

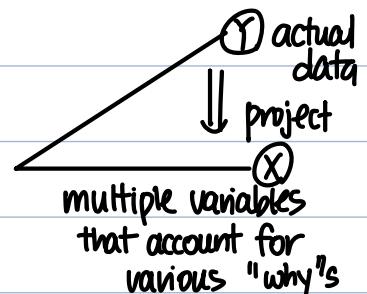
METRICS FLOW SO FAR

* Projection Theorem

- shortest distance between a "guess" and "actual data"
- way to minimize error / expected loss (risk)

* OLS (Bryan → "Least Sq. Fit") sample

- minimizes **Sample sq. error**
- use proj thm for LSF



* Linear Regression pop.

- generalization : pop. analogue of LSF
 - linear regression = some fnc of joint dist. of X & Y

↳ the only feature of pop. that matters for linear reg.

$\ddot{X} \cdot (cf)$ identification → population : whether mapped well into structural parameters
Bryan emphasizes diff btw pop. & sample!
statistics → sample : how subsets (sample) approach population

* Decision Theory

- introducing Loss, Risk, Decision rule for the first time
 - risk: trade-off btw bias & variance
 - "admissibility" notion becomes important
- Point estimation of θ as a Bayesian
 - ⇒ given a prior & likelihood, what is the posterior?
 - LB, ML, LFB are simply all decision rules
 - ⇒ what type of decision are we trying to make when we observe $X=x$?
 - Avg risk : posterior mean = $\arg \min$ Posterior Exp. Loss
avg w/o prior

* Prediction Problem

- trying to predict values of Y given X
 - = finding the conditional mean of $Y|X$
 - if we knew the joint distribution, wouldn't be a problem
 - but all we have is samples → LSF of Y
 - ✓ consistency
 - ✓ asymptotic normality
 - but $\bar{N} \rightarrow$ inadmissible
- trying to reduce the volatility of our fit

Suppose $Y|X$ is well described by an 8th-order polynomial

\Rightarrow LSF: 8th-order polynomial

{ JS est

Now we multiply unbiased LSF coefficients by 0.7 to reduce volatility

\Rightarrow then BIASED is GENERATED

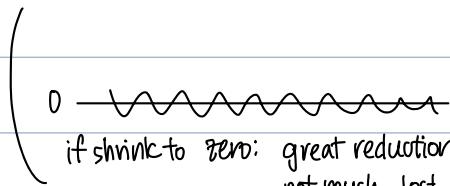
(b/c previously unbiased coeff. now biased towards 0)

Oracle est.

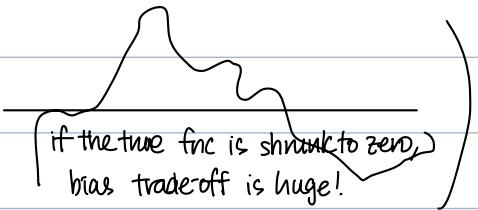
instead of shrinking all of the coefficients uniformly,

{ coeff $\neq 0 \rightarrow$ shrink to zero

coeff far from 0 \rightarrow leave them to capture the wiggly parts



if shrink to zero: great reduction in var,
not much lost for bias



Econ 240A - (2)

10/21 Monday

* Projection Theorem

- Vector space \mathcal{H} : an element of \mathcal{H} is a vector

[ex] \mathbb{R}^N : Euclidean space

→ usu. (finite) sample space

L^2 : (functions of) random variables (r.v.) w/ finite variance

→ usu. population space

null vector $\underline{0}$

$$\left\{ \begin{array}{l} \mathbb{R}^N : \underline{0} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \\ L^2 : \underline{0} = \text{degenerate r.v. which always equals zero} \end{array} \right.$$

- Hilbert Space: Vector Space + Inner Product

• Inner product

- [Def] $\langle \cdot, \cdot \rangle : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$

$h_1 \in \mathcal{H}, h_2 \in \mathcal{H} \quad \langle h_1, h_2 \rangle$

- [Properties] (i) Bi-linearity : $\langle ah_1 + bh_2, ch_3 + dh_4 \rangle$

$$= ac \langle h_1, h_3 \rangle + ad \langle h_1, h_4 \rangle + bc \langle h_2, h_3 \rangle + bd \langle h_2, h_4 \rangle$$

(ii) Symmetry : $\langle h_1, h_2 \rangle = \langle h_2, h_1 \rangle$

(iii) Positivity : $\langle h_1, h_1 \rangle \geq 0$ (inner product with itself : strictly positive if $h_1 \neq \underline{0}$)

[Examples]

(1) Euclidean space \mathbb{R}^N :

$\mathbf{X} = (X_1, \dots, X_N)', \mathbf{Y} = (Y_1, \dots, Y_N)'$ (X_i : years of completed schooling, Y_i = log earnings)

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \frac{\mathbf{X}' \mathbf{Y}}{N} = \frac{1}{N} \sum_{i=1}^N X_i Y_i \quad \begin{array}{l} \text{check it satisfies} \\ \text{i) bi-linearity, ii) symmetry} \end{array}$$

[cf] dot product

not by definition, but for the sake of statistical convenience

(2) L^2 space :

\mathbf{X} : years of completed schooling for a random draw from the target population

\mathbf{Y} : earnings from random draw

$$(X, Y) \sim f_{x,y} \text{ (population)}$$

$$\langle X, Y \rangle = \mathbb{E}[XY]$$

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \quad \begin{array}{l} \text{common to assume} \\ \mathbb{E}[r.v.] = 0 \text{ in } L^2 \end{array}$$

[cf] covariance inner product

- Norm : associated w/ any Inner Product is a norm
(measure of distance from the null vector)

[Def] $\|h\| = \sqrt{\langle h, h \rangle}$

[Properties] (i) $\|h\| = 0$ iff $h = 0$ (null vector)

(ii) $\|ah\| = |a| \cdot \|h\|$

(iii) Triangle Inequality (T. I.) : $\|h_1 + h_2\| \leq \|h_1\| + \|h_2\|$

[PF] Use Lemma 1 to prove T. I.

