(x_{j+1})}$$



$\text{---} = f(x)$, our target

iii) Exponentiate the piecewise components to get a set of (rotated) truncated exponentials.



4) Calculate the area of each exponential. This is used for the first step of sampling from the exponential - must first select a segment to use, with that sampling weighted by area under each segment.

$$A_i = \frac{\exp[h(x_i) - h'(x_i)]}{h'(x_i)} \left(\exp[h'(x_i)z_i] - \exp[h'(x_i)z_i] \right)$$

which is $\int_{z_i}^{z_{i+1}} \exp[h'(x_i)x + (h(x_i) - h'(x_i)x_i)] dx$
tan line fact

5) Sample to pick which exponential segment to use. Chooses an "area".

• Selected area = sample(1: length(Areas), size=1, replace=T, prob = Areas / sum(Areas))

• gives a selected area weighted by sample size.

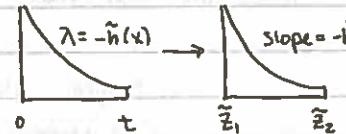
6) Draw an exponential sample (size1) from the selected exponential segment.

• Let $\tilde{h}(x)$ be the slope of the chosen segment.

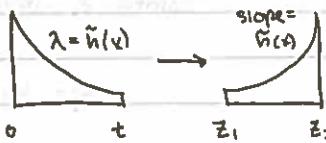
• Let \tilde{z}_1, \tilde{z}_2 be the z 's bounding the segment.

• Sample \tilde{x} from the truncated exponential. $\text{Exp}(\lambda)$ must have positive rate λ , so:

• If $\tilde{h}(x) < 0$, slope neg. Facing right direction, so use $r + \exp(n=1, m=-\tilde{h}(x), t=z_2-z_1) = \tilde{x}$
then have shifted result $x^* = z_1 + \tilde{x}$.



• If $\tilde{h}(x) > 0$, slope pos. Facing wrong direction, so use $r + \exp(n=1, m=\tilde{h}(x), t=z_2-z_1) = \tilde{x}$
to sample, then shift/flip, so result is $x^* = z_2 - \tilde{x}$.



7) Acceptance Test:

• Draw u from $U(0,1)$.

• Calc proposal: $g(x^*) = \exp\{h'(x_i)x^* + [h(x_i) - h'(x_i)x_i]\}$, value of tan line for selected segment, evaluated at x^* .

• If $u \leq f(x^*)/g(x^*)$, accept x^* as being from $f(x)$. Repeat process w/ same envelope.

• If not, reject x^* . Also, update vector of x points used to create piecewise linear envelope to include x^* . Sort the x vector, and start process again with new x vector.

Importance Sampling

• Goal: Estimate the integral $I = \int g(x) dx$

• Approach: Write I as $I = \int \frac{g(x)}{f(x)} \cdot f(x) dx$ where $f(x)$ is a density we can sample.

Then use Monte Carlo as usual to approximate the integral.

Steps:

i) Rewrite I as $I = \int \frac{g(x)}{f(x)} \cdot f(x) dx$

ii) Draw $x_i \stackrel{iid}{\sim} f(x)$

iii) Estimate $\hat{I}_N = \frac{1}{N} \sum_{i=1}^N \frac{g(x_i)}{f(x_i)}$ where $x_i \sim f(x)$. (Even though $x_i \sim f(x)$, still plug into the whole thing, $g(x)/f(x)$.)

• Key: Now we are introducing density $f(x)$ to help us, instead of beginning with some integral $I = \int g(x)f(x)dx$ which already contains $f(x)$.

Normalized Importance Sampling

• Want to estimate integral $I = E_f(g(x)) = \int g(x) f(x) dx$, but we only know $f(x)$ up to a constant of proportionality, call this $f^*(x)$.

• Can't use regular Monte Carlo integration, bc need the constant included for integration.

• Can use Monte Carlo with a modification: $f(x) = \frac{f^*(x)}{\int f^*(x) dx} \leftarrow \text{normalizing constant}$

i) Estimate normalizing constant:

$$\hat{I}_N^* = \frac{1}{N} \sum_{i=1}^N \frac{f^*(x_i)}{h(x_i)} \cdot h(x_i), \quad x_i \sim h(x)$$

• Using importance sampling, where $h(x)$ is a density we can sample.

ii) Estimate whole integral:

$$\hat{I}_N = \left[\frac{1}{N} \sum_{i=1}^N \frac{g(x_i)}{h(x_i)} \cdot f^*(x_i) \right] \frac{1}{\hat{I}_N^*}$$

using $x_i \stackrel{iid}{\sim} h(x)$.

• Can reuse same x 's; don't need to generate a new sample.

• Plug x 's from (i) into $\frac{g(x)}{h(x)} \cdot f^*(x)$

Ratio of Uniforms:

- Another way to sample from a density $f(x)$.

Algorithm:

- Calculate $a = b$.
 $a = \max_x \sqrt{f(x)}$
 $b = \max_x x\sqrt{f(x)}$
- \max_x roots of x which maximize $\sqrt{f(x)}$ and $x\sqrt{f(x)}$.

- Sample $u_1 \sim \text{unif}(0, a)$
- $u_2 \sim \text{unif}(0, b)$

- Accept $x = u_2/u_1$ as being from $f(x)$ if $u_1 \leq \sqrt{f(u_2/u_1)}$, else reject.

Why it works:

- Introducing a joint density to help sample $f(x)$: $f(x,y) \propto y \cdot \mathbb{1}(y < \sqrt{f(x)})$

$\int y dy = \frac{1}{2} f(x)$, so $\int f(x,y) dy = f(x)$.

- Bivariate transformation: $u_1 = y_1$, $u_2 = yx \rightarrow y = u_1$, $x = u_2/u_1$

- Then $f(u_1, u_2) \propto u_1 \cdot \mathbb{1}(0 < u_1 < \sqrt{f(u_2/u_1)})$ where $J = \begin{vmatrix} -u_2/u_1 & u_1 \\ 1 & 0 \end{vmatrix} = 1/u_1$

- To sample from the region, can sample from its bounding rectangle, then accept when (u_1, u_2) samples land in S.



- How to find the bounds of rectangle R?

- Range for u_1 is 0 to $\max_x \sqrt{f(x)}$ → sample u_1 from $u(0, \max_x \sqrt{f(x)})$

- For u_2 , if $x \geq 0$, then $u_1 \leq \max_x \sqrt{f(x)}$ and $y_1 = u_1, y_2 = u_2 \rightarrow u_2 \leq \max_x x\sqrt{f(x)}$

- If x can take values < 0 , just look at u_2^2 , and you get +/- same bound:

$$\min_x x\sqrt{f(x)} \leq u_2 \leq \max_x x\sqrt{f(x)}$$

Basic Markov Chain Theory

- Goal: Generate a sample from target density $f(x)$.

- Each element in a Markov chain depends only on previous element.

- Need a transition density $p(x_{n+1}|x_n)$ which satisfies the stationary condition:

$$f(x_{n+1}) = \int p(x_{n+1}|x_n) f(x_n) dx_n$$

- The stationary condition preserves $f(x)$ as the marginal density.

- Can also sample x_i 's then plug into $\hat{I}_N = \frac{1}{N} \sum_{i=1}^N g(x_i)$ and still converges, since x_i 's are marginally from $f(x)$.

- Catch: $E(\hat{I}_N) = I$ still, but $\text{Var}(\hat{I}_N)$ increases. This is b/c $\text{Cov}(x_{n+1}, x_n) > 0$.

- $\text{Var}(\hat{I}_N) = \sqrt{n}(\hat{I}_N - I) \rightarrow N(0, \sigma^2)$, $\sigma^2 = \text{Var}(g(x_i)) + 2 \sum_{i=1}^N \text{Cov}(g(x_i), g(x_{i+1}))$

- Reversibility: If $p(x'|x) f(x) = p(x|x') f(x')$ then chain is stationary.

Proof: Let $p(x'|x) f(x) = p(x|x') f(x)$.

Then $\int p(x'|x) f(x) dx = \int p(x|x') f(x') dx = f(x')$ so stationary cond is met.
density wrt x' , integrates to 1

Finding a transition density:

- want to find a transition → evaluate its covariance, $\text{Cov}(x_{n+1}, x_n)$.
- If we have a joint density, the product of the full conditionals is the transition.
- If don't have a joint, need to introduce one.

- If x is in discrete space, can directly form transition matrix P .

- If x continuous, can introduce joint via indicator funcs.

Ex 1: X in discrete space (HW 2, #2)

- The joint density $f(x, q) = q^x(1-q)^{1-x}$, $0 \leq q \leq 1$, $x \in \{0, 1\}$. $P(x'|x) = 2 \int_0^{x+1} q^{x+1} (1-q)^{2-x-x'} dq$.
- Find the stationary density.

Since X only takes values 0 or 1, evaluate $p(x'|x)$ at each combo of x, x' . (Plug each combo into $p(x'|x)$ & eval integral.)

$$\begin{array}{ll} p(x'=0|x=0) = 2/3 & p(x'=0|x=1) = 1/3 \\ p(x'=1|x=0) = 1/3 & p(x'=1|x=1) = 2/3 \end{array} \Rightarrow P = \begin{bmatrix} 2/3 & 1/3 \\ 1/3 & 2/3 \end{bmatrix}$$

We know $\pi^T = \pi^T P$ since want a stationary density.

We know $\pi = [\pi_0 \ \pi_1]$ where $\pi_0 + \pi_1 = 1$, since π_0 is prob for x_0 , π_1 is prob for x_1 .

↓

These give us a system of 2 eqns to solve:

$$(1) [\pi_1 \ \pi_0] P = [\pi_1 \ \pi_0]$$

The result of solving these gives $\pi_1 = \pi_0 = 1/2$, so

$$(2) \pi_1 + \pi_0 = 1$$

stationary is discrete Unif(1/2) on $\{0, 1\} = X$.

Ex 2: X in continuous space, no joint to start: (Ex 6, prob 2)

Find a transition which has $\pi(x) = e^{-x}$, $x > 0$ as the stationary.

Need to introduce a joint $\pi(u, x)$ s.t. $p(x'|x) = \int \pi(x'|u) \pi(u|x) dx$ where $\pi(x|u) + \pi(u|x)$ are full conditionals of $\pi(u, x)$.

Try $\pi(u, x) = \mathbb{1}(u < f(x))$, where $f(x) = e^{-x}$. Find full conditionals:

$$\pi(u|x) = U(0, e^{-x})$$

$$\pi(x|u) \text{ is unif on } x: f(x) > u \rightarrow e^{-x} > u \rightarrow -x > \log u \rightarrow x < -\log u$$

so $\pi(x|u) = U(0, -\log u)$.

Begin w/ any x , then iterate sampling u from $U(0, e^{-x})$
 x from $U(0, -\log u)$

Can also find the stationary which gives update $x' = \phi(x, u, v)$ where $u, v \stackrel{iid}{\sim} U(0, 1)$.

$$u \sim U(0, e^{-x}) \rightarrow u = e^{-x}v \quad (\text{just rescaling } v \sim U(0, 1)).$$

$$x \sim U(0, -\log u) \rightarrow u(0, -\log(e^{-x}v)) = U(0, x - \log v)$$

↓

$$x' = u \cdot (x - \log v) \text{ with } u, v \stackrel{iid}{\sim} U(0, 1)$$

why? ease of sampling, & can make finding $\text{cov}(x', x)$.

Covariance of Markov Chains

- For a stationary density, first find the transition, then find $\text{cov}(x_{n+1}, y_n) \approx \text{cov}(x', x)$.
- Two helpful tricks for finding covariance:

i) Find transition in terms of $x' = \phi(x, u, v)$ where $u, v \stackrel{iid}{\sim} U(0, 1)$

ii) Iterated expectation: $\text{cov}(x', x) = E(x'x) - E(x')E(x) = E(x \cdot E(x'|x)) - E(x) \cdot E(E(x'|x))$

Ex: Exercise 6, Prob 1:

- Find a transition which has $\pi(x) \sim \text{Ga}(x|a, 1)$ as stationary, s.t. $\pi(x) = \frac{x^{a-1} e^{-x}}{\Gamma(a)}$, $x > 0, a > 0$. Find $\text{cov}(x', x)$.

1) Introduce joint density & find transition:

Need a joint $\pi(u, x)$ so can use full conditionals $\pi(u|x)$, $\pi(x|u)$.

i) $(u < \frac{x^{a-1} e^{-x}}{\Gamma(a)})$ not easily invertible. Split into product of 2 facts & take intersection.

$$\pi(u, x) = \underbrace{\mathbb{1}\left(u < \frac{x^{a-1}}{\Gamma(a)}\right)}_{\pi_1(u, x)} \cdot \underbrace{\mathbb{1}(u < e^{-x})}_{\pi_2(u, x)}$$

Then need two pairs of full conditionals.

$$\text{a) For } \pi_1(u, x): \pi_1(u|x) = U(0, \frac{x^{a-1}}{\Gamma(a)}). \rightarrow u < \frac{x^{a-1}}{\Gamma(a)} \rightarrow \log(u/\Gamma(a)) < (a-1)\log x$$

$$\rightarrow \exp\left[\frac{\log(u/\Gamma(a))}{a-1}\right] < x$$

$$\pi_1(x|u) = U\left(\exp\left[\frac{\log u/\Gamma(a)}{a-1}\right], \infty\right)$$

$$\text{b) For } \pi_2(u, x): \pi_2(u|x) = U(0, e^{-x}) \rightarrow u < e^{-x} \rightarrow x < -\log u$$

$$\pi_2(x|u) = U(0, -\log u)$$

Can sample $u_1, u_2 \stackrel{iid}{\sim} U(0, 1)$ for the two bounds, then combine the above to get

$$x' \sim U\left(\exp\left[\frac{\log u_1/\Gamma(a)}{a-1}\right], -\log u_2\right)$$

- Put above result in form $x' = \phi(x, u, v)$ for easy covariance calcs. In this case, we will need a 3rd uniform:

$$x' = \phi(x, u, y, w) \text{ where } u, y, w \stackrel{iid}{\sim} U(0, 1)$$

$$\cdot \pi_1(u(x)) \rightarrow u_1 = v \frac{x^{a-1}}{r(a)} \quad \left. \right\} \text{ for } v, w \text{ s.t. } v(0,1)$$

$$\pi_2(u|x) \rightarrow u_2 = w e'$$

• Plug $u_1 + u_2$ into x bounds: $\pi(x|u_1, u_2) \sim U\left(\exp\left[\frac{\log u_1 r(a)}{a-1}\right], -\log u_2\right) \rightarrow U\left[x e^{\log r(a)/a-1}, x - \log u_2\right]$

• Add $\gamma \sim U(0, i)$. when $\gamma = 0$, want lower bound for x^* , ie $x^* \geq \log(a-1)$.

When $\gamma=1$, want upper bound for x^t , ie $x - \log w$.

$$x' = x e^{\log \sqrt{a-1}} + y [(x - \log w) - x e^{\log \sqrt{a-1}}] \quad \text{where } u, w, y \stackrel{iid}{\sim} U(0, 1)$$

3) Find Cov

$$\text{Cov}(X', X) = E(X'X) - E(X)E(X') = E(X \underbrace{E(X'|X)}_{= E(X)} - E(X) \cdot E(\underbrace{E(X'|X)}_{= E(X)}))$$

$$\bullet E(Y'|X) = E(X^{\gamma/a-1} + \gamma X - \gamma \log w - \gamma X^{\gamma/a-1} | X) \quad \text{by plugging in } x' \text{ expression}$$

$$= \underbrace{X E(V^{1/a-1})}_{\int_0^1 V^{1/a-1} (1) dV = \frac{a-1}{a}} + \underbrace{X E(Y)}_{=1/2} - E(Y) E(\log W) - \underbrace{X E(Y) E(V^{1/a-1})}_{=1/2} \quad \text{and recall } Y, V, W \stackrel{iid}{\sim} U(0,1)$$

$$= x \left(\frac{a-1}{a} \right) + \frac{x}{2} - \frac{1}{2}(-1) - x \left(\frac{1}{2} \right) \left(\frac{a-1}{a} \right) = \boxed{\frac{1}{2}x \left(\frac{a-1}{a} + 1 \right) + \frac{1}{2}}$$

$$\text{Cov}(x^i, x) = E\left[x \left(\frac{1}{2} \times \left(\frac{a-1}{a} + 1\right) + \frac{1}{2}\right)\right] - E\left[\frac{1}{2} \times \left(\frac{a-1}{a} + 1\right) + \frac{1}{2}\right] E(x)$$

$$= E \left[\frac{1}{2} X^2 \left(\frac{a-1}{a} + 1 \right) + \frac{X}{2} \right] - \left[\frac{1}{2} X \left(\frac{a-1}{a} + 1 \right) + \frac{1}{2} \right] E(X)$$

$$= \frac{a-1}{2a} E(X^2) + \frac{1}{2} E(X^2) + \frac{1}{2} E(X) - \frac{a-1}{2a} E(X)^2 - \frac{1}{2} E(X)^2 - \frac{1}{2} E(X)$$

$$= \left(1 - \frac{1}{2}a\right) \ln \psi(x)$$

- What we really want is AR1, i.e. autocorrelation:

$$p = \frac{\text{Cov}(x^1, x)}{\text{sd}(x^1) \text{sd}(x)} = \frac{(1 - \frac{1}{2}a) \text{Var}(x)}{[\text{sd}(x)]^2} = \frac{(1 - \frac{1}{2}a)}{\text{Var}(x)}$$

Metropolis-Hastings:

goal: Sample a chain from density $\pi(x)$

- Algorithm: • Draw x^* from proposal $q(x'|x)$.

• Draw u from $u(0,1)$.

If $u < \alpha(x^*|x)$, then $x' = x^*$, else $x' = x$.

$$\cdot d(x'|x) = \min \left\{ \frac{\pi(x')}{\pi(x)}, \frac{q(x|x')}{q(x'|x)}, 1 \right\}$$

$$p(x^*|x) = \alpha(x^*|x) q(x^*|x) + (1 - r(x)) \mathbb{1}(x^* = x); \quad r(x) = \int \alpha(x^*|x) q(x^*|x) dx^*$$

Proof: Show the stationary condition holds for target $\pi(x)$. ie sufficient to show reversability, ie $p(x'|x)\pi(x) = p(x|x')\pi(x')$

$$p(x'|x) \pi(x) = \pi(x) \cdot \left[r(x|x) q(x^*|x) + (1 - r(x)) \mathbb{1}(x' = x) \right]$$

$$= \pi(x) \left[\min \left\{ 1, \frac{\pi(x')}{\pi(x)} \cdot \frac{g(x|x')}{g(x'|x)} \right\} g(x'|x) + (1 - r(x')) \mathbf{1}(x = x') \right]$$

Multiply $\pi(x)$ through, or multiply $q(x'|x)$ into the min fn.

$$= \min \left\{ \pi(x) q(x'|x), \frac{\pi(x') \pi(x)}{\pi(x)} \frac{q(x|x') q(x'|x)}{q(x'|x)} \right\} + \underbrace{\pi(x) (1 - r(x')) \mathbf{1}(x' = x)}$$

$$= \pi(x^*) \left[\min \left\{ \frac{\pi(x) q(x^*|x)}{\pi(x^*) q(x^*|x)}, 1 \right\}^{q(x|x^*)} + (1 - r(x^*)) \mathbb{1}(x^* = x) \right]$$

by pulling $\pi(x')$ out and pulling $q(x|x')$ out of min fctn

$\pi(x^t) p(x|x^t)$. Therefore, reversability holds, so stationary holds. ■

How to pick proposal g ?

- Want density w. mass in similar areas as $\pi(x)$, w. same support.
- Center g around previous x value.
- If variance of g is too small, chain will get stuck around same x value, ie high acceptance, slow exploration of the space.
- If variance of g is too large, chain will bounce around, ie fast exploration of space, low acceptance.
- Plot the traces of chains to evaluate how proposal choice is mixing.

Example: (midterm, #3)

- Want to sample from the power law dist, $p(x=x) \propto x^{-\alpha}$, for $x \in \{1, 2, 3, \dots\}$, $\alpha > 1$.

Will use MH with proposal $g(x'|x) = \begin{cases} x' = x+1 & \text{with prob } 1/2 \\ x' = x-1 & \text{with prob } 1/2 \end{cases}$ and $g(2|1) = 1$.

Write the details of implementing MH.

First, define $a(x'|x)$:

$$\text{For } x > 1, a(x'|x) = \min \left\{ 1, \frac{f(x)g(x|x')}{f(x')g(x'|x)} \right\} = \min \left\{ \frac{x^{\alpha-1}(1/2)}{(x')^{\alpha-1}(1/2)}, 1 \right\} = \min \left\{ \frac{x^{\alpha-1}}{(x')^{\alpha-1}}, 1 \right\}$$

$$\text{For } x=1, a(x'|x) = \min \left\{ \frac{x^{\alpha-1}(0)}{(x')^{\alpha-1}(1)}, 1 \right\} = 0$$

Algorithm:

- Draw x^* from proposal.
 - If $x > 1$, $x^* = \text{Binom}(1, 1/2)$ outcome for $x^* = x+1$ or $x^* = x-1$.
 - If $x=1$, $x^* = 2$.
- Draw $u \sim U(0, 1)$
- Calc $a(x^*, x)$ based on above definition.
- If $u < a(x^*, x)$ then $x' = x^*$, else $x' = x$.

Note: Since $a=0$ if $x=1$, we will always accept $x'=x^*$ for the move to 2. ■

Gibbs Sampler

- MH is difficult for hi-dim proposals. The Gibbs Sampler is an alternative. Great for generating samples from a joint posterior, or any joint density.

- If $p(u, \lambda | x)$ is a joint density.

- Begin with some arbitrary u, λ .
- Sample u' from full conditional $\pi(u | \lambda, x)$
- Sample λ' from full conditional $\pi(\lambda | u', x)$
- Continue iterating. (Order doesn't matter)

Key: product of the iteratively sampled full conditionals is the transition density:

$$p(u', \lambda' | x, u, \lambda) = \pi(u' | \lambda', x) \cdot \pi(\lambda' | u, x)$$

Proof: Prove the stationary condition holds for Gibbs. Say stationary density is $p(u, \lambda | x)$. Transition is $p(u' | \lambda', x) p(\lambda' | u, x)$.

$$\begin{aligned} & \iint p(u' | \lambda', x) p(\lambda' | u, x) p(u, \lambda | x) dud\lambda \\ &= \int p(\lambda' | u, x) p(u, \lambda | x) du = \int \underbrace{p(\lambda', u, \lambda | x)}_{\text{density w.r.t. } u} du = p(\lambda', \lambda | x) \text{ Then plug in.} \\ & \int p(u' | \lambda', x) p(\lambda' | u, x) d\lambda = \int \underbrace{p(u', \lambda', \lambda | x)}_{\text{density w.r.t. } \lambda} d\lambda = p(u', \lambda' | x). \blacksquare \end{aligned}$$

* Notes:

- In univariate case, we were introducing "u" placeholders to create a joint. No need here.
- The $|x$ can be removed and all results would still hold.

Gibbs Sampler With MH Step:

- If one (or more) of the full conditionals is difficult to sample from, can introduce a MH step for that full conditional within the Gibbs Sampler.

Gibbs Sampler Proof for 3 variables.

- Stationary is $p(\mu, \lambda, \theta | x)$
- Update order: λ , then θ , then μ . (Doesn't matter; can do in any order.)
- Transition: $p(\mu' | \theta', \lambda', x) \cdot p(\theta' | \lambda', \mu, x) \cdot p(\lambda' | \mu, \theta, x)$

Proof that the stationary condition holds: (Begin w/ RHS)

$$\iiint p(\mu' | \theta', \lambda', x) p(\theta' | \lambda', \mu, x) p(\lambda' | \mu, \theta, x) p(\mu, \theta, \lambda | x) d\lambda d\theta d\mu$$

$$= \underbrace{\int p(\lambda', \mu, \theta, \lambda | x) d\lambda}_{\text{density wrt } \lambda} = p(\lambda', \mu, \theta | x). \text{ Plug back in:}$$

$$= \iint p(\mu' | \theta', \lambda', x) p(\theta' | \lambda', \mu, x) p(\lambda' | \mu, \theta | x) d\theta d\mu$$

$$= \underbrace{\int p(\theta', \lambda', \mu, \theta | x) d\theta}_{\text{density wrt } \theta} = p(\theta', \lambda', \mu | x). \text{ Plug back in:}$$

$$= \int p(\mu' | \theta', \lambda', x) p(\theta', \lambda', \mu | x) d\mu$$

$$= \underbrace{\int p(\mu' | \theta', \lambda', \mu | x) d\mu}_{\text{density wrt } \mu} = p(\mu', \theta', \lambda' | x). \text{ Stationary condition holds.}$$

2-Component Gaussian Mixture Model:

- Model: $g(x|\theta) = w N(x|\mu_1, \sigma_1^2) + (1-w) N(x|\mu_2, \sigma_2^2)$, $\theta = \{w, \mu_1, \mu_2, \sigma_1^2, \sigma_2^2\}$

- Problem: the likelihood is difficult to use directly bc it has a sum:

$$L(\theta | x) = \prod_{i=1}^n [w N(x_i | \mu_1, \sigma_1^2) + (1-w) N(x_i | \mu_2, \sigma_2^2)]$$

- Solution: Introduce one or more latent variables to get rid of the sum in the likelihood. Only need 1 latent var in this case.

$d_i \in \{1, 2\}$ to tell us which component ea. obs x_i belongs to.

- Can then rewrite $g(x|\theta)$ as $g(x, d | \theta) = w_d N(x | \mu_d, \sigma_d^2)$ where $w_d \Rightarrow w_1 = w, w_2 = (1-w)$.

- Lhood becomes $\prod_{d=1}^D N(x_i | \mu_d, \sigma_d^2) \cdot \prod_{d=2}^M N(x_i | \mu_d, \sigma_d^2)$, ie $\prod_{i=1}^n w_{d_i} N(x_i | \mu_{d_i}, \sigma_{d_i}^2)$

- Key: choose latent so when you integrate it out, you get original lhood.

General M-Component Gaussian Mixture Model:

- Model: $g(x|\theta) = \sum_{j=1}^M w_j N(x|\mu_j, \sigma_j^2)$ where $\sum_j w_j = 1$, $\theta = \{w_1, \dots, w_M, \mu_1, \dots, \mu_M, \sigma_1^2, \dots, \sigma_M^2\}$

- Priors of μ_j 's are iid Normals $N(m, \tau^2)$

- Prior for $\sigma^2 = IG(a, b)$ (Jeffreys, $\frac{1}{\sigma^2}$, would work also)

- Prior for $(w_1, \dots, w_M) \sim \text{Dirichlet}$:

$$p(w) \propto \prod_{j=1}^M w_j^{d_j-1} \text{ with } \sum_{j=1}^M w_j = 1 \rightarrow p(w) \propto \prod_{j=1}^{M-1} w_j^{d_j-1} (1-w_{M-1} - \dots - w_1)^{d_M-1}$$

- If $M=2$, this is Beta. Can use Beta(1,1).

- Again introduce latent d_i 's to identify class membership for each obs.

⊗

Mixture Model Example ($\mathbf{x} \in \mathbb{R}^n$, #1)

Model: $g(\mathbf{x}|w) = w N(\mathbf{x}|\mu_1, \sigma^2) + (1-w) N(\mathbf{x}|\mu_2, \sigma^2)$.

$p(\sigma^2) \sim IG(1, 1)$

$p(\mu_1), p(\mu_2) \sim N(0, v)$ ($v=0$)

$p(w) \sim Beta(1, 1)$.

Use Gibbs Sampler to generate samples from $\Theta = \{\mu_1, \mu_2, \sigma^2, w\}$. Note: latent variable must be included in the sampler.

Let $d = \{1, 2\}$ with probs $p(d=1) = w$ and $p(d=2) = (1-w)$.

Lhood: $\prod_{i=1}^{n_1} w d_i N(x_i | \mu_1, \sigma^2) \cdot \prod_{i=1}^{n_2} (1-w) d_i N(x_i | \mu_2, \sigma^2)$
 for $d_i = 1$ obs for $d_i = 2$ obs

Full conditionals:

$$1) (\mu_1 | \dots) \propto \exp\left[-\frac{1}{2v} \mu_1^2\right] \cdot \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^{n_1} (x_i - \mu_1)^2\right] \text{ using only obs in } d_1 \text{ group}$$

$$\propto \exp\left[-\frac{1}{2v} \mu_1^2 - \frac{1}{2\sigma^2} (n_1 \mu_1^2 - 2\mu_1 \sum_{i=1}^{n_1} x_i)\right]$$

$$\propto \exp\left[-\frac{1}{2v} \mu_1^2 - \frac{n_1}{2\sigma^2} (\mu_1^2 - 2\mu_1 \frac{1}{n_1} \sum_{i=1}^{n_1} x_i)\right]$$

\sim Normal with $\text{var} = \left[\frac{1}{v} + \frac{n_1}{\sigma^2}\right]^{-1}$, mean = $\text{Var} \cdot \left[\frac{n_1}{\sigma^2} \bar{x}_1\right]$

2) ($\mu_2 | \dots$) is same as μ_1 , but using only obs in d_2 group, and w .

$$3) (w | \dots) \propto \underbrace{\prod_{i=1}^{n_1} w}_{\text{lhood}} \cdot \underbrace{\prod_{i=1}^{n_2} (1-w)}_{\text{Beta}(1,1) \text{ prior}} \cdot w^{n_1} (1-w)^{n_2} \sim Beta(n_1+1, n_2+1)$$

$$4) (\sigma^2 | \dots) \propto (\sigma^2)^{-n_1/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^{n_1} (x_i - \mu_1)^2\right] \cdot \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^{n_2} (x_i - \mu_2)^2\right] (\sigma^2)^{-a-1} \exp[-b/\sigma^2]$$

$$\text{for } d_i = 1 \text{ obs} \quad \text{for } d_i = 2 \text{ obs}$$

$$\sim IG\left(\frac{n_1}{2} + a, b + \frac{1}{2} \sum_{i=1}^{n_1} (x_i - \mu_1)^2 + \frac{1}{2} \sum_{i=1}^{n_2} (x_i - \mu_2)^2\right)$$

$$5) p(d_i = 1) = \frac{w \cdot N(x_i | \mu_1, \sigma^2)}{w \cdot N(x_i | \mu_1, \sigma^2) + (1-w) N(x_i | \mu_2, \sigma^2)}$$

Modeling 1

Convergence in Probability:

An estimator $\hat{\theta}$ is consistent if $\hat{\theta} \xrightarrow{P} \theta$.

$$\hat{\theta} \xrightarrow{P} \theta \rightarrow E(\hat{\theta}) \xrightarrow{P} E(\theta)$$

Definition: $X_n \xrightarrow{P} X$ iff $\lim_{n \rightarrow \infty} P(|X_n - X| \geq \epsilon) = 0$

Tools:

1. Markov's Inequality: $P(g(x) \geq r) \leq \frac{E(g(x))}{r}$

2. Chebyshev's Inequality: $P(|X - E(X)| \geq t \text{Var}(X)) \leq \frac{1}{t^2}$ ie $P(|X - E(X)| \geq t) \leq \frac{\text{Var}(X)}{t^2}$
(x is some parametrl, ie $\hat{\theta}, S, \bar{X}$, etc)

3. WLLN: $\bar{X} \xrightarrow{P} \mu$. Also includes moments: $\frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{P} \mu_k$

4. Continuous Mapping Theorem: If $g(x)$ continuous and $X_n \xrightarrow{P} X$, then $g(X_n) \xrightarrow{P} g(X)$

Also works w. 2+ variables: Ex, if $g(a, b) = a/b$ and $X_n \xrightarrow{P} X$, $Y_n \xrightarrow{P} Y$,
then $X_n/Y_n \xrightarrow{P} X/Y$.