► Lemma 1 : Cauchy-Schwarz Inequality

if $h_1, h_2 \in \mathbb{H}$, $|\langle h_1, h_2 \rangle| \leq \|h_1\| \cdot \|h_2\|$

w/ equality iff $h_1 = \alpha h_2$ or $h_2 = 0$

$0 \leq \langle h_1 - \alpha h_2, h_1 - \alpha h_2 \rangle$ (positivity)

$= \langle h_1, h_1 \rangle - \alpha \langle h_1, h_2 \rangle - \alpha \langle h_2, h_1 \rangle + \alpha^2 \langle h_2, h_2 \rangle$ (bi-linearity)

$= \|h_1\|^2 - 2\alpha \langle h_1, h_2 \rangle + \alpha^2 \|h_2\|^2$ (symmetry, def. of norm)

Set $\alpha = \frac{\langle h_1, h_2 \rangle}{\|h_2\|^2}$

$0 \leq \|h_1\|^2 - \frac{\langle h_1, h_2 \rangle^2}{\|h_2\|^2} \Leftrightarrow \langle h_1, h_2 \rangle^2 \leq \|h_1\|^2 \cdot \|h_2\|^2$

[Ex] use CS inequality to show that correlation coefficient : $-1 \leq \rho \leq 1$.

(Two cases) $\rho = \frac{\text{Cov}(X, Y)}{\sqrt{V(X)} \cdot \sqrt{V(Y)}}$ ← use L^2 , inner product

sample correlation coeff. ← use \mathbb{R}^N , inner product

► [Lemma 2] $\|h_1 + h_2\| \leq \|h_1\| + \|h_2\|$

$$\|h_1 + h_2\|^2 \leq \langle h_1 + h_2, h_1 + h_2 \rangle$$

$$= \|h_1\|^2 + 2\langle h_1, h_2 \rangle + \|h_2\|^2$$

$$\leq \|h_1\|^2 + 2 |\langle h_1, h_2 \rangle| + \|h_2\|^2$$

CS $\leq \|h_1\|^2 + 2 \|h_1\| \cdot \|h_2\| + \|h_2\|^2$

$$= (\|h_1\| + \|h_2\|)^2 \blacksquare$$

- Orthogonal : X, Y are orthogonal ($X \perp Y$) if their inner product is zero.

[cf] $X \perp Y$: X is independent of Y

► [Thm 1 : Pythagorean Thm]

If $\mathbf{h}_1 \perp \mathbf{h}_2$, then $\|\mathbf{h}_1\|^2 + \|\mathbf{h}_2\|^2 = \|\mathbf{h}_1 + \mathbf{h}_2\|^2$

[Ex] \mathbb{R}^2 + connect it to geometric property $\triangle : \mathbb{R}^2$

• Other things to think about

- $\|\mathbf{h}_1 - \mathbf{h}_2\|$ is a measure of distance btw \mathbf{h}_1 and \mathbf{h}_2

$$[\mathbb{R}^2]: \mathbf{h}_1 = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \mathbf{h}_2 = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix}$$

$$\|\mathbf{h}_1 - \mathbf{h}_2\| = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2}$$

$$[L^2]: \mathbf{h}_1 = X, \mathbf{h}_2 = Y$$

$$\|\mathbf{h}_1 - \mathbf{h}_2\| = E[(X - Y)^2]^{\frac{1}{2}} = \text{Root Mean Squared Error (RMSE)}$$

• Projection Thm

Let \mathcal{L} be some linear subspace of \mathcal{H} .

Projection of $\mathbf{Y} \in \mathcal{H}$ onto the subspace \mathcal{L} .

- Projection operator : $\Pi(\cdot | \mathcal{L}) : \mathcal{H} \rightarrow \mathcal{L}$

where $\Pi(\mathbf{Y} | \mathcal{L})$ is the element of $\hat{\mathbf{Y}} \in \mathcal{L}$

that achieves $\min_{\mathbf{Y} \in \mathcal{L}} \|\mathbf{Y} - \hat{\mathbf{Y}}\|$

$$\bullet (\text{Ex: } \mathbb{R}^N) \quad \left[\begin{array}{c} \mathbf{Y} \\ \mathbf{1}_{N \times 1}, \mathbf{X} \\ \mathbf{1}_{N \times 1} \end{array} , \quad \mathcal{L} = \text{linear span of } \mathbf{1}_{N \times 1} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \mathbf{1}_N, \mathbf{X} \right]$$

↳ (i.e. vectors of the form $\alpha \mathbf{1} + \beta \mathbf{X}$)

$$\min_{(\alpha, \beta) \in \mathbb{R}^2} \|\mathbf{Y} - \alpha \mathbf{1} - \beta \mathbf{X}\|^2$$

$$= \min_{(\alpha, \beta) \in \mathbb{R}^2} \frac{1}{N} \sum_{i=1}^N (Y_i - \alpha - \beta X_i)^2$$

→ corresponds to finding the OLS fit of \mathbf{Y} onto a constant and \mathbf{X} .

$$\bullet (\text{Ex: } L^2) \quad \left[(\mathbf{X}, \mathbf{Y}) \sim F_{\mathbf{X}, \mathbf{Y}}, \mathcal{L} \text{ all linear functions of } \mathbf{X}; \alpha + \beta \mathbf{X} \right]$$

$$\min_{(\alpha, \beta) \in \mathbb{R}^2} \|\mathbf{Y} - \alpha - \beta \mathbf{X}\|^2 = \min_{(\alpha, \beta) \in \mathbb{R}^2} E[(Y - \alpha - \beta X)^2]$$

Mean Squared Error

population
linear regression
function

► [Thm 2: Projection Thm]

Let \mathcal{H} be a vector space with an inner product and associated norm

and \mathcal{L} a subspace of \mathcal{H}

Then for γ an arbitrary element of \mathcal{L} ,

if there exists a vector $\hat{\gamma} \in \mathcal{L}$ s.t.

$$\|\gamma - \hat{\gamma}\| \leq \|\gamma - \tilde{\gamma}\| \quad \text{for } \forall \tilde{\gamma} \in \mathcal{L}$$

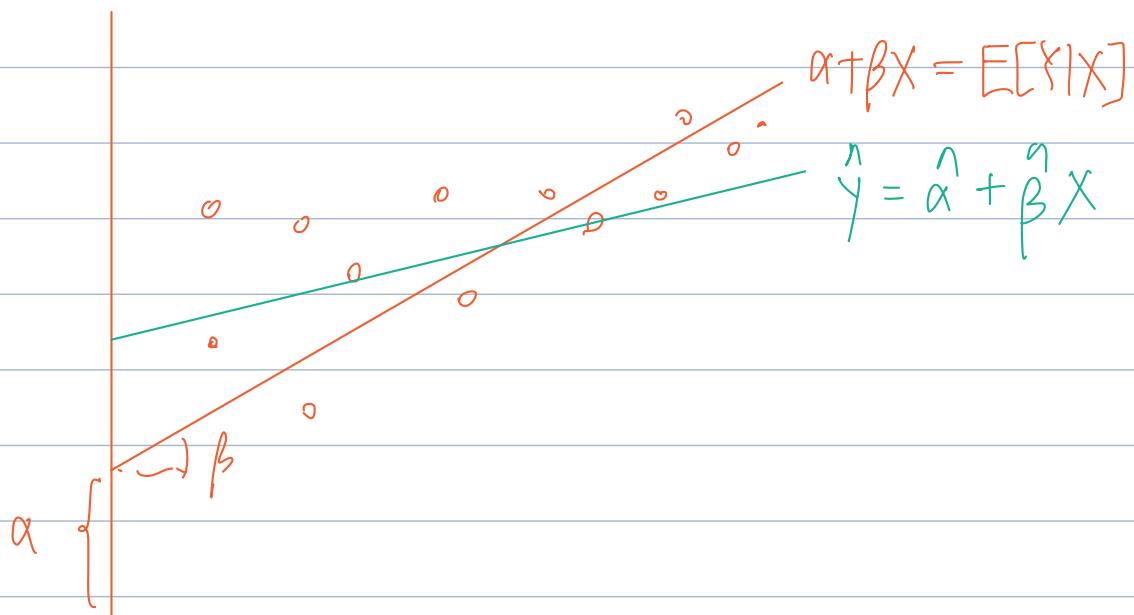
then (i) $\hat{\gamma} = \underset{\text{projection}}{\pi}(\gamma | \mathcal{L})$ is unique

(ii) nec & suff condition for $\hat{\gamma}$ to be the solution:

$$\langle \gamma - \hat{\gamma}, \tilde{\gamma} \rangle = 0 \quad \text{for } \forall \tilde{\gamma} \in \mathcal{L}$$

$$(\text{or } \gamma - \pi(\gamma | \mathcal{L}) \perp \tilde{\gamma} \quad \text{for } \forall \tilde{\gamma} \in \mathcal{L})$$

↑
a lot of the power of regression
comes from the orthogonality



* Projection Thm (cont'd)

* Recap

- [\mathcal{H} : vector space w/ inner product + norm]
- [\mathcal{L} : subspace, $Y \in \mathcal{H}$, $\hat{Y} \in \mathcal{L}$]
- If $\exists \hat{Y} \in \mathcal{L}$, $\|Y - \hat{Y}\| \leq \|Y - \tilde{Y}\| \quad \forall \tilde{Y} \in \mathcal{L}$

- { (i) $\hat{Y} = \Pi(Y|\mathcal{L})$ is unique
(ii) nec + suff: $\langle Y - \hat{Y}, \tilde{Y} \rangle = 0 \quad \forall \tilde{Y} \in \mathcal{L}$

\hookrightarrow $\begin{array}{c} Y - \Pi(Y|\mathcal{L}) \\ \text{target} \quad \text{approximated} \\ \text{approx. error} \end{array} \perp Y : \text{orthogonal}$

[PF] in the notes (* forthcoming)

* Properties of Projections

(i) Projections are Linear operators

[Set-up] $X = \underbrace{\Pi(X|\mathcal{L})}_{\text{vector}} + \underbrace{Ux}_{\text{projection onto } \mathcal{L}}, \quad Ux \perp \mathcal{L}, \quad Ux \perp \tilde{Y} \quad \forall \tilde{Y} \in \mathcal{L}$

all vectors can be expressed in this form

$$Y = \Pi(Y|\mathcal{L}) + Uy, \quad Uy \perp \mathcal{L}$$

$$\textcircled{1} \quad aX + bY = a\Pi(X|\mathcal{L}) + b\Pi(Y|\mathcal{L}) + aUx + bUy \quad (a, b) \text{ scalar}$$

Linear combination of two vectors:

$$\textcircled{2} \quad W \in \mathcal{L}$$

• by bilinearity $\Rightarrow \langle aUx + bUy, w \rangle = a \underbrace{\langle Ux, w \rangle}_{=0} + b \underbrace{\langle Uy, w \rangle}_{=0} \quad w \in \mathcal{L}$

$\boxed{\Pi(aX + bY) = a\Pi(X|\mathcal{L}) + b\Pi(Y|\mathcal{L})} \Rightarrow \text{Projections are linear operators}$

[Following Property] $E[aX + bY | W] = aE[X|W] + bE[Y|W]$

$$= a \int x \cdot f_{X|W}(x|W) dx + b \int y \cdot f_{Y|W}(y|W) dy$$

* (notation) $Y \sim F_Y$: generic random draw $\rightarrow Y_i$: i-th random draw

$E[Y] = E[Y_i] = E[Y_1] = E[Y_{3,61}]$: under sample random sampling

$Y \in \mathbb{Y}$: support of r.v. Y

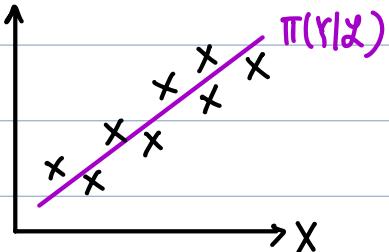
y : specific element of \mathbb{Y} (sample space)

$$\begin{cases} \mathbb{E}[Y|X] = g(x) & (\text{r.v.} : \text{function of r.v.}) \\ \mathbb{E}[Y|X=x] = g(x) & \text{for } x \in \mathbb{X} \\ \mathbb{E}[Y|x] & (\text{Goldberger notation}) \end{cases}$$

$$(X, Y) \sim F_{X,Y}$$

(ii) Idempotency

Q. "The projection of a projection is itself." (T/F question possible) $\Rightarrow \text{True}$



[WTS] $\pi(\pi(Y|Z)|Z) = \pi(Y|Z)$

$$0 = \langle \pi(\pi(Y|Z)|Z) - \pi(Y|Z), \tilde{Y} \rangle \quad \tilde{Y} \in Z \quad \text{by NSC of PT}$$

\downarrow
 $Y - Y$ target

$$= \underbrace{\langle Y - \pi(Y|Z), \tilde{Y} \rangle}_{=0 \text{ P.T.}} - \underbrace{\langle Y - \pi(\pi(Y|Z)|Z), \tilde{Y} \rangle}_{\text{NSC of PT}} \quad \text{by bi-linearity, uniqueness of PT.}$$