Ex: Generate $(z_1, x_1), \dots, (z_n, x_n)$ as $z_i = \begin{cases} \text{blue w. prob } \theta \\ \text{red w. prob } 1-\theta \end{cases}$ $x_i = \begin{cases} \text{Pois}(1) \text{ if } z_i \text{ blue} \\ \text{Pois}(4) \text{ if } z_i \text{ red} \end{cases}$

Given only the x_i 's, find a consistent estimator of θ .

$$E(x_i) = E(E(x_i | z_i)) = E(1\theta + 4(1-\theta)) = 4-3\theta \rightarrow E(x_i) = 4-3\hat{\theta} \rightarrow \hat{\theta} = \frac{4-\bar{x}}{3}$$

Consistent b/c by WLLN, $\bar{x} \xrightarrow{P} E(x_i)$, and:

$$\lim_{n \rightarrow \infty} P(|\hat{\theta} - \theta| \geq \epsilon) = \lim_{n \rightarrow \infty} P\left(\left|\frac{4-\bar{x}}{3} - \frac{4-E(x_i)}{3}\right| \geq \epsilon\right) = \lim_{n \rightarrow \infty} P(|\bar{x} - E(x_i)| \geq \epsilon) = 0$$

by WLLN & continuous mapping thrm.

Consistency \neq Unbiased:

consistency means $\hat{\theta} \xrightarrow{P} \theta$ as $n \rightarrow \infty$. True if sampling dist of $\hat{\theta}$'s var $\rightarrow 0$, ie becomes
↑ concentrated as $n \rightarrow \infty$.

unbiased applies to any sample size. 'On avg, $\hat{\theta}$ hits true θ '.

$\hat{\sigma}^2 = \frac{1}{n-1} \sum (x_i - \bar{x})^2$ is biased but consistent.

Convergence in Distribution:

Definition: $X_n \xrightarrow{d} X$ if $\lim_{n \rightarrow \infty} F_n(t) = F(t)$, ie if $\lim_{n \rightarrow \infty} P(X_n \leq t) = P(X \leq t)$

Tools:

1. Convergence in prob \rightarrow convergence in dist. (Converse false)

2. Slutsky's Lemma: If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{P} c$, a constant, then
 $X_n Y_n \xrightarrow{d} Xc$
 $X_n + Y_n \xrightarrow{d} X + c$

3. CLT: If $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ then $\bar{X} \xrightarrow{d} N(\mu, \sigma^2/n)$

4. Continuous Mapping Thm: If $g(x)$ continuous and $X_n \xrightarrow{d} X$, then $g(X_n) \xrightarrow{d} g(X)$.

5. Delta Method: If $\frac{\sqrt{n}(Y_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$ and $g(x)$ is a fctn where $g'(x) \neq 0$, then

$$g(X_n) \xrightarrow{d} N(g(\mu), \frac{\sigma^2}{n} g'(\mu)^2)$$

How to begin: $P(X_n \leq t)$ and begin rearranging.

Ex: Let $X_1, \dots, X_n \sim U(0, 1)$. Show that $n(\hat{\theta} - \theta) \xrightarrow{d} \text{Exp}(\lambda)$, where $\hat{\theta}$ is MLE of θ . Find λ .

MLE: $\prod_{i=1}^n \frac{1}{1-\theta} = \left(\frac{1}{1-\theta}\right)^n$ is maximized for sample min, $X_{(1)}$. So $\hat{\theta}_{MLE} = X_{(1)}$.

$$P(n(\hat{\theta} - \theta) \leq t) = P(\hat{\theta} - \theta \leq t/n) = P(\min(X_i) \leq \theta + t/n) = 1 - P(\text{all } X_1, \dots, X_n \geq \theta + t/n)$$

$$= 1 - P(X_1 \geq \theta + t/n) \cdots P(X_n \geq \theta + t/n) \text{ since } X_i \text{'s iid} = 1 - P(X \geq \theta + t/n)^n \text{ since } X \text{'s iid}$$

$$= 1 - [1 - P(X \geq \theta + t/n)]^n \text{ where } X \sim U(0, 1) \text{ so cdf is } \frac{x-a}{b-a}, X = \theta + t/n, a = \theta, b = 1.$$

$$= 1 - [1 - \frac{\theta + t/n - \theta}{1-\theta}]^n = 1 - [1 - \frac{t}{n(1-\theta)}]^n \text{ and } \lim_{n \rightarrow \infty} [1 - \frac{t}{n}]^{1/n} = e^{-t}, \text{ with } X \sim n, b=1, a=\frac{t}{1-\theta},$$

$$so = 1 - \exp[-t/(1-\theta)] \text{ which is the cdf for } \text{Exp}(\lambda) \text{ where } \lambda = 1/(1-\theta)$$

MLE Review:

- To find MLE:

$$\cdot \text{Lhood} = \prod_{i=1}^n f(x_i)$$

Take log. Take $\frac{\partial}{\partial \theta}$ log L, set = 0, solve for param(s) θ .

MLE Properties:

- Consistency: $\hat{\theta}_{\text{MLE}} \xrightarrow{P} \theta$
- Asymptotically normal. (EXCEPT: 1. if θ is in support of distribution
2. if # params grows as n grows (Neyman Scott))
- Invariance: If $\hat{\theta}$ is MLE of θ , $g(\hat{\theta})$ is MLE of $g(\theta)$.
- Asymptotic Optimality: Smallest variance among all unbiased estimators.

Asymptotic Rel. Efficiency: ratio of variances of 2 estimators.

smaller variance = \uparrow efficient.

Fisher Info:

Asymptotic distribution of $\hat{\theta}_{\text{MLE}}$ is $N(\theta, I(\theta))$ where $I(\theta)$ = Fisher's Info.

$$I(\theta) = -E\left[\frac{d^2 \log f(x_i | \theta)}{d\theta^2}\right] = -E\left[\frac{\partial^2 \log f}{\partial \theta^2}\right]$$

Ex: $x_1, \dots, x_n \sim \text{Bernoulli}(p)$. $L = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum x_i} (1-p)^{n-\sum x_i}$

$$\log L = \sum x_i \log p + (n - \sum x_i) \log(1-p)$$

$$\frac{\partial}{\partial p} \log L = \frac{\sum x_i}{p} - \frac{(n - \sum x_i)}{1-p} \rightarrow \frac{d^2 \log L}{d^2 p} = \frac{-\sum x_i}{p^2} - \frac{n - \sum x_i}{(1-p)^2} = -\frac{n\bar{x}}{p^2} - \frac{n-n\bar{x}}{(1-p)^2}$$

$$-E\left[\frac{d^2 \log L}{d^2 p}\right] = -E\left[\frac{-n\bar{x}}{p^2} - \frac{n-n\bar{x}}{(1-p)^2}\right] = \frac{np}{p^2} + \frac{n(1-p)}{(1-p)^2} = \frac{n}{p(1-p)} = I_n(p)$$

Matrix form: elements are $\{I(\theta)\}_{i,j} = -E\left[\frac{\partial^2 \log L}{\partial \theta_i \partial \theta_j}\right]$

Jeffreys Prior: $p(\theta) = \sqrt{\det I(\theta)}$

Cramer-Rao Lower Bound: Formalize property that MLE has smallest var of all unbiased estimators.

$$\text{Var}(\hat{\theta}) \geq I(\theta)^{-1}$$

Recall, $\text{Var}(\hat{\theta}_{\text{MLE}}) = I(\theta)^{-1}$

When MLEs Don't Exist:

- Sometimes when θ is in support. (May exist still, but not asympt. normal.
Ex: $x \sim U(0, \theta)$. $\hat{\theta} = \max x$ is not in interval, so MLE doesn't exist. Would if $U(0, B]$.)
- when p grows as n grows (Neyman Scott Example.)

Consider n groups of normals, w. different means + shared variance.

$$m_1, \dots, m_k \sim N(\mu, \sigma^2), \dots, m_{k+1}, \dots, m_n \sim N(\mu_k, \sigma^2)$$

Goal: estimate σ^2 . Seems smart to take mean of sample vars from each group.

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n s_i^2 \text{ where } s_i^2 = \frac{1}{k} \sum_{j=1}^k (x_{ij} - \bar{x}_j)^2$$

We know $\frac{(k-1)s_i^2}{\sigma^2} \sim \chi^2_{k-1}$ so has mean = df = k-1.

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \frac{(k-1)s_i^2}{k} = \frac{\sigma^2}{kn} \sum_{i=1}^n \frac{(k-1)s_i^2}{\sigma^2} \xrightarrow{\text{WLLN}} \frac{\sigma^2}{k} E\left[\frac{(k-1)s_i^2}{\sigma^2}\right] = \frac{\sigma^2 \sigma^2 (k-1)}{k^2} = \frac{\sigma^2 (k-1)}{k} \neq \sigma^2 \frac{n}{n}$$

Bootstrap

1. Nonparametric bootstrap:

Let $x_1, \dots, x_n \sim f(x|\theta)$ where θ unknown. To estimate θ :

For b in $1:b$ {

- x_b^* = Sample of size n , with repl., from x_1, \dots, x_n
- $\hat{\theta}_b$ = mle of θ , calculated using x_b^* .

}

Then have $\hat{\theta}_1, \dots, \hat{\theta}_B$. calc $\hat{\theta} = \text{mean}(\hat{\theta}_1, \dots, \hat{\theta}_B)$

$$\text{var}(\hat{\theta}) = \text{var}(\hat{\theta}_1, \dots, \hat{\theta}_B)$$

Obtain CI using quantiles of $\hat{\theta}_1, \dots, \hat{\theta}_B$, or $\hat{\theta} \pm 1.96 \sqrt{\text{var}(\hat{\theta})}$

2. Parametric Bootstrap:

Similar to non-parametric, except:

- Begin by calc'g $\hat{\theta}_{\text{mle}}$ from orig. sample x_1, \dots, x_n .
- For each bootstrap iteration, instead of drawing sample from x_1, \dots, x_n , generate n random obs from $f(x|\hat{\theta}_{\text{mle}})$.

3. Little Bag of Bootstraps:

Way to parallelize when data set too large for a single bootstrap.

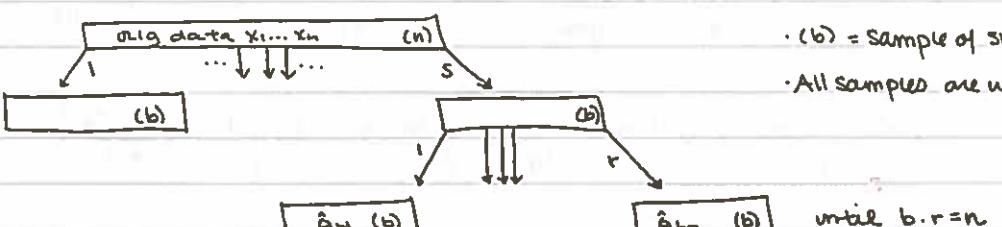
Sample without replacement the orig. data, into s samples of size b .

For each sample:

Resample & compute $\hat{\theta}$, until total sample size $= n$.

Take average of $\hat{\theta}_1, \dots, \hat{\theta}_s$

Careful of scaling: Ex: $\hat{\text{var}} = \frac{\sum(x_i - \bar{x})^2}{n}$, but here, $\frac{\sum(x_i - \bar{x})^2}{b}$, so mult'g by (b/n)



Estimating $\hat{\text{Var}}\theta$:

Have $\text{Var}_1 = \text{Var}(\theta_1, \dots, \theta_{1:b})$ to $\text{Var}_s = \text{Var}(\theta_{s1}, \dots, \theta_{sb})$. $\hat{\text{Var}}\theta = \text{mean}(\text{Var}_1, \dots, \text{Var}_s)$

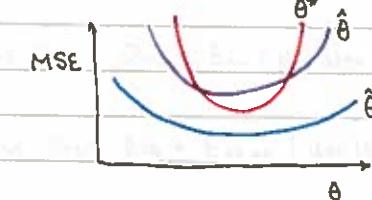
Scaling: $\text{Var}_1 = \frac{1}{b} \sum_{i=1}^b (\theta_{ii} - \bar{\theta}_1)^2$, which is why need $\text{Var}_1 = \text{Var}_1 \left(\frac{b}{n} \right)$ to ensure right scale.

Shrinkage Estimators

- ### 1. Dominance:
- Consider $\hat{\theta}$, an estimate of θ . This $\hat{\theta}$ dominates another estimator if $\hat{\theta}$ has smaller or equal MSE for all values $\theta \in \Theta$.

$$\text{MSE}(\hat{\theta}) \leq \text{MSE}(\tilde{\theta}), \text{ ie } E[(\hat{\theta}-\theta)^2] \leq E[(\tilde{\theta}-\theta)^2]$$

- ### 2. Admissible:
- $\hat{\theta}$ is admissible if no other estimator dominates it.



$\tilde{\theta}$ dominates both $\hat{\theta}$ and θ^* . If these are the only 3 estimators, $\tilde{\theta}$ is admissible.
 $\hat{\theta}, \theta^*$ do not dominate ea other - neither has lower MSE over the whole param space of θ .

To prove: Let $\tilde{\theta}$ be some arbitrary estimator. Show $E[(\tilde{\theta}-\theta)^2] \leq E[(\hat{\theta}-\theta)^2] \forall \theta$.

3. James Stein Estimator

Let $y \sim N(\mu, \sigma^2 I)$. Want to estimate μ for a single obs.

If $n=1$, $\hat{\mu} = \bar{y} = y_1$; just the value of that point.

$$\hat{\mu}_{\text{JS}} = \left[1 - \frac{(p-2)}{\|y\|^2} \right] y_1$$

Dominates $\hat{\mu}_{\text{mle}}$ for $p \geq 3$

Adds a bit of bias in exchange for lower variance; shrink est. of μ towards 0.

4. Empirical Bayes Estimator:

Let $y|\theta \sim N(\theta, I)$ Do normal-normal update:

$$\theta \sim N(0, \tau^2 I)$$

$$\text{posterior} \propto \exp\left[-\frac{1}{2}(y-\theta)^T I(y-\theta)\right] (\tau^2)^{-p/2} \exp\left[-\frac{1}{2\tau^2}\theta^T I \theta\right]$$

$$= (\tau^2)^{-p/2} \exp\left[-\frac{1}{2}\left\{\theta^T(I + \frac{1}{\tau^2}I)\theta - 2\theta^T y\right\}\right]$$

$$\sim N(m, v) \text{ with } v = (1 + \frac{1}{\tau^2})^{-1} I = \left(1 + \frac{1}{\tau^2}\right)^{-1} I, \quad m = \left(\frac{1}{1 + \frac{1}{\tau^2}}\right) y$$

"empirical" bc estimate $m = v$ using $\hat{\tau}^2$ from data.

m (posterior mean) is a shrinkage estimator. If plug in $\hat{\tau}^2$ mle, get James Stein.

Bias-Variance Tradeoff Proof:

- Let $\hat{\theta}$ be an estimator of θ .
- $MSE = E[(\hat{\theta} - \theta)^2]$
- $Var(\hat{\theta}) = E[(\hat{\theta} - E(\hat{\theta}))^2]$
- $Bias(\hat{\theta}) = E(\hat{\theta}) - \theta$

Prove that $MSE = Var + Bias^2$

$$\begin{aligned} MSE &= E[(\hat{\theta} - \theta)^2] = E[(\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta)^2] \\ &= E[(\hat{\theta} - E(\hat{\theta}))^2] + E[(E(\hat{\theta}) - \theta)^2] + E[(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)] \\ &\stackrel{\text{Var}}{=} (E(\hat{\theta}) - \theta)^2 = 0 \text{ bc } E[\hat{\theta} - E(\hat{\theta})] = E(\hat{\theta}) - E(\hat{\theta}) = 0 \\ &\stackrel{\text{Bias}^2}{=} \end{aligned}$$

so $MSE = Var + Bias^2$. ■

Gauss-Markov Proof (Repeat from Linear)

- Prove that $\hat{\beta}_{LS}$ is BLUE among all linear predictors of θ , with "best" = smallest variance.
- Let $a'y$ be an arbitrary linear predictor.
- We know: $E(y) = X\beta$, $\text{Var}(y) = \sigma^2 I$, $\hat{\beta} = p'My$ (or $\hat{\beta} = p'X\beta$). $\hat{\beta} \sim N(X\beta, \sigma^2(X'X)^{-1})$
- $\text{Var}(a'y) = \text{Var}(a'y - \hat{\beta} + \hat{\beta}) = \text{Var}(a'y - p'My + p'My)$

$$= \underbrace{\text{Var}(a'y - p'My)}_{\geq 0 \text{ b/c } a \text{ variance}} + \underbrace{\text{Var}(p'My)}_{\text{Var}(\hat{\beta}_{LS})} + \underbrace{2\text{Cov}(a'y - p'My, p'My)}_{= 0}$$

$$\text{Cov}(a'y - p'My, p'My) = (a' - p'M) \text{Cov}(y) M p = \sigma^2 (a' - p'M) M p = \sigma^2 (a'M - p'M)p$$

$$\text{Then } E(a'y) = \frac{p'My}{p'M} = p'M\beta \text{ b/c unbiased, so } a' = p', \text{ and } (a'M - p'M) = 0.$$

$$E(a'y) = a'E(y) = a'E\beta \quad \downarrow a'M = p'M$$

Therefore, $\text{Var}(a'y) \geq \text{Var}(\hat{\beta}_{LS})$ and $\hat{\beta}_{LS}$ is BLUE.