[Norm Reducing]

Let 1 denote the constant vector

$$\text{(a) } \mathcal{H} = \mathbb{R}^N \text{ Euclidean space ; } 1 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = 1_N$$

$$\text{(b) } \mathcal{H} = L^+ \quad ; \quad 1 = 1$$

$$\text{(a) } \pi(Y|1) = \frac{1}{N} \cdot \sum_{i=1}^N Y_i = \bar{Y} \quad \langle X, Y \rangle = \frac{1}{N} \sum_i X_i Y_i = \frac{1}{N} X' Y$$

$$\text{(b) } \pi(Y|1) = E[Y] \quad \langle X, Y \rangle = E[XY]$$

$$\text{(a) } \underbrace{\|Y - \pi(Y|1)\|^2}_{\text{projection error}} = \frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})^2$$

$$\text{(b) } \|Y - \pi(Y|1)\|^2 = V(Y) \quad \text{Var}(Y)$$

(iii) [Lemma 3: Analysis of Variance]

If the linear subspace Z contains a constant vector,

$$\|Y - \pi(Y|1)\|^2 = \|Y - \pi(Y|Z)\|^2 + \|\pi(Y|Z)\|^2 - \|\pi(Y|1)\|^2$$

projection error from
projecting Y onto a constant vector
contains 1

$$\text{(cf)} \quad \text{TSS} = \text{RSS} + \text{ESS}$$

$$V(Y) = \mathbb{E}[V(Y|X)] + V(\mathbb{E}[Y|X])$$

$$\begin{aligned}
 - [\text{PF}] \quad \|Y - \pi(Y|Z)\|^2 &= \|Y - \pi(Y|1) - [\pi(Y|Z) - \pi(Y|1)]\|^2 \\
 &\stackrel{\substack{\text{by def. of norm } \|Y\| = \langle Y, Y \rangle^{\frac{1}{2}} \\ \& \text{by bi-linearity}}}{=} \|Y - \pi(Y|1)\|^2 - 2 \langle Y - \pi(Y|1), \pi(Y|Z) - \pi(Y|1) \rangle + \|\pi(Y|Z) - \pi(Y|1)\|^2
 \end{aligned}$$

(Show: $(*) = \|\pi(Y|Z) - \pi(Y|1)\|^2$ HW)

$$= \|Y - \pi(Y|1)\|^2 - \|\pi(Y|Z) - \pi(Y|1)\|^2$$

[Implication] Projections are norm-reducing:

$$\|Y - \pi(Y|Z)\| \leq \|Y - \pi(Y|1)\|$$

\Rightarrow [Recap]

- (i) Projections are linear operators
- (ii) Idempotency
- (iii) ANOVA, Projections are norm-reducing

* OLS

- Set-up
 - $Y_{N \times 1}$: log earnings for a random sample of N adult males
 - $X_{N \times K}$: covariates — e.g. 1, Yrs Sch, AFQT, demographic controls ...
 - Euclidean space

Assume : columns of $X = \begin{pmatrix} 1 & \text{Yrs Sch}_1 & \text{AFQT}_1 \\ \vdots & \vdots & \vdots \\ 1 & \text{Yrs Sch}_N & \text{AFQT}_N \end{pmatrix} \rightarrow$ assume linearly independent
 $\therefore \text{rank}(X) = K$

$\mathcal{L} = C(X) := \{ X\beta : \beta \text{ is a } K \times 1 \text{ vector of real numbers} \}$

$\text{column space of } X \quad \underbrace{(N \times K)}_{N \times 1} \quad \underbrace{(K \times 1)}_{N \times K}$

- Find: Projection of $Y_{N \times 1}$ onto the column space of $X_{N \times K}$ (i.e. onto $\mathcal{L} = C(X)$)

- Use the Projection Thm:

$$\begin{aligned}
 (\text{NSC}) \quad & \langle Y - \pi(Y|Z), X\beta \rangle = 0 \\
 & = X\hat{\beta} \quad \uparrow \text{an element of } \mathcal{L} = C(X)
 \end{aligned}$$

$$(*) \quad \langle Y - X\hat{\beta}, X\beta \rangle = 0$$

↑ projection error

$$(\#) \quad \frac{1}{N} \sum_i^n (Y_i - X_i' \hat{\beta}) X_i' \beta = 0$$

$$\mathbb{X}_{N \times K} \mathbb{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1K-1} \\ 1 & X_{21} & \cdots & X_{2K-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{N1} & \cdots & X_{NK-1} \end{pmatrix} \quad \mathbb{X}_i = \begin{pmatrix} 1 \\ X_{i1} \\ \vdots \\ X_{iK-1} \end{pmatrix} \Rightarrow \mathbb{X} = \begin{pmatrix} X_1' \\ X_2' \\ \vdots \\ X_N' \end{pmatrix}$$

$$\mathbb{X} \beta = \begin{pmatrix} X_1' \\ X_2' \\ \vdots \\ X_N' \end{pmatrix} \beta = \begin{pmatrix} X_1' \beta \\ X_2' \beta \\ \vdots \\ X_N' \beta \end{pmatrix}$$

~~$$X_i' \beta = (1, X_{i1}, \dots, X_{iK-1}) \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{K-1} \end{pmatrix} \rightarrow (X_{i1}, \dots, X_{iK}) \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix} \text{ where } X_1 = 1 \text{ (constant)}$$~~

$$= \sum_{k=1}^K X_{ik} \beta_k$$

$$\Rightarrow (\#) \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \beta_k X_{ik} (Y_i - X_i' \hat{\beta}) = 0 \quad \text{for } \forall \beta \in \mathbb{R}^K \quad (\text{NSC or PT})$$

- Next, set $\begin{cases} \beta_k = 1 \\ \beta_j = 0 \text{ for } j \neq k \end{cases}$

$$\Rightarrow \frac{1}{N} \sum_{i=1}^N X_{ik} (Y_i - X_i' \hat{\beta}) = 0 \quad \text{for } k=1, \dots, K$$

- K equations \times K unknowns : $\hat{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix}$

- Stacking the K equations : $\frac{1}{N} \sum_{i=1}^N X_i (Y_i - X_i' \hat{\beta}) = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} \quad [\text{FOC}]$

$$= \frac{1}{N} \sum_{i=1}^N X_i Y_i - \left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right] \hat{\beta} = 0$$

$$\Rightarrow \underbrace{\left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right]}_{K \times K} \hat{\beta} = \underbrace{\left[\frac{1}{N} \sum_{i=1}^N X_i Y_i \right]}_{K \times 1}$$

$$\Rightarrow \hat{\beta} = \left[\frac{1}{N} \sum_{i=1}^N X_i X_i' \right]^{-1} \cdot \left[\frac{1}{N} \sum_{i=1}^N X_i Y_i \right]$$

$$= (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' \mathbb{Y} \quad \text{where } \mathbb{X}_{N \times K} \quad \mathbb{Y}_{N \times 1}$$

$\hat{\beta}$ unique
(in this class)
 β unique?

Hence, the projection of $\mathbb{Y}_{N \times 1}$ onto the column space of $\mathbb{X}_{N \times K}$ is :

$$\Pi(\mathbb{Y} | \mathbb{X}) = \mathbb{X} \hat{\beta} = \underbrace{\mathbb{X} (\mathbb{X}' \mathbb{X})^{-1} \mathbb{X}' \mathbb{Y}}_{\hat{\beta}} : \text{the Ordinary Least Squares fit of } Y_l \text{ onto } X_l \text{ for } l=1, \dots, N$$

- More on OLS coefficient $\hat{\beta}$

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^k} \frac{1}{N} \sum (Y_i - X_i' \beta)^2$$

$$= \arg \min_{\beta \in \mathbb{R}^k} \| Y - X' \beta \|_2^2$$

since $\| Y - \hat{Y} \| \leq \| Y - \tilde{Y} \| \quad \forall \tilde{Y} \in C(X)$

$$\| Y - X \hat{\beta} \| \leq \| Y - X \beta \| \quad \forall \beta \in \mathbb{R}^k$$

* Linear Regression

• Set-up

$$\begin{cases} X, Y \sim F_{X,Y} \\ d = \{X'\beta - \beta \text{ is a } (K \times 1) \text{ vector of reals}\} \\ \|X\|_d = \left[\sum_{k=1}^K X_k^2 \right]^{1/2} : \text{Euclidean norm} \end{cases}$$

• Assumption on $F_{X,Y}$ ("regularity conditions") (\rightarrow conditions on moments)

$$(A.1) \begin{cases} \text{(i)} \mathbb{E}[Y^2] < \infty & \rightarrow \text{vector in } L^2 \text{ space} \\ \text{(ii)} \mathbb{E}[\|X\|_d^2] < \infty \\ \text{(iii)} \mathbb{E}[(\alpha'X)^2] > 0 & (\alpha \text{ is a non-zero vector of constants}) \end{cases}$$

\hookrightarrow implies: $\mathbb{E}[(\alpha'X)^2] = \alpha' \underbrace{\mathbb{E}[XX']}_{\text{positive def.}} \alpha > 0$

\Rightarrow ensures uniqueness of coefficient vector indexing the projection
 $(\because \mathbb{E}[XX']^{-1} \text{ exists})$

$\Rightarrow \mathbb{E}[XX']^{-1}_{(K \times K)}$ is well-defined (*)

• Goal: Find Projection of Y onto \mathcal{L}

$$\min_{\beta \in \mathbb{R}^K} \|Y - X'\beta\|^2 = \min_{\beta \in \mathbb{R}^K} \mathbb{E}[(Y - X'\beta)^2]$$

↑ cov.
inner product

mean squared error (mse)
(heteroskedasticity)

in L^2 : $\langle X, Y \rangle = \mathbb{E}[XY]$

$\|Y\| = \mathbb{E}[Y^2]^{\frac{1}{2}}$

- Use the NSC of Projection Thm:

$$\begin{cases} \langle Y - X'\beta_0, X'\beta \rangle = 0 & \forall \beta \in \mathbb{R}^K \\ \mathbb{E}[(Y - X'\beta_0) \cdot \underbrace{X'\beta}_{\text{r.v.}}] = 0 & \Leftrightarrow \mathbb{E}[(Y - X'\beta_0) X] = 0 \end{cases}$$

where $\langle Y - \hat{Y}, \tilde{Y} \rangle = 0$ for $\tilde{Y} \in \mathcal{L}$

$X'\beta_0$ $X'\beta$

\therefore zero covariance condition

$$\begin{aligned} \mathbb{E}[XY] - \mathbb{E}[XX'] \beta_0 &= 0 && (\text{by (*) from A1-iii}) \\ \mathbb{E}[XY] &= \mathbb{E}[XX'] \beta_0 \\ \Rightarrow \beta_0 &= \mathbb{E}[XX']^{-1} \mathbb{E}[XY] \end{aligned}$$

Linear Regression: $\pi(Y|\mathcal{L}) = \mathbb{E}^*[Y|\mathcal{X}] = X'\beta_0$

linear reg. of Y on X
(not expectation operator)

$\beta_0 = \mathbb{E}[XX']^{-1} \cdot \mathbb{E}[XY]$

\ddot{X} (Note) $X \sim F_X$: $\left\{ \begin{array}{l} X \text{ is a random draw from } F_X \\ x \text{ is a specific configuration} \\ X \text{ is the support of r.v. } X \end{array} \right.$

* Decision Theory (Laplace)

- Notation / Introduction

- $\{P_\theta : \theta \in \Theta\}$
 - \uparrow statistical model
 - \uparrow parameter
 - \uparrow parameter space (all states of nature)
- Z_i : a r.v. w/ distribution P_θ for some $\theta \in \Theta$.
- Z : sample space
- Experiment : the process that generates Z
- After observing Z , the econometrician takes an action a in some action space A .
- Our payoff depends on our action a , and the "true" state-of-nature θ .
- Loss function $L(\theta, a)$ is real-valued on $A \times \Theta$.
- Let $d: Z \rightarrow A$: after observing $Z = z$, we choose action $d(z) \in A$, $d(z)$ = decision rule
- Risk (-Expected Utility) : $R(\theta, d) = E_\theta [L(\theta, d(Z))]$

decision-rule:
fully-specified plan of
what to do for $\forall z$

- Estimate a Bernoulli Success problem

- (i) Set-up

$Z = (z_1, \dots, z_N)'$: a sequence of independent Bernoulli draws

w/ an unknown success probability of θ .