Proof that Training Error Always Less Than Test Error:

Let: • Expected test error = $R_{TE} = E_{y_{te}|y_{tr}} \left[\sum_{i=1}^n (y_{i,te} - y_{i,tr}|\hat{\beta})^2 \right] = E_{y_{te}} \left[(y_{te} - X\hat{\beta})^T (y_{te} - X\hat{\beta}) \right]$

• Expected training error = $R_{TR} = E_{tr} \left[(y_{tr} - X\hat{\beta})^T (y_{tr} - X\hat{\beta}) \right]$ where $\hat{\beta} = (X'X)^{-1}X'y_{tr}$

• Built model on y_{tr} , then use same $X, \hat{\beta}$ in test model for new y_{te} .

• Goal: Show $R_{TE} > R_{TR}$.

$$1) \text{ Expand } R_{TR}: R_{TR} = E_{tr} (y_{tr}^T y_{tr} - 2y_{tr}^T X\hat{\beta} + \hat{\beta}^T X^T X \hat{\beta}) = E_{tr} (y_{tr}^T y_{tr}) - E_{tr} (2y_{tr}^T X \hat{\beta}) + E_{tr} (\hat{\beta}^T X^T X \hat{\beta})$$

$$2) \text{ Expand } R_{TE}: R_{TE} = E_{y_{te}|y_{tr}} (y_{te}^T y_{te} - 2y_{te}^T X\hat{\beta} + \hat{\beta}^T X^T X \hat{\beta}) = E_{te} (y_{te}^T y_{te}) - E_{te} (2y_{te}^T X \hat{\beta}) + E_{te} (\hat{\beta}^T X^T X \hat{\beta})$$

$$\begin{aligned} 3) R_{TR} - R_{TE} &= E_{tr} (y_{tr}^T y_{tr}) - E_{tr} (2y_{tr}^T X \hat{\beta}) + E_{tr} (\hat{\beta}^T X^T X \hat{\beta}) - \left\{ E_{te} (y_{te}^T y_{te}) - 2E_{te} (y_{te}^T X \hat{\beta}) + E_{te} (\hat{\beta}^T X^T X \hat{\beta}) \right\} \\ &= -2 E_{tr} (y_{tr}^T X \hat{\beta}) - 2 E_{te} (y_{te}^T X \hat{\beta}) \end{aligned}$$

This has form $\text{cov}(u, v) = E(uv^T) - E(u)E(v^T)$, which = $E[(u - E(u))(v - E(v))^T]$

$$E(y) = X\beta$$

$$E(X\hat{\beta}) = XE(\hat{\beta}) = X(X'X)^{-1}X'E(y) = X(X'X)^{-1}X'\beta = X\beta$$

$$R_{TR} - R_{TE} = -2 E[(X\hat{\beta} - X\beta)(Y - X\beta)^T] = -2 E[\underbrace{(X(X'X)^{-1}X'y - X\beta)(Y - X\beta)^T}_M]$$

$$= -2 E[(My - X\beta)(Y - X\beta)^T]$$

$$= -2 \text{tr}\{ E[(y - X\beta)M(Y - X\beta)^T] \} = -2 \text{tr}\{ E[M(y - X\beta)(Y - X\beta)^T] \}$$

and $E(y - X\beta)(Y - X\beta)^T = \text{cov}(y - X\beta)$. Since $y \sim N(X\beta, \sigma^2 I)$, $y - X\beta \sim N(0, \sigma^2 I)$

$$= -2 \sigma^2 \text{tr}(M) = -2 \sigma^2 p$$

Therefore, $R_{TR} - R_{TE} = -2 \sigma^2 p$, so $R_{TE} > R_{TR}$. ■

Model Selection:

1. $AIC = -2\log L + 2p$ Favors models w/ more params.
 2. $BIC = -2\log L + p\log n$ Favors models w/ fewer params (larger penalty than AIC)
- Use for comparing models w/ same # of params.
• Bad for figuring out how many params to include in model. For this, need
Forward Selection, Backward Selection, CV or LOOCV.
3. Best Subsets: Look at all possible models w/ k features. Use AIC/BIC to compare.
This is not realistic bc $\binom{n}{k}$ models, And still need to know k , # of params.

4. Forward Selection:

- Begin w/ intercept-only model.
- Fit all models w/ 1 param; keep term w/ most significant β .
- Keep adding terms 1 by 1, until no more sig. predictors to add.
- Computationally expensive.

5. Backward Selection:

- Like forward, but begin with full model, - remove least sig. predictors 1 by 1, until no more insignificant predictors to remove.

6. Cross-Validation (CV):

- Use CV to pick k , the # of params to include.
- Also good for tuning parameters.

Cross-Validation:

• Use to estimate best # of predictors in a model. Est. test error for ea. k .

1. Split data randomly. 5 or 10 folds common.

c1	c2	c3	c4	c5	Test Set
----	----	----	----	----	----------

TRAIN set

2. Pick a model size, say $k=6$.

a. For fold C1, fit best model with $k=6$ params. Pick using best subset.
Fit model using C2-C5.

b. Calc $RSS^{(c1)} = (\mathbf{y} - \hat{\mathbf{x}}\hat{\beta})^T(\mathbf{y} - \hat{\mathbf{x}}\hat{\beta})$ using C1 as test data.

c. Repeat a,b for other folds, so have $RSS^{(c1)}, \dots, RSS^{(c5)}$

d. $RSS_k = \text{avg}(RSS^{(c1)}, \dots, RSS^{(c5)})$, where RSS_k is the estimated RSS for $k=6$.

3. Repeat (2) for all other model sizes you're interested in.

4. Pick the model size w/ the smallest RSS. Then estimate test error by fitting this model to the test set, ie calculate RSS_{test} . Note: fit the model with k-best params on ALL training data to get $\hat{\beta}$, then $RSS_{\text{test}} = (\mathbf{y}_{\text{te}} - \hat{\mathbf{x}}\hat{\beta})^T(\mathbf{y}_{\text{te}} - \hat{\mathbf{x}}\hat{\beta})$.

Caution! Wrong:
1. Screen for a good set of predictors.

2. Build a model using this subset.

3. Use CV to estimate any unknown tuning params & test error for final model.

Predictors have an unfair advantage! They were chosen in 1, seeing all the data.

- LOOCV: Same as above, but do n folds, leaving out 1 obs ea. time.

• Source of CV Bias:

We decreased n by removing a fold, so not fitting model using full n . (Want to estimate the test error using $\hat{\beta}$'s fitted w/ all n points.)

Can reduce bias w/ LOOCV, but increases var bc sets (folds) strongly correlated, since they share almost all data points.

Ridge Regression (Repeat from Linear)

• Like OLS, but adds ℓ_2 norm penalty to diagonal of $X^T X$ before inversion.

$$\text{loss fn: } \min_{\beta} \left\{ (y - X\beta)^T (y - X\beta) + \lambda \beta^T \beta \right\} = \min_{\beta} \left\{ \sum_{i=1}^n (y_i - x_i^T \beta)^2 - \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

• Scale and center data; not invariant to covariates w/ differing scales.

$$\text{Closed-form soln: } \hat{\beta}_R = (X^T X + \lambda I)^{-1} X^T y$$

• Motivation: Originally proposed to handle collinearity. If X not full rank, or close to strongly corr predictors, adding λ to diag makes $X^T X$ invertible.

$$\begin{aligned} \text{MSE Intuition: } \text{MSE} &= E[(\hat{\beta} - \beta)^T (\hat{\beta} - \beta)] = \text{tr}[\sigma^2 (X^T X)^{-1}] \quad \text{since} \\ &= \sigma^2 \sum_{j=1}^p \frac{1}{\phi_j^2} \text{ where } \phi_j^2 = \text{eigvals.} \end{aligned}$$

• If any eigenvalues too close to 0, MSE blows up.

• Adding λ , a penalty, to RR \rightarrow MSE = $\sigma^2 \sum_{j=1}^p \frac{1}{\phi_j^2 + \lambda}$, prevents explosion.

Bayesian Interpretation:

$$\begin{aligned} \text{Lhood: } p(y|\beta, \sigma^2) &\sim N(X\beta, \sigma^2 I) \quad \text{Then } \hat{\beta}_R = \text{posterior mean} \\ p(\beta|\sigma^2) &\sim N(0, \frac{\sigma^2}{\lambda} I) \end{aligned}$$

• SVD connection: $X = UDV^T$ as usual. $\hat{\beta}_R = (X^T X + \lambda I)^{-1} X^T y$

$$\Rightarrow \hat{y} = X\hat{\beta} = UDV^T (V D V^T U D V^T + \lambda I)^{-1} V D V^T y = U D (D^2 + \lambda I)^{-1} D V^T y$$

$$= \frac{UD^2U^Ty}{D^2 + \lambda I} = \sum_{j=1}^p u_j \left(\frac{\phi_j}{\phi_j + \lambda} \right) u_j^T y$$

scalar, ϕ_j = eigenval bc $D = \text{diag}(x_j)$, $\lambda_j^2 = \phi_j$

• So ridge computes coords of y wrt orthonormal basis U (like OLS), but also shrinks by factors $\left(\frac{\phi_j}{\phi_j + \lambda} \right)$, so \downarrow shrinkage applied to basis vecs w/ smaller eigenvals.

• Tuning penalty λ : Use CV to tune λ .

Generalized Linear Models:

$$1. \text{ logistic regression: } \log \left(\frac{p(y=1)}{1-p(y=1)} \right) = X\beta + e \quad \cdot \text{ i.e. } \Phi(y) = X\beta + e; \quad \Phi(y) = \text{logit}(y)$$

• models log-odds

• can fit using IRLS (equiv to Newton's method), SGD, etc.

$$\cdot p(y=1|X\beta) = \frac{e^{X\beta}}{1+e^{X\beta}}, \text{ call it "p". Then } \nabla \log(p) = \sum_{i=1}^n y_i (y_i - p).$$

• For $y = \text{binary response only}$

• Can set up as a latent variable model, with latent $z_i = X\beta + e$. Similar to probit approach.

* "link" fn, logit in this case, describes relationship b/w expected value of the response and the linear prediction of the covariates

$$2. \text{ Poisson glm: } E(y|x) = e^{X\beta}, \text{ so model is } \log \theta = X\beta + e, \text{ where } y = \text{counts}$$

• e^{β_i} represents effect of a 1-unit increase in x on θ .

Naive Bayes

- Classify $y \in \{0, 1\}$ binary points using Bayes Rule.
- "Naive" bc uses assumption of conditional independence.

$$\text{Bayes Rule: } P(y|x) = \frac{P(x|y)P(y)}{\sum_i P(x|i)P(i)}$$

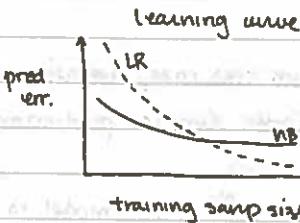
Reduces # of predictors by assuming condit indep

Regular Bayes Rule needs $2(2^k - 1) + 1$ where $k = \# \text{ predictors}$.
 \uparrow
 # levels for $P(y)$ prior

Naive Bayes only needs $(2k+1)$ params estimated.

NB vs Logistic Regression:

LR is a ↑ complex model - takes longer to converge.



Example (a) Discrete X : (HW 3, 2.2)

prior: $P(y=1) = .5$ Say have 2 binary covariates, $x_1 = x_2$

$$\begin{aligned} \text{given: } P(x_1=1|y=1) &= .8 & P(x_1=0|y=0) &= .7 \\ P(x_2=1|y=1) &= .5 & P(x_2=0|y=0) &= .9 \end{aligned}$$

(1) Calc decision rule for y , for every combo of (x_1, x_2) . If $P(y=1|x_1, x_2) \geq .5$, classify as $y=1$.

$$\text{Ex: } P(y=1|x_1=0, x_2=0) = P(x_1=0|y=1)P(x_2=0|y=1)P(y=1)$$

$$P(x_1=0|y=1)P(x_2=0|y=1)P(y=1) + P(x_1=0|y=0)P(x_2=0|y=0)P(y=0)$$

$$= .2(.5)(.5) / [.2(.5)(.5) + .7(.9)(.9)] = .14, \text{ so classify as } y=0.$$

Repeat for 3 other (x_1, x_2) combos.

(2) Error rate = sum $(a+b+c+d)$ where $a = P(x_1=0, x_2=0|y=1)$, ie of seeing these x 's w/o opp y .

$$a = P(x_1=0|y=1)P(x_2=0|y=1)P(y=1) = .5(.2)(.5) = .05$$

For b,c,d, repeat for each combo of x 's.

Example from Midterm:

$$f(x,y) = \text{bivariate uniform, ie } f(x,y) = \frac{1}{(x_n-x_1)(y_n-y_1)}$$

Data:	x	0	1	1	2	8	:	6	7	5	6	0
y	2	0	5	3	4	:	6	4	7	3	6	
class	1	1	1	1	1	:	2	2	2	2	2	

a) What are mles for x_1, x_2, y_1, y_2, y_3 ? $x_{(1)}, x_{(n)}, y_{(1)}, y_{(n)}$

b) Use Bayes Rule to estimate the point $(0,1)$

$$\begin{aligned} P(\text{class}=1|x=0, y=1) &= \frac{P(x=0, y=1|\text{class}=1)P(\text{class}=1)}{P(x=0, y=1|\text{class}=0)P(\text{class}=0)} \\ P(\text{class}=0|x=0, y=1) &= \frac{P(x=0, y=1|\text{class}=0)P(\text{class}=0)}{P(x=0, y=1|\text{class}=1)P(\text{class}=1)} \end{aligned}$$

$$\textcircled{1} \quad f(x, y | \text{class}=1) = \frac{1(0 < x < 8, 0 < y < 5)}{(8-0)(5-0)} = \frac{1}{40} = 1/40$$

$$\textcircled{2} \quad f(x, y | \text{class}=0) = \frac{1(0 < x < 7, 3 < y < 7)}{(7-0)(7-3)} = \frac{1}{28} = 1/28$$

$(0,1)$ outside support of class 2, so classified as class 1.

c) Use Bayes Rule to estimate point $(2,4)$

$$\frac{P(\text{class}=1|x=2, y=4)}{P(\text{class}=2|x=2, y=4)} = \frac{\frac{1}{40}}{\frac{1}{28}} = \frac{28}{40} = \frac{7}{10} < 1, \text{ so classify as class 2.}$$

K-Means Clustering:

- Given number of clusters k , want to classify obs $y = y_1, \dots, y_n$ into clusters.
- No y labels; unsupervised.
- Objective func: $\min_{\{m_1, \dots, m_k\}} \sum_{i=1}^n \|x - m_i\|^2$, i.e. want to divide obs into sets which minimize the within-cluster variance.
- m_k = cluster mean k .

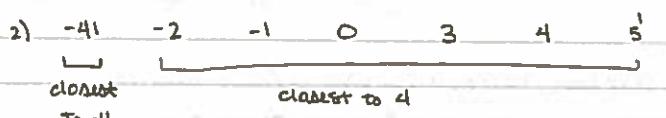
Algorithm:

- Randomly set k cluster centers to begin.
- Assign every point to its closest cluster.
- Recalc cluster mean as avg of all points in clusters.
- Repeat.

Ex: Say have $n=7$ points, $\{-4, -2, -1, 0, 3, 4, 5\}$. Use $K=2$, with initial centers at $m_1 = -1$, $m_2 = 4$.



New centers: $-1, 4$



New centers: $-1, 3\frac{1}{2}$.

No changes in future iter, so we are done.

Robust Statistics:

1. Estimating location parameters w outliers:

α -trimmed mean: trim α proportion of min/max obs from each tail, then compute the mean.

α -Winsorized mean: like α -trimmed, but replaces the points in the tails with the most extreme non-tail points, then calc mean.

Ex: Data set is $\{2, 4, 9, 10, 100\}$. 20% trimmed mean: $(4+5+10)/3 = 6.33$

20% Winsorized mean: $(4+4+5+10+10)/5 = 6.6$

2. Sensitivity Curve:

Measures effect of changing 1 data point in the set.

T = the sample statistic

$$SC_n(x, T) = \frac{T(x_1, \dots, x_{n-1}, x_{new}) - T(x_1, \dots, x_{n-1})}{x_n}$$

3. Influence Function:

$$\lim_{n \rightarrow \infty} SC_n(x, T)$$

If bounded (not $-\infty$ or ∞), we call the statistic T 'robust'.

Ex: With prob $(1-\alpha)$, points from $U(0, \theta)$. With prob α , points from other dist, i.e. noise.

(a) Say $\alpha=0$. What is mle of θ ? $\hat{\theta} = x_{(n)}$, the sample max.

$$(b) \text{ Sensitivity curve for } \hat{\theta}: SC_n(\hat{\theta}, x) = n [\max(x_1, \dots, x_{n-1}, x_{new}) - \max(x_1, \dots, x_{n-1})]$$

$$= n [\max\{x_{(n)}, x_{new}\} - x_{(n)}]$$

(c) Based on SC_n , is $\hat{\theta}$ appropriate if there are outliers, i.e. $\alpha > 0$?

No. If set $x_{new} = \infty$, $SC_n = n[\infty - x_{(n)}] = \infty$

(a) If know $\alpha=1$, suggest a robust variant of $\hat{\theta}$.

α -Trimmed sample max. Trim α off of top of sorted data set (and bottom, though won't make a difference here). Then compute sample max.

M-Estimators:

Generalized MLE. Instead of setting $\frac{\partial \log L}{\partial \theta} = 0$ & solving for $\hat{\theta}$, we instead minimize some loss fcn $p(x_i, \theta)$.

Instead of "Score fcn" $\frac{\partial \log L}{\partial \theta} = 0$, use $\psi(z) = \frac{d p(z)}{dz}$

Setting is $\hat{\theta} = \arg \min_{\theta} \left\{ \sum_{i=1}^n p(x_i, \theta) \right\}$ → solve for $\hat{\theta}$ using $\sum_{i=1}^n p(x_i, \theta) = 0$

Example: When $\theta = \text{mean } \mu$, use squared loss fcn $(x_i - \theta)^2$.

$$p(x_i, \theta) = (x_i - \theta)^2, \text{ so } \psi(x_i, \theta) = -2(x_i - \theta) \rightarrow \sum_{i=1}^n -2(x_i - \theta) = 0 \rightarrow \sum x_i - n\theta = 0 \rightarrow \hat{\theta} = \bar{x}$$

Example: When $\theta = \text{median}$, use absolute loss fcn $|x_i - \theta|$.

$$p(x_i, \theta) = |x_i - \theta|, \text{ so } \psi(x_i, \theta) = \text{sign}(x_i - \theta). \text{ Then } \sum_{i=1}^n \text{sign}(x_i - \theta) = 0 \rightarrow \hat{\theta} = \text{sample median } \tilde{x}$$

If $\psi(x, \theta)$ bounded for all x , then $\hat{\theta}$ is a robust estimator.

Asymptotic Relative Efficiency:

Ratio of asymptotic means of two estimators.

Estimator w. smaller variance is † efficient.

Ex: Mean vs. Median.

$$\bar{x} \sim N(\mu, \sigma^2/n)$$

$$\tilde{x} \sim N(\mu, \frac{\sigma^2}{2n})$$

$$\text{ARE}(\bar{x}, \tilde{x}) = \frac{\sigma^2 \pi / 2n}{\sigma^2 / n} = \frac{\pi/2}{1/2} = \frac{\pi}{2} > 1, \text{ so } \bar{x} \text{ more efficient than } \tilde{x}.$$

EM Algorithm: For A 2-Component Gaussian Mixture

$f(x|\theta) = \pi N(x|\mu_1, \sigma^2) + (1-\pi) N(x|\mu_2, \sigma^2)$ (say σ^2 known)
 $\theta = \{\pi, \mu_1, \mu_2\}$

Introduce latent z_i var to indicate class membership for each obs x_i . $z_i = 1 \text{ or } 0$.

$$\text{Original likelihood: } \prod_{i=1}^n \left[\underbrace{\pi N(x_i|\mu_1, \sigma^2)}_{\phi_1(x_i)} + (1-\pi) \underbrace{N(x_i|\mu_2, \sigma^2)}_{\phi_2(x_i)} \right] = \prod_{i=1}^n \left[\pi \phi_1(x_i) + (1-\pi) \phi_2(x_i) \right]$$

The latent z_i 's allow us to eliminate the summation in the likelihood.

$$\text{Augmented likelihood: } \prod_{i=1}^n \frac{\phi_1(x_i)^{z_i} \phi_2(x_i)^{1-z_i}}{f(x_i|z_i, \theta)} \frac{\pi^{z_i} (1-\pi)^{1-z_i}}{p(z_i)} \text{ so } \prod_{i=1}^n f(x_i|z_i, \theta) p(z_i)$$

$$\text{Log Likelihood: } \sum_{i=1}^n [z_i \log \phi_1(x_i) + (1-z_i) \log \phi_2(x_i) + z_i \log \pi + (1-z_i) \log (1-\pi)] \quad (\log L)$$

E Step: Using current params, update z_i 's, using their expected values.

$$P(z_i=1) = z_i = E(z_i|x_i, \theta) = \frac{\pi \phi_1(x_i)}{\pi \phi_1(x_i) + (1-\pi) \phi_2(x_i)}$$

Then classify $z_i=1$ if $P(z_i=1) > .5$, else $z_i=0$.

M Step: Using updated z_i 's, estimate updated MLEs for each param (π, μ_1, μ_2) .

$$(i) \frac{\partial \log L}{\partial \pi} = 0 \rightarrow \frac{\sum z_i}{\pi} - \frac{\sum (1-z_i)}{1-\pi} = 0 \rightarrow \frac{\pi}{1-\pi} = \frac{\sum z_i}{n - \sum z_i} \xrightarrow{\text{divide by } n} \frac{\pi}{1-\pi} = \frac{\sum z_i/n}{1 - \sum z_i/n}$$

$$\text{So } \hat{\pi} = \frac{\sum z_i}{n} \text{ where } z_i = 1 \text{ or } 0 \text{ based on E step updates.}$$

$$(ii) \frac{\partial \log L}{\partial \mu_1} = 0 \rightarrow \sum_{i=1}^n -2z_i(x_i - \mu_1) = 0 \rightarrow 2\sum z_i x_i - 2\sum z_i \mu_1 = 0 \rightarrow \hat{\mu}_1 = \frac{\sum z_i x_i}{\sum z_i}$$

which is the mean of the obs in updated group 1.

$$\cdot \text{ Comes from likelihood } L = \prod_{i=1}^n (\sigma^2)^{-2/2} \exp\left[-\frac{1}{2\sigma^2} (x_i - \mu_1)^{2z_i} (1-\pi)^{1-2z_i}\right] \pi^{z_i} (1-\pi)^{1-z_i}$$

$$(iii) \frac{\partial \log L}{\partial \mu_2} = 0 \rightarrow \sum_{i=1}^n -2(1-z_i)(x_i - \mu_2) = 0 \rightarrow \sum_{i=1}^n (1-z_i)x_i - \sum_{i=1}^n (1-z_i)\mu_2 = 0 \rightarrow \hat{\mu}_2 = \frac{\sum_{i=1}^n (1-z_i)x_i}{\sum_{i=1}^n (1-z_i)}$$

EM-Related Exponential Example

- Say there are 2 sets of (independent) bulbs.
- $\{y_1 \dots y_n\}$ We observe these the whole time, so know failure times.
- $\{E_1 \dots E_m\}$ We do not observe failure times, latent $\{z_1 \dots z_m\}$ - just whether bulbs in set 2 are on or off at time T . $\{E_i \dots E_m\}$ are 0 or 1, 1=on at T .
- $y \sim \text{Exp}(\theta)$, $z \sim \text{Exp}(\theta)$, where $f(x|\theta) = \frac{1}{\theta} e^{-x/\theta}$. So $F(x|\theta) = 1 - e^{-x/\theta}$
- Goal: estimate θ . If we had z^{obs} , $\hat{\theta} = \left(\frac{\sum_i y_i + \sum_j z_j}{n+m} \right)$
- Since don't have z^{obs} , we can replace them with their expected values, $E(z_i | E_i, \theta)$.
- Each z_i only depends on its own E_i and θ .
- $E_i \in 0 \text{ or } 1 \text{ only}$:

 - If $E_i = 1$, $E(z_i | E_i, \theta) = T + \theta$, since mean of $\text{Exp}(\theta) = \theta$, and exponential dist. is memoryless.
 - If $E_i = 0$, we can solve for this:

$$E(z_i) = E[z_i | E_i=0, \theta] P(E_i=0) + E[z_i | E_i=1, \theta] P(E_i=1)$$

θ ? prob of failure before T , ie
 bc mean of $\text{Exp}(\theta)$ from above prob of failure after time T , ie
 $e^{-T/\theta}$ 1-cdf eval at T :
 $e^{-T/\theta}$

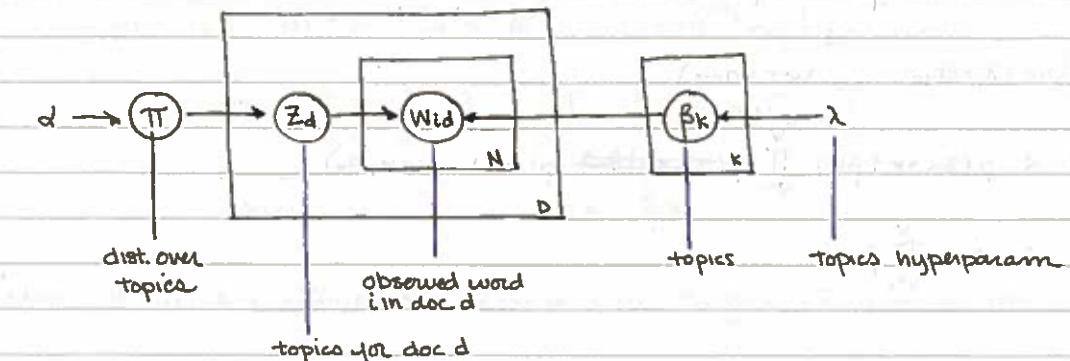
$$\theta = E[z_i | E_i=0, \theta] (1 - e^{-T/\theta}) + (T + \theta) e^{-T/\theta}$$

$$E[z_i | E_i=0, \theta] = \frac{\theta - (T + \theta) e^{-T/\theta}}{1 - e^{-T/\theta}}$$

- Plug the expressions for $E[z_i | E_i=1, \theta]$ and $E[z_i | E_i=0, \theta]$ into the expression for $\hat{\theta}$, in place of z^{obs} , based on E_i for each obs.

Topic Modeling Review

- Goal: Classify a document into a topic based on its words. Assume naive: 1 topic per doc
- How model is generated:
 - Randomly pick a topic z_d , from dist. over topics $\pi(z_1 \dots z_k) \sim \text{Dir}(d_1 \dots d_k)$
 - Topic z_d is just a distribution over words, fully specified by probs of words, $\beta = (\beta_1 \dots \beta_v)$ where $v = \# \text{ of words in vocab}$.
 - Words w_{id} doc d are randomly selected from vocab with probs β .



- $z_d = \text{topic of document } d, \text{ for } d \in \{1 \dots N\}$ (N total docs to classify)
- $z_d \in \{1 \dots k\}$, where $k = \# \text{ of topics}$
- $\pi \sim \text{Dir}(d_1 \dots d_k)$ is dist. of topics, so $P(z_d=k) = \pi_k$ for $d \in 1 \dots N$
- $w_{id} = \text{word } i \text{ in doc } d. P(w_{id} | z_d=k) = \beta_{k,i}$
- $\beta_k = (\beta_{k,1} \dots \beta_{k,v}) \sim \text{Dir}(\lambda_1 \dots \lambda_v)$ is probs for each word in vocab, in topic k . (Each topic has its own dist. of words.)

- Can derive full conditionals & use a Gibbs Sampler.

$$\begin{aligned}
 ① p(\pi) &\propto p(\pi) \cdot \prod_{d=1}^D p(z_d=k | \pi) \quad \text{with } p(\pi) \sim \text{Dir}(d_1 \dots d_k) \text{ and } p(z_d=k) = \pi_k \\
 &= \pi_1^{d_1-1} \dots \pi_k^{d_k-1} \cdot \prod_{d=1}^D \pi_1^{\mathbb{I}(z_d=1)} \dots \pi_k^{\mathbb{I}(z_d=k)} \\
 &= \pi_1^{d_1-1} \dots \pi_k^{d_k-1} \cdot \prod_{d=1}^D \frac{\pi_1^{\mathbb{I}(z_d=1)}}{\sum_{k=1}^D \pi_k^{\mathbb{I}(z_d=k)}} \dots \frac{\pi_k^{\mathbb{I}(z_d=k)}}{\sum_{k=1}^D \pi_k^{\mathbb{I}(z_d=k)}} \\
 &= \pi_1^{n_{1,d}-1} \dots \pi_k^{n_{k,d}-1} \sim \text{Dir}(n_1+d_1, \dots, n_k+d_k)
 \end{aligned}$$

Let $\sum_{d=1}^D \mathbb{I}(z_d=k) = n_k$, # docs classified as topic k .

$$\textcircled{2} \quad p(\beta_k | \dots) \propto p(\beta_k) \prod_{d=1}^P \prod_{i=1}^{n_d} p(w_{id} | \beta_k)$$

$$= \beta_{k1}^{\lambda_1-1} \cdots \beta_{kv}^{\lambda_v-1} \cdot \prod_{d=1}^P \prod_{i=1}^{n_d} \prod_{v=1}^V \beta_{kv}^{1(w_{id}=v)}$$

$$= \beta_{k1}^{\lambda_1-1} \cdots \beta_{kv}^{\lambda_v-1} \cdot \prod_{v=1}^V \beta_{kv}^{\sum_i \sum_d \mathbb{1}(w_{id}=v)} = m_{kv}, \text{ the # of times word } v \text{ appears in topic } k \text{ across all docs.}$$

$$= \beta_{k1}^{\lambda_1-1} \cdots \beta_{kv}^{\lambda_v-1} \cdot \beta_{k1}^{m_{k1}} \cdots \beta_{kv}^{m_{kv}}$$

$$\sim \text{Dir}(\lambda_1 + m_{k1}, \dots, \lambda_v + m_{kv})$$

$$\textcircled{3} \quad p(z_d=k | \dots) \propto p(z_d=k | \beta_k) \prod_{v=1}^V p(w_{id} | z_d=k, \beta_k)$$

$$= \pi_k \cdot \prod_{v=1}^V \beta_{kv}$$

Sufficient Statistics

1. **Factorization Theorem:** Let $f(x|\theta)$ be the joint density of sample x_1, \dots, x_n . $T(x)$ is a sufficient statistic for θ iff can factor $f(x|\theta)$ as

$$f(x|\theta) = g(T(x)|\theta) \cdot h(x)$$

Useful for finding a sufficient stat.

2. **Ex 1. Normal Mean:** $x_1, \dots, x_n \sim N(\theta, \sigma^2)$. Find $T(x)$, a sufficient stat for θ .

$$f(x|\theta) = \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp\left[-\frac{1}{2\sigma^2}(x_i - \theta)^2\right] = (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \theta)^2\right]$$

$$= (2\pi\sigma^2)^{-n/2} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2\right] \cdot \underbrace{\exp\left[-\frac{1}{2\sigma^2} \left\{ n\theta^2 - 2\theta \sum_{i=1}^n x_i \right\}\right]}_{h(x)}$$

$$T(x) = \sum_{i=1}^n x_i$$

Note: If $T(x)$ is a sufficient stat, $c(x)$ is also. So $\frac{1}{n}T(x) = \bar{x}$ is a sufficient stat.