- (ii) [Goal] Construct a point-estimate of θ based on Z .

$$\begin{aligned} \text{expected loss} \Rightarrow \text{risk} & \quad L(\theta, a) = (\theta - a)^2 & \rightarrow \text{cannot choose } a \text{ to min } L \text{ since } \theta \text{ unknown} \\ & \quad R(\theta, d) = E_\theta [L(\theta, d(Z))] \end{aligned}$$

$$\begin{aligned} &= E_\theta [(\theta - d(Z))^2] \\ &= E_\theta \{ [(\theta - E_\theta [d(Z)]) - (d(Z) - E_\theta [d(Z)])]^2 \} \\ &= (\theta - E_\theta [d(Z)])^2 - 2[\theta - E_\theta [d(Z)]][E_\theta [d(Z)] - E_\theta [d(Z)]] = 0 \\ &\quad + E_\theta \{ (d(Z) - E_\theta [d(Z)])^2 \} \end{aligned}$$

$$\Rightarrow R(\theta, d) = \frac{(\theta - \mathbb{E}_\theta[d(Z)])^2}{\text{bias squared}} + \frac{V_\theta(d(Z))}{\text{variance}}$$

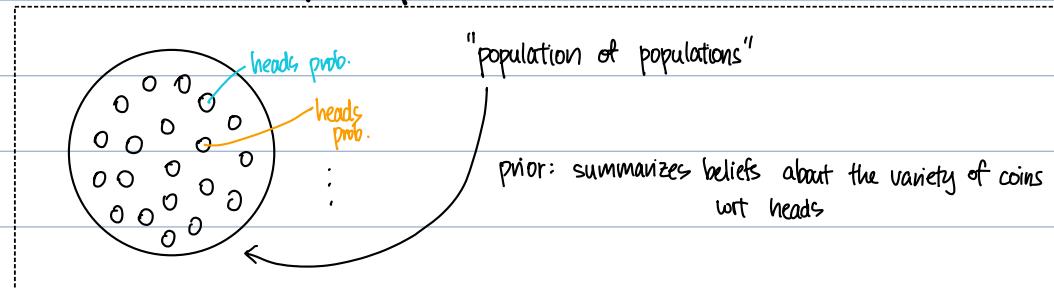
: low-risk decision rules involve a variance trade-off

(iii) [Goal] Find a decision rule $d(Z)$ that "predicts" θ "well"

\Rightarrow use Bayes' Rule

Bayes' Rule

Prior distribution: summarizes our pre-experiment beliefs about θ



$$\Rightarrow \text{Prior : } \pi(\theta) = \begin{cases} 1 & 0 < \theta < 1 \\ 0 & \text{otherwise} \end{cases}$$

\leftarrow controversial b/c consequential

$$-\text{Likelihood: } l(z_i | \theta) = \prod_{i=1}^N \theta^{z_i} (1-\theta)^{1-z_i}$$

$$= \theta^{S_N} (1-\theta)^{N-S_N}$$

$(S_N = \sum_{i=1}^N z_i : \text{sum of success of heads})$

$$(Z_i = (0, 1, 0, 0, \dots, 1))$$

Posterior Distribution $(Z = z = (0, 1, 1, 0, 0, 0, \dots))$

$$\pi(\theta | z) = \frac{l(z|\theta) \cdot \pi(\theta)}{\int_{t=0}^{t=1} l(z|t) \pi(t) dt} = \frac{\theta^{S_N} (1-\theta)^{N-S_N}}{\int_{t=0}^{t=1} t^{S_N} (1-t)^{N-S_N} dt}$$

Beta fnc: $\beta(S_N+1, N-S_N+1)$

$$(\because \beta(\alpha_1, \alpha_2) = \int_{t=0}^{t=1} t^{\alpha_1-1} (1-t)^{\alpha_2-1} dt)$$

$$\Rightarrow \theta | Z = z \sim \text{beta}(S_N+1, N-S_N+1) : \text{posterior dist.}$$

* Note: median of a beta (α_1, α_2) distributed r.v. = $\frac{\alpha_1 - 1/2}{\alpha_1 + \alpha_2 - 2/2}$

I believe there is a 50% chance that θ is less than $\frac{S_N + 2/3}{N + 4/3}$

or { subjective prob. statements
degree-of-belief statements

- Posterior mean & variance :

$$\left\{ \begin{array}{l} E[\theta | Z_1 = z] = \frac{s_N + 1}{N + 2} \\ V[\theta | Z_1 = z] = \frac{(s_N + 1)(N - s_N + 1)}{(N + 2)^2(N + 3)} \end{array} \right.$$

- Decision rule of "best guess" for θ

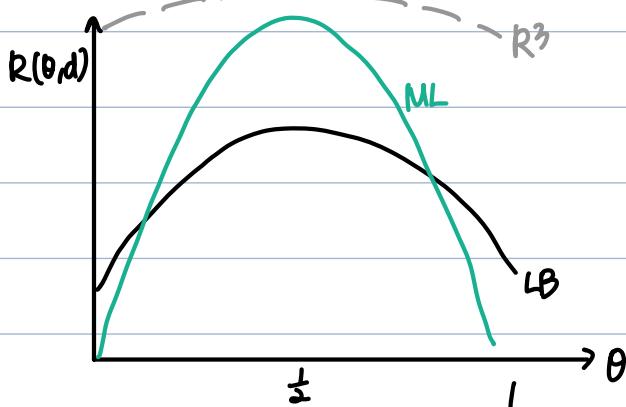
$$① d_{LB}(Z_1) = \frac{s_N + 1}{N + 2}$$

↑
fix which
depends
on Z_1

$$\left\{ \begin{array}{l} E_\theta[d_{LB}(Z_1)] = \frac{N\theta + 1}{N + 2} \\ V_\theta[d_{LB}(Z_1)] = \frac{N\theta(1-\theta)}{(N+2)^2} \end{array} \right.$$

$$\Rightarrow R(\theta, d_{LB}) = \underbrace{\left(\frac{2\theta - 1}{N+2} \right)^2}_{\text{bias sq.}} + \underbrace{\frac{N\theta(1-\theta)}{(N+2)^2}}_{\text{var.}}$$

$$② d_{ML}(Z_1) = \frac{s_N}{N}, \quad R(\theta, d_{ML}) = \frac{\theta(1-\theta)}{N}$$



: not admissible \rightarrow worse than another rule in every state of world (strictly dominated)

\rightarrow both rules are admissible (cross!)

if true value $\hat{\theta} = 0.5$: ML

$\hat{\theta} = 0$ or 1 : LB

(\because neither dominates the other)