3. **Bernoulli Ex:** $x_1, \dots, x_n \sim \text{Bernoulli}(p)$. Find a sufficient stat for p .

$$f(x|p) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum x_i} (1-p)^{n-\sum x_i} \text{ so } h(x) \equiv 1, T(x) = \sum_{i=1}^n x_i$$

4. **Weibull Ex:** β is a parameter, γ is known.

$$f(x|\beta) = \prod_{i=1}^n \frac{x_i}{\beta} \cdot x_i^{\gamma-1} \exp\left[-\frac{x_i^\gamma}{\beta}\right] = \left(\frac{x}{\beta}\right)^n \exp\left[-\sum x_i^\gamma / \beta\right] \cdot \prod_{i=1}^n x_i^{\gamma-1}$$

$$\text{so } h(x) = \prod_{i=1}^n x_i^{\gamma-1} \text{ and } T(x) = \sum_{i=1}^n x_i^\gamma$$

- ④ For exponential families, $f(x|\theta) = h(x) c(\theta) \exp\left(\sum_{i=1}^n w_i(\theta) t_i(x)\right)$

Then $T(x) = [\sum_{i=1}^n t_1(x_i), \dots, \sum_{i=1}^n t_K(x_i)]$ is a suff stat for θ .

Minimal Sufficient Statistic:

- A sufficient statistic $T(X)$ is minimal if for any other stat $\tilde{T}(x)$, $T(x)$ is a func of $\tilde{T}(X)$.
- Ex: \bar{X} is the min suff stat. The whole data set, $x_1 \dots x_n$, is a suff stat.
 \bar{X} is a func of $x_1 \dots x_n$.

- $cT(X)$ is a min suff stat if $T(X) \neq 0$, where c = constant.

How to find a min. suff. stat? Thm 6.2.13 (Berger + Casella).

Say $\exists T(X) \geq 0$ so that $f(x|\theta)/f(y|\theta)$ is constant wrt θ . Then $T(X)$ is min suff stat.

Ex: Find a min suff stat for Weibull, for β , where γ is known.

$$f(x|\beta) = \prod_{i=1}^n \frac{\gamma}{\beta} x_i^{\gamma-1} \exp[-x_i^\gamma/\beta] = \left(\frac{\gamma}{\beta}\right)^n \exp\left[-\sum_i x_i^\gamma/\beta\right] \prod_{i=1}^n x_i^{\gamma-1}$$

$$\text{Ratio: } \frac{\left(\frac{\gamma}{\beta}\right)^n \exp\left[-\sum_i x_i^\gamma/\beta\right] \prod_{i=1}^n x_i^{\gamma-1}}{\left(\frac{\gamma}{\beta}\right)^n \exp\left[-\sum_i y_i^\gamma/\beta\right] \prod_{i=1}^n y_i^{\gamma-1}}$$

• Need $\sum_i x_i^\gamma = \sum_i y_i^\gamma$ so β 's cancel, i.e. ratio is constant wrt β .

• $T(X) = \sum_i x_i^\gamma$ is min suff stat.

Ex: Find a min sufficient stat for θ in $U[0, \theta]$.

$$\frac{f(x|\theta)}{f(y|\theta)} = \frac{\prod_{i=1}^n \frac{1}{1-\theta} I_{[0,\theta]}(x_i)}{\prod_{i=1}^n \frac{1}{1-\theta} I_{[0,\theta]}(y_i)} = \frac{\left(\frac{1}{\theta}\right)^n I(X)_{[0,\theta]}}{\left(\frac{1}{\theta}\right)^n I(Y)_{[0,\theta]}}$$

• Can split indicators:

• If $\max(x) > \theta$, then all $x < \theta$. If $\min(x) > 0$, all $x > 0$, but doesn't depend on θ .

$$\frac{\prod_{i=1}^n I(x_i)_{(-\infty, 0]} I(x_i)_{[0, \theta]}}{\prod_{i=1}^n I(y_i)_{(-\infty, 0]} I(y_i)_{[0, \theta]}}$$

$$\frac{I(X_m)_{(-\infty, 0]}}{I(Y_m)_{(-\infty, 0]}}$$

so X_m is min suff stat.

$$\frac{\prod_{i=1}^n I(y_i)_{(-\infty, 0]} I(y_i)_{[0, \theta]}}{\prod_{i=1}^n I(y_i)_{(-\infty, 0]}}$$

$$\frac{I(Y_m)_{(-\infty, 0]}}{I(Y_m)_{(-\infty, 0]}}$$

Chapter 1 - Basics

Calculus Review: U-substitution

$$1) \int (2x+5)(x^2+5x)^3 dx \quad \text{Let } u = x^2+5x, \text{ so } du = (2x+5)dx. \quad -\int u^3 du = \frac{u^4}{4} = \frac{(x^2+5x)^4}{4}$$

$$2) \int (3-x)^{10} dx \quad \text{Let } u = 3-x, \text{ du} = -dx. \quad -\int u^{10} du = -\frac{u^{11}}{11} = -\frac{(3-x)^{11}}{11}$$

$$3) \int (7x+9)^{1/2} dx \quad \text{Let } u = 7x+9, \text{ du} = 7dx. \quad \frac{1}{7} \int u^{1/2} du = \frac{1}{7} \cdot \frac{u^{3/2}}{3/2} = \frac{1}{21} u^{3/2} = \frac{2(7x+9)^{3/2}}{21}$$

$$4) \int x^3(1+x^4)^{-1/3} dx \quad \text{Let } u = 1+x^4, \text{ du} = 4x^3 dx. \quad \frac{1}{4} \int u^{-1/3} du = \frac{1}{4} \cdot \frac{u^{2/3}}{2/3} = \frac{3}{8} u^{2/3} = \frac{3(1+x^4)^{2/3}}{8}$$

$$5) \int e^{5x+2} dx. \quad \text{Let } u = 5x+2, \text{ du} = 5dx. \quad \frac{1}{5} \int e^u du = \frac{1}{5} e^u = \frac{1}{5} e^{5x+2}$$

$$6) \int 4 \cos(3x) dx. \quad \text{Let } u = 3x, \text{ du} = 3dx. \quad \int \frac{4}{3} \cos(u) du = \frac{4}{3} \sin(u) = \frac{4}{3} \sin(3x)$$

$$7) \int \frac{\sin(\ln(x))}{x} dx. \quad \text{Let } u = \ln(x), \text{ du} = \frac{1}{x} dx. \quad \int \sin(u) du = -\cos(u) = -\cos(\ln(x)).$$

Calculus Review: Integration by Parts

$$\int u dv = uv - \int v du$$

$$1) \int x e^x dx. \quad \text{Let } u=x, \text{ dv} = e^x dx. \quad \text{Then } du=dx, \quad v=\int e^x dx = e^x.$$

$$= x e^x - \int e^x dx = x e^x - e^x$$

$$2) \int x \ln(x) dx. \quad \text{Let } u=\ln(x), \text{ dv} = x dx. \quad \text{Then } du=\frac{1}{x} dx, \quad v=\int x dx = \frac{x^2}{2}$$

$$= uv - \int v du = \frac{x^2}{2} \ln(x) - \int \frac{x^2}{2} \left(\frac{1}{x}\right) dx = \frac{x^2}{2} \ln(x) - \int \frac{x}{2} dx = \frac{x^2}{2} \ln(x) - \frac{x^2}{4}$$

$$3) \int x \cos(3x) dx. \quad \text{Let } u=x, \text{ dv} = \cos(3x) dx. \quad \text{Then } du=dx, \quad v=\int \cos(3x) dx. \quad \text{Let } u'=3x, \text{ du}'=3dx$$

$$\text{so } v = \int \cos(u') du' = \frac{1}{3} \sin(u') = \frac{1}{3} \sin(3x).$$

$$\text{Then } uv - \int v du = \frac{x}{3} \sin(3x) - \int \frac{1}{3} \sin(3x) dx = \frac{x}{3} + \underbrace{\left(\frac{1}{9}\right) \cos(3x)}$$

another u-sub

$$4) \int_{x^2} \ln(x) dx = \int \ln(x) \cdot x^2 dx. \text{ Let } u = \ln(x), dv = x^2 dx. \text{ Then } du = \frac{1}{x} dx, v = \int x^2 dx = \frac{x^3}{3}$$

$$= uv - \int v du = -\frac{\ln(x)}{4} - \int \frac{-x^4}{4} \left(\frac{1}{x}\right) dx = -\frac{\ln(x)}{4} + \int \frac{1}{4} x^3 dx = -\frac{\ln(x)}{4} - \frac{1}{4} \frac{x^4}{4} = -\frac{\ln(x)}{4} - \frac{1}{16} x^4$$

$$5) \int \ln(x) dx. \text{ Let } u = \ln(x), dv = (1) dx. \text{ Then } du = \frac{1}{x} dx, v = \int 1 dx = x.$$

$$uv - \int v du = x \ln(x) - \int x \left(\frac{1}{x}\right) dx = x \ln(x) - \int (1) dx = x \ln(x) - x.$$

$$6) \int x^2 e^{3x} dx. \text{ Let } u = x^2, dv = e^{3x} dx. \text{ Then } du = 2x dx, v = \int e^{3x} dx = \frac{e^{3x}}{3}$$

$$uv - \int v du = \frac{x^2 e^{3x}}{3} - \int \frac{e^{3x}}{3} (2x) dx$$

Needs a second int. by parts:

$$\text{Let } u = 2x, dv = \frac{e^{3x}}{3} dx, \text{ then } du = 2dx, v = \int \frac{e^{3x}}{3} dx = \frac{e^{3x}}{9}$$

$$uv - \int v du = \frac{2x e^{3x}}{9} - \int \frac{2}{9} e^{3x} dx = \frac{2x e^{3x}}{9} - \frac{2 e^{3x}}{27} = \frac{2x e^{3x}}{9} - \frac{2 e^{3x}}{27}$$

$$\text{Put back together: } = \frac{x^2 e^{3x}}{3} - \frac{2x e^{3x}}{9} + \frac{2 e^{3x}}{27}$$

CH 1 Basic Formulas:

$$1) \text{ Bayes' Thrm: } P(Y|x) = \frac{P(x|Y)P(Y)}{P(x)}, \text{ where } P(x) = \int P(x|y)P(y) dy \text{ or } \sum_i P(x|Y_i)P(Y_i)$$

$$2) \text{ Def. of conditional probability: } P(Y|x) = P(X,Y)/P(x)$$

3) To prove a form is a pdf:

1. $f(x)$ is non-negative; $f(x) \geq 0 \forall x$

2. $f(x)$ integrates (or sums, for discrete) to 1: $\int_{-\infty}^{\infty} f(x) dx = 1$ or $\sum_x f(x) = 1$.

CH 2 Transformations + Expectations

• Thm 2.1.3 - Cdf of a transformation: X has cdf $F_x(x)$. Let $Y = g(x)$.

• Check if $g(x)$ increasing ($\frac{d}{dx}g(x) > 0$) or decreasing.

• If $g(x)$ increasing, $F_y(y) = F_x(g^{-1}(y))$

• If $g(x)$ decreasing, $F_y(y) = 1 - F_x(g^{-1}(y))$

• Ex: $X \sim U(0,1)$. Let $Y = g(x) = -\log(x)$. $y = -\log(x) \rightarrow x = e^{-y} \rightarrow g^{-1}(y) = e^{-y}$

• Check $g'(x) = \frac{d}{dx} -\log(x) = \frac{1}{x} < 0$ since $x \in (0,1)$. So decreasing.

• $F_y(y) = 1 - e^{-y}$, which we recognize as the cdf of an exponential (1).

• Thm 2.1.5 - Pdf of a transformation: Let X , and $Y = g(x)$.

• $f_y(y) = f_x(g^{-1}(y)) |J|$ where $|J| = \left| \frac{d}{dy} g^{-1}(y) \right|$

• Ex: $f(x) = \frac{b-a}{r(a)} x^{a-1} e^{-bx}$, $0 < x < \infty$. Find pdf of $Y = 1/x$. ($X = g(a,b)$).

• Find $g^{-1}(y) = Y = 1/x \rightarrow x = 1/y$, so $g^{-1}(y) = 1/y = y^{-1}$

$$|J| = \left| -y^{-2} \right| = 1/y^2$$

$$f_y(y) = f_x(g^{-1}(y)) |J| = \frac{b-a}{r(a)} \left(\frac{1}{y}\right)^{a-1} e^{-b(1/y)} \left(\frac{1}{y^2}\right) = \frac{b-a}{r(a)} y^{-a-1} e^{-by} \sim Ig(a,b)$$

• Ex 2.1.7, Square Transformation: Let X be continuous, and $Y = X^2$. Find cdf of Y .

$$P(Y \leq y) = P(X^2 \leq y) = P(-\sqrt{y} \leq X \leq \sqrt{y}) = P(X \leq \sqrt{y}) - P(X \leq -\sqrt{y})$$

• Ex 2.3.9 - Normal-Chi-square Relationship: Let $X \sim N(0,1)$. $Y = X^2$. Find pdf of Y .

• Since $Y = X^2 \rightarrow X = \sqrt{Y}, -X = \sqrt{Y}$, split into 2 partitions.

• A₁: $X \in (-\infty, 0)$: $Y = X^2 \rightarrow X = -\sqrt{Y}$ so $g^{-1}(y) = -\sqrt{y}$. } Support of Y is $y \geq 0$ for

• A₂: $X \in (0, \infty)$: $Y = X^2 \rightarrow X = \sqrt{Y}$ so $g^{-1}(y) = \sqrt{y}$ } both

$$\text{Jacobian: } A_1: |J| = \left| \frac{d}{dy} g^{-1}(y) \right| = \left| \frac{1}{2} y^{-1/2} \right| = \frac{1}{2\sqrt{y}}$$

$$A_2: |J| = \left| \frac{d}{dy} g^{-1}(y) \right| = \left| \frac{1}{2} y^{-1/2} \right| = \frac{1}{2\sqrt{y}}$$

$$\text{pdf: } f(y) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2} (\sqrt{y})^2\right] \left(\frac{1}{2\sqrt{y}}\right) + \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2} (-\sqrt{y})^2\right] \left(\frac{1}{2\sqrt{y}}\right)$$

$$= \frac{1}{\sqrt{2\pi}} e^{-y/2} \frac{1}{2\sqrt{y}} + \frac{1}{\sqrt{2\pi}} e^{-y/2} \frac{1}{2\sqrt{y}} = \frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{y}} e^{-y/2} \sim \chi^2_1$$

CHAPTER 3 - Approximations + Relationships

Thm 2.10 - Probability Integral Transfer Theorem:

- X has cdf $F(x)$. Define Y as $Y = F(X)$. Then $Y \sim U(0,1)$
- Used for random number generation.

2.2: Expected Values:

- $E(g(x)) = \int g(x) f(x) dx$, or $\sum g(x) f(x)$ for discrete.

Ex: Exponential mean. $X \sim \text{Exp}(\lambda)$. Find $E(X)$. $f(x) = \frac{1}{\lambda} e^{-x/\lambda}$.

$$E(X) = \int x \left(\frac{1}{\lambda}\right) e^{-x/\lambda} dx. \text{ Int by parts: } \int u dv = uv - \int v du.$$

$$\text{Let } u = x, dv = \frac{-e^{-x/\lambda}}{\lambda} dx \rightarrow du = dx, v = -e^{-x/\lambda} \text{ bc } v = \int \frac{1}{\lambda} e^{-x/\lambda} dx = -e^{-x/\lambda}$$

$$\text{Then } E(X) = -x e^{-x/\lambda} \Big|_0^\infty + \int_0^\infty e^{-x/\lambda} dx$$

$\exp(\lambda)$ kernel, int. to $1/\lambda$

$$\lim_{x \rightarrow \infty} (-x e^{-x/\lambda}) - \lim_{x \rightarrow 0} (-x e^{-x/\lambda}) + \frac{1}{\lambda} = \frac{1}{\lambda}, \text{ so } E(X) = 1/\lambda.$$

Ex: Poisson mean: $Y \sim \text{Pois}(\lambda)$, so $p(Y=y) = \frac{\lambda^y e^{-\lambda}}{y!}$. Find $E(Y)$.

$$E(Y) = \sum_{y=0}^{\infty} y \frac{\lambda^y e^{-\lambda}}{y!} \text{ Then } y=0 \text{ term is } 0, \text{ so } = \sum_{y=1}^{\infty} y \frac{\lambda^y e^{-\lambda}}{y!} = \sum_{y=1}^{\infty} \frac{\lambda^y e^{-\lambda}}{(y-1)!}$$

$$= \sum_{y=1}^{\infty} \frac{\lambda^y e^{-\lambda}}{(y-1)!} \rightarrow \lambda \sum_{y=1}^{\infty} \frac{\lambda^{y-1} e^{-\lambda}}{(y-1)!} \text{ Then let } X = y-1: \lambda \sum_{x=0}^{\infty} \frac{\lambda^x e^{-\lambda}}{x!} = \lambda.$$

So $E(Y) = \lambda$.