- Another example:

(Prior ①) $\theta \sim \text{Beta}(\alpha_1, \alpha_2)$

$$\Rightarrow \theta | Z_1 = z \sim \text{Beta}(s_N + \alpha_1, N - s_N - \alpha_2)$$

Imagine that you observed $\alpha_1 + \alpha_2$ coin flips of the coin in hand already and of these α_1 were heads

$\Rightarrow \alpha_1, \alpha_2$: pre-data "pseudo-counts"

$$\theta \sim U[0,1] = \text{Beta}(1,1)$$

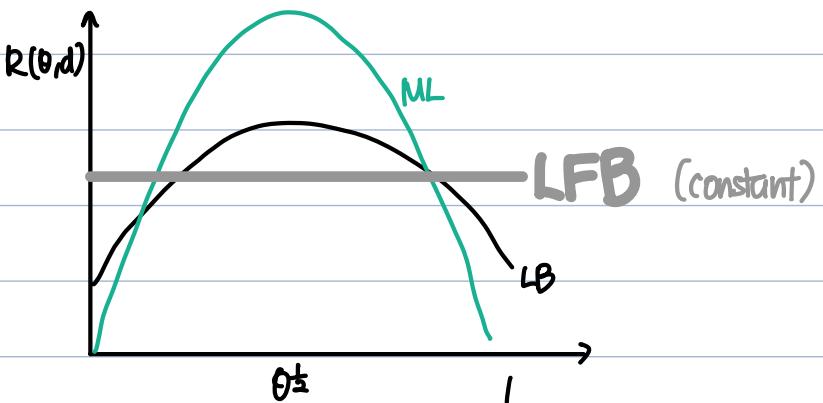
Laplace-Bayes

$\left. \begin{array}{l} \theta \sim \text{Beta}(\alpha_1, \alpha_2) \\ \text{as } \alpha_1, \alpha_2 \rightarrow 0 : \text{improper prior, non-informative prior} \end{array} \right)$

(Prior ②) Least Favorable Bayes :

$$\theta \sim \text{Beta}(\alpha_1, \alpha_2), \quad \alpha_1 = \alpha_2 = \sqrt{\frac{N}{2}}$$

$$\left\{ \begin{array}{l} d_{LFB}(\mathbf{z}) = \frac{s_N + \frac{1}{2}\sqrt{N}}{N + \sqrt{N}} \\ R(\theta, d_{LFB}) = \frac{N}{4(N + \sqrt{N})^2} \end{array} \right.$$



- Under a general $\theta \sim \text{Beta}(\alpha_1, \alpha_2)$, the posterior mean decision rule

$$d_{BB}(\mathbf{z}) = \frac{s_N + \alpha_1}{N + \alpha_1 + \alpha_2}$$

$$= \frac{N}{N + \alpha_1 + \alpha_2} \left(\frac{s_N}{N} \right) + \frac{\alpha_1 + \alpha_2}{N + \alpha_1 + \alpha_2} \left(\frac{\alpha_1}{\alpha_1 + \alpha_2} \right)$$

wavy underline MLE / sample mean wavy underline prior mean

: weighted avg
btw sample mean &
prior mean

(var of BB < var of MLE, but biased)

Q. Why is R highest for $\hat{\theta}_{ML}$ around $\theta = \frac{1}{2}$?

- risk = mean sq. error = bias² + variance
- MLE unbiased \rightarrow risk entirely in variance only
 \Rightarrow coin flip 50%. (50%. H, 50%. T): much more uncertainty in the middle \rightarrow huge variance
 (cf) if Heads = 95%. \rightarrow next flip almost surely H \rightarrow almost no uncertainty

Q. What is LFB est. and why is R_{LFB} constant?

- generally, R changes w/ θ .
- LFB is "designed" s.t. variance - bias² tradeoff remains constant.

10/70 Lecture

X. Recap

[Likelihood] $Z_i | \theta \sim \text{Ber}(\theta)$, $i=1, \dots, N$, $Z = (Z_1, \dots, Z_N)'$

[Prior Dist.] $\theta \sim \text{Beta}(\alpha_1, \alpha_2)$

(pulling a coin \rightarrow beta dist.)

\hookrightarrow start flipping the coin \rightarrow Bernoulli dist.

[Posterior Dist] $\theta | Z \sim \text{Beta}(S_N + \alpha_1, N - S_N + \alpha_2)$ ($S_N = \sum_{i=1}^N Z_i$)

\Rightarrow 2 ways of thinking about this:

{ (i) Subjective Bayesian (Savage) : Posterior as belief summary (θ being b/w 0.45 \sim 0.55)

(ii) Frequentist [don't like making probability statements about θ]

: θ is a fixed population parameter (θ is in an interval either w/ prob. 1 or 0)

(but might still use this procedure) want to use Bayesian Apparatus to construct procedures w/ good frequentist properties

* Point Estimation as a Bayesian

- [Goal] What is our best guess of θ as a Bayesian?

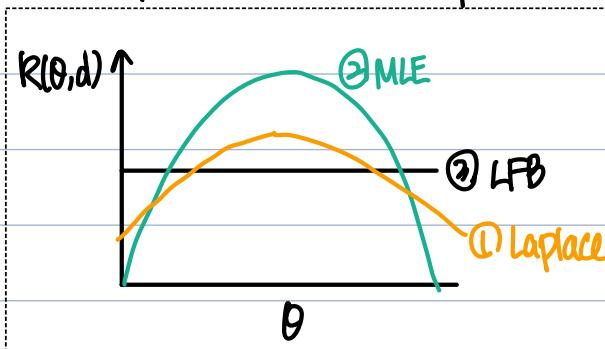
- [Idea] { (i) we need a loss function (this gives content to "best")

(ii) Recap: ① $\theta \sim \text{Beta}(1,1) = U[0,1]$ (Laplace)

- ② $\theta \sim \text{Beta}(0,0)$ (improper prior) \Rightarrow "proper posterior"

③ $\theta \sim \text{Beta}(\sqrt{\frac{N}{4}}, \sqrt{\frac{N}{4}})$ (Least Favorable Bayes) : minimax

\Rightarrow use posterior mean as point estimate



\Rightarrow risk-functions cross : admissibility

- an admissible estimator is not uniformly dominated by another estimator

- Choose estimators w/ good average risk properties

\hookrightarrow over possible values of θ ... but what distribution?