MGFs:

- 1) $M_X(t) = E(e^{xt})$.

- 2) If $X \perp Y$, $M_{X+Y}(t) = M_X(t)M_Y(t)$. (B if $Z = X+Y$, $M_Z(t) = M_X(t)M_Y(t)$).

- 3) How to find moments: n^{th} moment of X : $M_X(t)^{(n)} = \frac{d^{(n)}}{dt^{(n)}} M_X(t)$ eval at $t=0$.

(Intⁿ deriv of $M_X(t)$ wrt t , eval at $t=0$.)

- 4) $M_{X+b}(t) = E(e^{(Ax+b)t}) = e^{bt} E(e^{At}) = e^{bt} M_X(at)$

Binom Approximation of Poisson:

- If $X \sim \text{Bin}(n, p)$ and $Y \sim \text{Pois}(\lambda)$, then $P(X=x) \approx P(Y=x)$ for large n , small np .

Normal Approximation to Binomial:

$\text{Binom}(n, p) \approx N(np, n(1-p))$

ok if $\min(np, n(1-p)) \geq 5$

Continuity correction: $P(X \leq x) \approx P(Y \leq x + 1/2)$

$P(X \geq x) \approx P(Y \geq x - 1/2)$

Biased Estimation:

Bias: $E(\hat{\theta}) - \theta$ for estimator $\hat{\theta}$ of θ .

Can derive unbiased estimators. Example:

Ex: $y_1, \dots, y_n \stackrel{iid}{\sim} \text{Exp}(\theta)$. Find an unbiased estimate for $\gamma = \lambda^{-1}$, the sample min.

$$f(x) = \frac{1}{\theta} e^{-x/\theta}, x \geq 0, \theta > 0$$

$$F(x) = 1 - e^{-x/\theta}$$

$$f_{Y_{(1)}}(x) = \frac{n!}{(n-1)!} \frac{1}{\theta} e^{-x/\theta} [1 - e^{-x/\theta}]^{n-1} = \frac{n}{\theta} e^{-x/\theta} e^{-x(n-1)/\theta} = \frac{n}{\theta} e^{-x/\theta} \text{ So } f(y) = \frac{n}{\theta} e^{-y/\theta}$$

$$\text{Then Bias} = E(Y) - \theta, \text{ so } E(Y) = \int y \frac{n}{\theta} e^{-y/\theta} dy = \frac{n}{\theta} \int y e^{-y/\theta} dy$$

$$\text{Let } u = y, du = dy, \\ du = e^{-y/\theta}, \\ v = e^{-y/\theta}, dv = -\frac{1}{\theta} e^{-y/\theta}$$

$$= \left(\frac{n}{\theta}\right) \left(\frac{\theta}{n}\right)^2 = \frac{\theta}{n}$$

$$\text{or, } \int y e^{-y/\theta} dy \sim$$

$$g(z, \frac{n}{\theta}) \text{ so integrated to } r(z) \left(\frac{n}{\theta}\right)^2 = \left(\frac{\theta}{n}\right)^2$$

Estimator is biased: $E(\hat{\theta}) - \theta = \frac{\theta}{n} - \theta \neq 0$

To find unbiased: $a(\frac{\theta}{n}) - \theta = 0 \rightarrow a = \theta(\frac{n}{\theta}) \rightarrow a = n$, so

use $\hat{\theta} = ny = n\bar{y}$ as the unbiased estimator.

- Ex: $Z_i = \begin{cases} \text{"blue"} \text{ w. prob } \theta & X_i \sim \text{Pois}(1) \text{ if } Z_i = \text{"blue"} \\ \text{"red"} \text{ w. prob } 1-\theta & \text{Pois}(4) \text{ if } Z_i = \text{"red"} \end{cases}$

Only observe X_i 's. Find consistent estimator of θ .

$$(\hat{\theta} \xrightarrow{P} \theta)$$

$$E(X_i) = E(E(X_i|Z)) = 1\theta + 4(1-\theta) = 4-3\theta, \text{ so } E(X_i) = 4-3\theta$$

So since $\bar{X} \xrightarrow{P} E(X)$ by the WLN, use $\bar{X} = 4-3\theta \rightarrow \hat{\theta} = \frac{4-\bar{X}}{3}$,

so by the WLN, $\hat{\theta} \xrightarrow{P} \theta$.

Max Likelihood Estimation Review:

Binom(n, p) example: $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Bin}(n, p)$. Find mle of p .

Lhood = $\prod_{i=1}^n \binom{n}{x_i} p^{x_i} (1-p)^{n-x_i} \propto p^{\sum x_i} (1-p)^{n-\sum x_i}$? Nope - we always add binomial RVs together to get 1 Binomial RV, with $n = \text{sum of } n$'s.

So: Lhood for $x \sim \text{Binom}(n, p) = \binom{n}{x} p^x (1-p)^{n-x}$, ie $\propto p^{\sum x_i} (1-p)^{n-\sum x_i}$ for $x_i = \text{Bernoulli trials}$.

$$\log L \propto \sum_i x_i \log p + (n - \sum x_i) \log(1-p)$$

$$\frac{\partial}{\partial p} \log L = \frac{\sum x_i}{p} - \frac{(n - \sum x_i)}{1-p} = 0 \rightarrow \hat{p} = \frac{\sum x_i}{n - \sum x_i} \rightarrow \hat{p} = \frac{\bar{x}}{1-\bar{x}} \text{ so } \hat{p} = \bar{x}$$

Jeffreys Prior

A vague prior, which can lead to a proper or improper prior.

Equal to sqrt(Fisher's info)

Jeffreys Prior $\propto \sqrt{\det I(\theta)}$

$$1) \text{ Univariate case: } \pi(\theta) \propto \sqrt{\det I(\theta)}, \quad I = -E_x \left[\frac{\partial^2}{\partial \theta^2} \log L \right]$$

$$2) \text{ Multivariate case: } \pi(\theta) \propto \sqrt{\det(I(\theta))}, \quad I = -E_x \left[\frac{\partial^2}{\partial \theta_i \partial \theta_j} \log L \right]_{p \times p}$$

Ex 1: $X \sim \text{Binom}(n, p)$ for a single obs x , with n known. $\log L \propto x \log p + (n-x) \log(1-p)$

$$I(p) = -E_x \left[\frac{\partial^2 \log L}{\partial p^2} \right] \Rightarrow \frac{d \log L}{dp} = \frac{x}{p} - \frac{(n-x)}{1-p} \rightarrow \frac{\partial^2 \log L}{\partial p^2} = \frac{x}{p^2} + \frac{(n-x)}{(1-p)^2}$$

$$\text{so } I(p) = -E_x \left[\frac{x}{p^2} + \frac{n-x}{(1-p)^2} \right] \text{ and } E(x) = np, \text{ so } - \left[\frac{np}{p^2} + \frac{n-np}{(1-p)^2} \right] = \frac{n}{p} + \frac{n}{(1-p)} = \frac{n}{p(1-p)}$$

$$\text{Then Jeffreys } \pi(\theta) \propto \sqrt{\frac{n}{p(1-p)}} \propto p^{-1/2} (1-p)^{-1/2} \sim \text{Beta}(1/2, 1/2)$$

Ex 2: $x | \theta \sim N(\mu, \sigma^2)$ so $\theta = (\mu, \sigma^2)$. Again, let X be a single obs.

$$\log L \propto -\log(\sigma) - \frac{1}{2\sigma^2} (X-\mu)^2$$

2nd partial derivs:

$$(1st) \rightarrow \frac{\partial}{\partial \mu} \log L = +2(X-\mu) = \frac{(X-\mu)}{\sigma^2}$$

$$\frac{\partial}{\partial \sigma} \log L = -\frac{1}{\sigma} + \frac{(X-\mu)^2}{\sigma^3}$$

$$(2nd) \rightarrow \frac{\partial^2 \log L}{\partial \mu^2} = -\frac{1}{\sigma^2} \quad \frac{\partial^2 \log L}{\partial \mu \partial \sigma} = \frac{\partial^2 \log L}{\partial \sigma \partial \mu} = \frac{-2(X-\mu)}{\sigma^3} \quad \frac{\partial^2 \log L}{\partial \sigma^2} = \frac{1}{\sigma^2} - \frac{3(X-\mu)^2}{\sigma^4}$$

$$\cdot \text{Expected vals: } E_x \left(\frac{\partial^2 \log L}{\partial \mu^2} \right) = -1/\sigma^2 \quad E_x \left(\frac{\partial^2 \log L}{\partial \mu \partial \sigma} \right) = 0 \quad E_x \left(\frac{\partial^2 \log L}{\partial \sigma^2} \right) = -2/\sigma^2$$

$$\cdot \text{Fisher's info: } I(\theta) = \begin{bmatrix} 1/\sigma^2 & 0 \\ 0 & 2/\sigma^2 \end{bmatrix}$$

$$\cdot \text{Jeffreys: } \pi(\theta) = \sqrt{\det I(\theta)} = \sqrt{(1/\sigma^2)(2/\sigma^2) - 0} = \sqrt{2/\sigma^4} = \sqrt{2}/\sigma^2 \propto 1/\sigma^2, \text{ so}$$

$$\pi(\theta) \propto \frac{1}{\sigma^2} \text{ for } \sigma > 0.$$

Deriving Densities of Order Stats $X_{(1)} \dots X_{(n)}$:

For any order stat, pdf is: $f(x_{(j)}) = \frac{n!}{(j-1)!(n-j)!} f(x) F(x)^{j-1} [1-F(x)]^{n-j}$

Can easily derive the CDFs for $X_{(1)} \dots X_{(n)}$: Say $X_1 \dots X_n \stackrel{iid}{\sim} F(x)$.

1) $Y = X_{(1)}$, sample min: $P(Y \leq y) = P(\min(X_1 \dots X_n) \leq y) = 1 - P(\text{all } X_1 \dots X_n \geq y)$

Then since X_i 's iid, $= 1 - (P(X \geq y))^n = 1 - [1 - P(X \leq y)]^n$ is cdf $F_{X_{(1)}}(y) = 1 - [1 - F(y)]^n$

For pdf, $f_{X_{(1)}}(y) = \frac{d}{dy} F_{X_{(1)}}(y) = -n [1 - F(y)]^{n-1} (-f_y(y)) = n [1 - F(y)]^{n-1} f_y(y)$

2) $Y = X_{(n)}$, sample max: $P(Y \leq y) = P(\text{all } X_i \leq y) = [P(X \leq y)]^n = F_y(y)$

Then pdf is $n F_y(y)^{n-1} f_y(y)$

Transformations Review: Non-monotone:

$f(x) = \frac{1}{2} e^{-|x|}$, $-\infty < x < \infty$. $y = |x|^3$. Find $f(y)$.

$A_1 = (0, \infty)$: $f(x) = \frac{1}{2} e^{-x} \rightarrow g(x) = x^3$, so $g^{-1}(y) = y^{1/3}$
 $A_2 = (-\infty, 0)$: $f(x) = \frac{1}{2} e^{-x} \rightarrow g(x) = x^3$, so $g^{-1}(y) = -y^{1/3}$

for A_1 : $\left| \frac{d}{dy} g^{-1}(y) \right| = \frac{1}{3} y^{-2/3}$ for A_2 : $\left| \frac{d}{dy} g^{-1}(y) \right| = -\frac{1}{3} y^{-2/3} = \frac{1}{3} y^{-2/3}$

$f_y(y) = \frac{1}{2} e^{-\left(\frac{1}{3} y^{1/3}\right)} \left(\frac{1}{3} y^{-2/3} \right) + \frac{1}{2} e^{-\left(-\frac{1}{3} y^{1/3}\right)} \left(-\frac{1}{3} y^{-2/3} \right) = \frac{1}{3} y^{2/3} e^{-y^{1/3}}$ for $y > 0$.

Mgf's of Transformations:

For $X \sim \exp(\beta)$, we can find mgf: $M_X(t) = E(e^{xt}) = \int_0^\infty e^{xt} \frac{1}{\beta} e^{-x/\beta} dx = \int_0^\infty \frac{1}{\beta} e^{-x(1/\beta - t)} dx$

Let $u = \frac{x}{\beta} + t$ or recognize $e^{-x(1/\beta - t)}$ kernel, so int to $1/(1/\beta - t)$.

$$= \frac{1}{\beta} \left(\frac{1}{1/\beta - t} \right) = \frac{1}{1 - \beta t}$$

Say $X \sim \text{Ga}(a, \beta)$ with $M_X(t) = \left(\frac{1}{1 - \beta t} \right)^a$. Find dist of $W = \frac{2X}{\beta}$ using mgfs.

$$M_W(t) = E[e^{2xt/\beta}] = \dots \text{ (solve the usual way now), } \int e^{2xt/\beta} f(x) dx$$

Ex: Sum of two exponentials: $X_1, X_2 \sim \text{Exp}(\beta)$. Find dist of $W = X_1 + X_2$ using mgfs.

$$M_W(t) = M_{X_1+X_2}(t) = E[e^{(X_1+X_2)t}] = E(e^{X_1 t} e^{X_2 t}) = M_{X_1}(t) M_{X_2}(t) = \left(\frac{1}{1 - \beta t} \right)^2 \sim \text{Ga}(2, \beta)$$

Joint Distributions:

Consider $f(y_1, y_2) = 4y_1 y_2$ for $0 \leq y_1 \leq 1$, $0 \leq y_2 \leq 1$.

$$P(y_1 \leq 1/2, y_2 \leq 3/4) = \int_0^{1/2} \int_0^{3/4} 4y_1 y_2 dy_1 dy_2 = \int_0^{1/2} \frac{4y_1^2}{2} y_2 dy_2 = \int_0^{1/2} \frac{1}{2} y_2 dy_2 = \frac{1}{2} \cdot \frac{3}{2}^2 = \frac{1}{4} \left(\frac{3}{4}\right)^2 = \frac{9}{64}$$

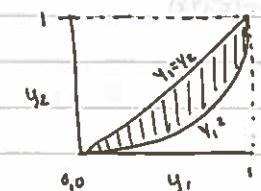
$$P(y_1 \leq 1/2 \mid y_2 \leq 3/4) = \frac{P(y_1 \leq 1/2, y_2 \leq 3/4)}{P(y_2 \leq 3/4)}$$

Numerator is $9/64$.

$$P(y_2 \leq 3/4) = \int_0^{3/4} \int_0^1 4y_1 y_2 dy_1 dy_2 = \int_0^{3/4} 4y_2 \frac{y_1^2}{2} dy_2 = \int_0^{3/4} \frac{4}{2} y_2 y_1^2 dy_2 = \frac{3}{4}^2 = 3/16$$

$$\frac{9/64}{3/16} = \frac{9}{64} \cdot \frac{4}{3} = \frac{12}{64} = \frac{3}{16}$$

$$P(y_1^2 \leq y_2 \leq y_1) = \int_0^1 \int_{y_1^2}^{y_1} 4y_1 y_2 dy_2 dy_1 = 1/6$$



Bivariate Transformations

1. Check that the two fctns being transformed are 1-1 + onto.

Say, considering a transform from $(x, y) \rightarrow (u, v)$.

• (1-1): Let $u_1, v_1 = u_2, v_2$, show implies $x_1, y_1 = x_2, y_2$. (ie $x_1 = x_2$ and $y_1 = y_2$)

• Onto: (u, v) yields a unique (x, y) .

$$2. |J| = \begin{vmatrix} \frac{\partial h_1}{\partial u} & \frac{\partial h_1}{\partial v} \\ \frac{\partial h_2}{\partial u} & \frac{\partial h_2}{\partial v} \end{vmatrix} \quad \text{where } h_1 + h_2 \text{ are like } g_1(y)$$

$$3. f(u, v) = f(h_1, h_2) |J| \quad \text{where if vars indep, } = f(h_1) f(h_2) |J|.$$

EX 1: Product of Two Beta. Let $X \sim \text{Beta}(a, b)$, $Y \sim \text{Beta}(a+b, \gamma)$. $X \perp\!\!\!\perp Y$.

Find pdf of XY . Note: $0 < x < 1$, $0 < y < 1$.

• Let $U = XY$, $V = Y$. Then $V = Y$, and $x = U/V$. $h_1 = U/V$, $h_2 = V$

• Both fctns are 1-1 + onto:

$$\cdot \text{Jacobian: } |J| = \begin{vmatrix} \frac{\partial h_1}{\partial u} & \frac{\partial h_1}{\partial v} \\ \frac{\partial h_2}{\partial u} & \frac{\partial h_2}{\partial v} \end{vmatrix} = \begin{vmatrix} 1/v & -1/v^2 \\ 0 & 1 \end{vmatrix} = 1/v$$

$$\cdot \text{pdf: } f_{UV}(u, v) = \frac{1}{B(a, b)} (u/v)^{a-1} (1-u/v)^{b-1} \cdot \frac{1}{B(a+b, \gamma)} v^{a+b-1} (1-v)^{\gamma-1} v^{-1}$$

$$= \frac{\Gamma(a+b+\gamma)}{\Gamma(a)\Gamma(b)\Gamma(\gamma)} \frac{(u/v)^{a-1} (1-u/v)^{b-1} v^{a+b-1} (1-v)^{\gamma-1}}{v^{a+b+\gamma}}$$

$$\cdot f_U = \int_u^1 \frac{\Gamma(a+b+\gamma)}{\Gamma(a)\Gamma(b)\Gamma(\gamma)} u^{a-1} v^{a-1} v^{a+b-1} (1-u/v)^{b-1} (1-v)^{\gamma-1} dv \quad \text{and so forth.}$$

Since $v = u/y$
in $[0, 1]$

Ex 2. $X \perp\!\!\!\perp Y \sim \text{Gamma}$, with $X \sim \text{Ga}(r, 1)$, $Y \sim \text{Ga}(s, 1)$. Let $z_1 = X+Y$, $z_2 = \frac{X}{X+Y}$.

Find dists of z_1, z_2 .

• Check 1-1 + onto ✓

$$\cdot z_1 = X+Y \rightarrow X = z_1 - z_2 \rightarrow z_2 = \frac{z_1 - y}{z_1 - z_2} \rightarrow z_1 z_2 = z_1 - y \rightarrow y = z_1 - z_1 z_2$$

$$\text{Then } z_1 = X \sim Z_1 - Z_1 z_2 \rightarrow X = Z_1 z_2. \text{ So } h_1(z_1, z_2) = Z_1 z_2, h_2(z_1, z_2) = Z_1 - Z_1 z_2$$

$$\cdot |J| = \begin{vmatrix} \frac{\partial h_1}{\partial z_1} & \frac{\partial h_1}{\partial z_2} \\ \frac{\partial h_2}{\partial z_1} & \frac{\partial h_2}{\partial z_2} \end{vmatrix} = \begin{vmatrix} Z_2 & Z_1 \\ (1-z_2) & -Z_1 \end{vmatrix} = -Z_1 z_2 - Z_1 (1-z_2) = Z_1$$

$$\cdot \text{pdf: } f(z_1, z_2) = \frac{1}{\Gamma(r)} (z_1 z_2)^{r-1} e^{-z_1 z_2} \cdot \frac{1}{\Gamma(s)} (z_1 - z_1 z_2)^{s-1} e^{-z_1 + z_1 z_2} \cdot z_1$$

$$= \frac{1}{\Gamma(r)\Gamma(s)} z_1^r z_2^{r-1} e^{-z_1} z_1^{s-1} (1-z_2)^{s-1} = \frac{1}{\Gamma(r+s)} z_1^{r+s-1} e^{-z_1} \cdot z_2^{r-1} (1-z_2)^{s-1}$$

$Z_1 \sim \text{Ga}(r+s, 1)$

$Z_2 \sim \text{Beta}(r, s)$.

• Convolutions: If working with the transformation $Z = X+Y$,

$$f(z) = \int_{-\infty}^z f_X(z-y) f_Y(y) dy \quad \text{where } f_X, f_Y \text{ are pdfs for } X + Y.$$

• Ex: Sum of two exponentials w. same parameter.

$X \sim \text{Exp}(\lambda)$, $Y \sim \text{Exp}(\lambda)$, $X \perp\!\!\!\perp Y$. Let $Z = X+Y$. Find pdf of Z .

$$\cdot f_X = \lambda e^{-\lambda x}, \quad f_{Y|X} = \lambda e^{-\lambda y}, \quad f_Z(z) = \int_0^z \frac{\lambda e^{-\lambda(z-y)}}{f_X(z-y)} \cdot \frac{\lambda e^{-\lambda y}}{f_{Y|X}(y)} dy$$

$$= \int_0^z \lambda^2 e^{-\lambda z} dy = z \lambda^2 e^{-\lambda z} \Big|_{y=0}^z = \lambda^2 z e^{-\lambda z} \sim \text{Ga}(2, \lambda)$$

Math Stats Flashcard Facts:

1. Delta Method: If $\sqrt{n}(Y_n - \theta) \sim N(0, \sigma^2/n)$ for $Y_n = g(X_n)$, then can approx the mean + var of Y as:

$$\sqrt{n}(g(X_n) - g(\theta)) \sim N(0, \sigma^2 [g'(\theta)]^2) \quad \text{ie } g(X_n) \xrightarrow{d} N(g(\theta), \frac{\sigma^2}{n} (g'(\theta))^2)$$

for $g'(\theta)$ differentiable, ie $\neq 0$. Can use $n=1$ for a single RV.

- Why it works: can write the 2nd order Taylor series of $g(X_n)$ around θ :

$$g(X_n) \approx g(\theta) + g'(\theta)(X_n - \theta) + \varepsilon \quad \varepsilon \rightarrow 0 \text{ as } n \rightarrow \infty$$

$$\text{So } \sqrt{n}(g(X_n) - g(\theta)) \approx g'(\theta)(X_n - \theta) \sim N(0, \text{Var}(X_n) \cdot g'(\theta)^2)$$

- Ex: $X \sim U(0, 1)$. $Y = \sqrt{X}$. Approximate $E(Y) \sim \text{Var}(Y)$. $E(X) = 1/2$, $\text{Var}(X) = 1/12$.

$$E(Y) \approx \sqrt{1/2}$$

$$\text{Var}(Y) \approx \frac{1}{12} \cdot \left[\frac{1}{2(1/2)^{1/2}} \right]^2 = \frac{1}{24}$$

$$g'(x) = \frac{d}{dx} x^{1/2} = \frac{1}{2x^{1/2}}$$

2. Formula for Taylor Series: of $f(x)$ around θ : $f(x) \approx f(\theta) + \frac{f'(\theta)}{1!}(x-\theta) + \frac{f''(\theta)}{2!}(x-\theta)^2 + \dots + e$

$$3. \text{ Limit definition of } e: \lim_{x \rightarrow \infty} \frac{(1+\frac{a}{x})^{ax}}{x} = e^{ab}$$

$$4. e^x \text{ in summation form: } \sum_{k=0}^{\infty} \frac{x^k}{k!} = e^x$$

$$5. \text{ Infinite Geometric Series: } \sum_{k=0}^{\infty} x^k = \frac{1}{1-x}$$

$$6. \text{ Binomial Summation: } \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k = (x+y)^n$$

$$7. \text{ Gamma fn: } \Gamma(x) = (x-1)!$$

$$8. \text{ Binomial Coefficient: } \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Inequalities + Convergence:

- Markov's Inequality: $P(g(X) \geq r) \leq \frac{E(g(X))}{r}$

- Chernoff's Inequality: $P(|X-\mu| > t) \leq \frac{\sigma^2}{t^2}$ or $P(|X-\mu| \geq \sigma t) \leq \frac{1}{t^2}$, $\mu = E(X)$, $\sigma^2 = \text{Var}(X)$.

- Derive from Markov's: Let $t^2 = r$, $g(x) = (x-\mu)^2/\sigma^2$. Then

$$P(g(X) \geq r) \leq \frac{E(g(X))}{r} \rightarrow P\left(\frac{(X-\mu)^2}{\sigma^2} \geq \frac{r}{\sigma^2}\right) \leq \frac{E[(X-\mu)^2]}{\sigma^2 r}$$

$$\rightarrow P(|X-\mu| \geq \sigma t) \leq \frac{\sigma^2}{\sigma^2 t^2} \text{ by multiplying LHS in P(.) by } \sigma^2 \text{ and taking sqrt.}$$

- Slootsky's Lemma: Let $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{P} c$, c is constant. Then

- $X_n + Y_n \xrightarrow{d} X + c$
- $X_n Y_n \xrightarrow{d} Xc$

Cauchy-Schwarz Inequality:

$$(x'y)^2 \leq (x'x)(y'y) \quad \text{also } E(xy)^2 \leq E(x^2)E(y^2)$$

- Stem's Lemma: $\text{Cov}(g(x), y) = E(g'(x)) \text{Cov}(x, y)$ when x, y are jointly normal.

- Central Limit Theorem: $\bar{X} \xrightarrow{d} N(\mu, \sigma^2/n)$ for $\mu = E(X)$, $\sigma^2 = \text{Var}(X)$

- WLLN: $\bar{X}_n \xrightarrow{P} \mu$, ie $\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \varepsilon) = 1$

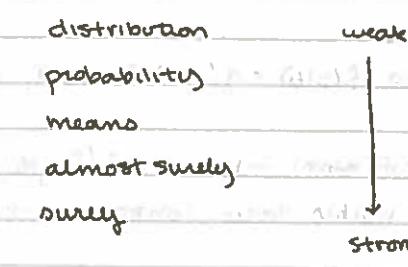
- SLLN: $\bar{X}_n \xrightarrow{a.s.} \mu$, a.s. = almost surely

- Convergence in Prob: $x_1, \dots, x_n \xrightarrow{P} x$, ie $x_n \xrightarrow{P} x$, iff

$$\lim_{n \rightarrow \infty} P(|x_n - x| \geq \varepsilon) = 0 \quad \equiv \quad \lim_{n \rightarrow \infty} P(|x_n - x| < \varepsilon) = 1$$

- Convergence in Dist: $x_n \xrightarrow{D} x$ iff $\lim_{n \rightarrow \infty} F_{x_n}(x) = F(x)$

Order of Strength of Convergence:



4. Useful Definitions + Tricks:

• Chain Rule of Probability: $p(a,b,c,d) = p(a)p(b|a)p(c|b,a)p(d|c,b,a)$

• Cov + Sample Cov: $\text{Cov}(X,Y) = E[(X-\mu_X)(Y-\mu_Y)^T] = E(XY^T) - E(X)E(Y)^T$

$$\hat{\text{Cov}}(X,Y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

• Conditional Independence: If $X \perp\!\!\!\perp Y | Z$, then

$$P(X|Z)P(Y|Z) = P(X,Y|Z)$$

$$P(X|Y,Z) = P(X|Z)$$

• Conditional Expectation + Variance: For ANY X and Y ,

$$E(X) = E(E(X|Y))$$

$$\text{Var}(X) = E(\text{Var}(X|Y)) + \text{Var}(E(X|Y)) \quad \text{"EVUE"}$$

This also holds if conditioning on something else, ie $E(X|Y) = E(E(X|Y,Z))$

Covariance is a good place to use this:

$$\text{Cov}(X,Y) = E(XY) - E(X)E(Y) = E(X \cdot E(Y|X)) - E(X)E(E(Y|X))$$

Easiest way to find $E(X)$: No transform, just $\int x f(x) dx$.

Independence: If $X \perp\!\!\!\perp Y$, then $g(x) \perp\!\!\!\perp h(y)$.

• Switching order of sums + integrals: If $h(x,\theta)$ continuous in θ , and $\sum_{i=0}^{\infty} h(x,\theta)$ converges uniformly on $[a,b]$, then

$$\sum_{i=0}^{\infty} \int_a^b h(x,\theta) d\theta = \int_a^b \sum_{i=0}^{\infty} h(\theta,x) d\theta$$

5. Finding sufficient statistics:

If can factor likelihood as $f(x|\theta) = g(T(x)|\theta) \cdot h(x)$, then $T(x)$ is a suff stat.

To find a minimal suff stat, take $f(x|\theta)/f(y|\theta)$, ratio of likelihoods. The suff stat is term which makes ratio constant wrt θ . (Doesn't have to be 1)

6. Distributional Relationships:

• Normal + Chi-sq: $Z \sim N(0,1)$ then $Z^2 \sim \chi^2_1$

$Z_1, \dots, Z_n \sim N(0,1)$ then $Z_1 + \dots + Z_n \sim \chi^2_n$

$$\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$$

• $F = \frac{U}{U+V}$ where $U \sim \chi^2_u$ + $V \sim \chi^2_v$. Then $1/F \sim F_{v,u}$

• Exp(λ) cdf: $F(x) = 1 - e^{-\lambda x}$

• Exp(λ) to Pois($t\lambda$): If $t \sim \text{Exp}(\lambda)$ is time b/w successive arrivals, let N be # of arrivals in a fixed period t . Then $N \sim \text{Pois}(t\lambda)$

• Sum of Exponentials + Gammas:

• Exp(t) is $\text{Ga}(1,\lambda)$

• Sum of $X_1, \dots, X_n \sim \text{Ga}(a_i, b) = \text{Ga}(\sum a_i, b)$

• Sum of $X_1, \dots, X_n \sim \text{Exp}(\lambda) = \text{Ga}(n, \lambda)$

• Distribution of an order stat $X_{(j)}$ $f(x_{(j)}) = \frac{n!}{(j-1)!(n-j)!} f(x) f(x)^{j-1} [1-F(x)]^{n-j}$

• t and F relationship: If $X \sim t_n$, $X^2 \sim F(1,n)$

• Uniform cdf for $U(a,b)$: $\frac{x-a}{b-a}$

• Sum of Indep Poissons: $X \sim \text{Pois}(\lambda)$, $Y \sim \text{Pois}(\theta)$, then $X+Y \sim \text{Pois}(\lambda+\theta)$

• Ratio of 2 std normals $X/Y \sim \text{Cauchy}(0,1)$

• Sum of Chi-Squareds: χ^2 , wrt df = sum of dfs

• $T = Z/\sqrt{U/n}$ for $T \sim t_n$, $Z \sim N(0,1)$, $U \sim \chi^2_n$

• Dist of $\frac{(\bar{X}-\mu)}{S/\sqrt{n}} \sim t_{n-1}$

• Inverse Gamma mean + cov: $X \sim \text{IG}(a,b)$. $E(X) = \frac{b}{a-1}$ $\text{Var}(X) = \frac{b^2}{(a-1)^2(a-2)}$

7. MGFA

- Thm 2.3.11 - Moments defining cdfs.
 - If $X \sim Y$ have bounded support, then $X=Y$ iff $E(X^r) = E(Y^r)$, w all moments same.
 - If mgfs exist & are equal in neighborhood of $t=0$, mgfs are same.

If $X \perp\!\!\!\perp Y$, $Z = X+Y$, then $M_Z(t) = M_X(t)M_Y(t)$

Mgf of $aX+b$: $E(e^{(ax+bt)}) = e^{bt}E(e^{ax}) = e^{bt}M_X(at)$

Mgf of \bar{X} : $M_{\bar{X}}(t) = [M_X(t/n)]^n$

8. Exponential families: $f(x|\theta) = h(x) \cdot c(\theta) \exp\left[\sum_{i=1}^k w_i(\theta)t_i(x)\right]$

If $k=\text{length}(\theta)$, is full exp family.

If $k < \text{length}(\theta)$, is curved exp family.

$T(x) = \left[\sum_{i=1}^k t_1(x_i), \dots, \sum_{i=1}^k t_k(x_i) \right]$ is a sufficient statistic.

If full exp family, $T(x)$ is a complete sufficient statistic.

NOT exp families: t, F, Cauchy, Binom if n unknown,

NBinom if r unknown, Weibull if shape unknown.