\Rightarrow use prior!

X. DIGRESSION

Complete Class Thm: Any admissible estimator is a Bayes rule
(or a limit of Bayes rule) wrt some prior

Minimax Procedure

(admissible +
constant risk
 \Rightarrow minimax)

$$\inf_{d \in \mathcal{D}} \sup_{\theta \in \Theta} R(\theta, d)$$

$$L(\theta, a) = (\theta - a)^2, \quad R(\theta, d) = \mathbb{E}_{\mathbf{z}|\theta} [L(\theta, d(\mathbf{z})) | \theta]$$

↑ target ↑ action

$$= \int L(\theta, d(\mathbf{z})) \cdot f(\mathbf{z}|\theta) dm(\mathbf{z})$$

multivariate integral
($\mathbf{z} = \text{vector}$)
for discrete & continuous

Average risk (Bayes risk) \leftarrow avg wrt prior

$$r(\pi, d) = \int R(\theta, d) \cdot \pi(\theta) dm(\theta)$$

depends on prior
avg over all possible states of the world

$$= \int \int L(\theta, d(\mathbf{z})) \cdot f(\mathbf{z}|\theta) dm(\mathbf{z}) \cdot \pi(\theta) dm(\theta)$$

$$= \int \left[\int L(\theta, d(\mathbf{z})) \cdot f(\mathbf{z}|\theta) \cdot \pi(\theta) dm(\theta) \right] dm(\mathbf{z})$$

① \curvearrowright switching order of integration

Consider we observe $\mathbf{Z} = \mathbf{z}$ and choose $d_m(\mathbf{z})$ to minimize ①

$$d\pi(\mathbf{z}) = \arg \min_{a \in \mathcal{A}} \int L(\theta, a) \underline{f(\mathbf{z}|\theta) \pi(\theta)} dm(\theta)$$

minimizing ①
= minimizing posterior expected loss
(\because posterior \propto likelihood · prior)

$$\pi(\theta|\mathbf{z}) \propto f(\mathbf{z}|\theta) \cdot \pi(\theta)$$

posterior likelihood · prior

Bayes Rule:

$$\begin{cases} \pi(\theta|\mathbf{z}) = \frac{f(\mathbf{z}|\theta) \cdot \pi(\theta)}{f(\mathbf{z})} \\ f(\mathbf{z}) = \int f(\mathbf{z}(t)) \cdot \pi(t) dt \end{cases}$$

$$= \arg \min_{a \in \mathcal{A}} \int L(\theta, a) \cdot \underline{\pi(\theta|\mathbf{z}) dm(\theta)}$$

: choosing $d\pi(\mathbf{z})$ to minimize ①

\Rightarrow minimize posterior expected loss

Solving: $d\pi(\mathbf{z}) = \arg \min \int L(\theta, a) \cdot \pi(\theta|\mathbf{z}) dm(\theta)$

$$= \mathbb{E}_{\theta|\mathbf{z}} [(θ - a)^2 | \mathbf{Z} = \mathbf{z}]$$

loss func
averaging over θ conditioned on \mathbf{z}

$$\begin{aligned}
 &= \mathbb{E}_{\theta|Z} \left[(\underbrace{\theta - \mathbb{E}[\theta|Z=z]}_{\text{(a)}} - \underbrace{(\alpha - \mathbb{E}[\theta|Z=z])^2}_{\text{(b)}}) \right] | Z=z \\
 &= (\alpha - \mathbb{E}[\theta|Z=z])^2 + V(\theta|Z=z) \quad \left. \begin{array}{l} \text{depends on our choice var. } \alpha \\ \text{posterior conditional variance} \end{array} \right\} \text{doesn't depend on } \alpha \\
 &\quad + 2(\alpha - \mathbb{E}[\theta|Z=z]) \underbrace{\mathbb{E}[\theta - \mathbb{E}[\theta|Z=z]]|Z=z}_0
 \end{aligned}$$

$$\Rightarrow d\pi(z) = \alpha^* = \mathbb{E}[\theta|Z=z]$$

Under squared error loss,

the posterior mean minimizes posterior expected loss.

Plan of Action: To minimize Bayes risk,

- | (i) we use the likelihood + our prior to compute the posterior,
- | (ii) then we use the posterior mean as our decision ("point estimate")

• Main Points

- { ① Bayes risk allows us to rank decision rules
- ② Ranking depends on our prior

\Rightarrow Alternative idea: look for minimax procedure (\leftarrow no prior needed here)

$$\inf_{d \in D} \sup_{\theta \in \Theta} R(\theta, d) \quad \text{: LFB}$$

* K-normal means (James-Stein estimator)

• Set-up

X, Y unknown $\left\{ \begin{array}{l} Y : \text{some outcome of interest (earnings)} \\ X : \text{a vector of covariates} \end{array} \right.$

Training sample : $\left\{ \begin{array}{l} Y_{(N \times 1)} = (Y_1, \dots, Y_N)' \\ X_{(N \times K)} = (X_1, \dots, X_N)'_{(K \times 1)} \end{array} \right.$

\Rightarrow use training sample to construct good predictions of new values of Y

- Assume that the covariate (feature) values for any future draw of Y will be in X

Condition on $X \rightarrow$ treat X as non-stochastic

- Notation

$$\underset{(N \times 1)}{m} = \mathbb{E}[Y] = \mathbb{E}[Y | X]$$

but depends on
specific values of X

↑
non-stochastic

$\hat{m}_i = \hat{m}(X_i)$: our prediction for Y when $X = X_i$

$$\underset{(N \times 1)}{\hat{m}} = \begin{pmatrix} \hat{m}(X_1) \\ \vdots \\ \hat{m}(X_N) \end{pmatrix} \rightarrow \text{in general, } \hat{m} \text{ will depend on } \underset{(N \times 1)}{Y}$$

$\left(\begin{array}{l} \hat{m} = d(Y) \\ = d(Y, X) \end{array} \right)$:: non-stochastic

- Squared Error Loss $(\|m\| = \sum_{i=1}^N m_i^2)^{\frac{1}{2}}$

- Risk : $\mathbb{E}[\|\hat{m} - m\|^2]$ \Rightarrow want to choose \hat{m} to be close to $m = \mathbb{E}[Y] = \mathbb{E}[Y | X]$

$$= \mathbb{E} \left[\underbrace{\sum_{i=1}^N (\hat{m}(X_i) - m(X_i))^2}_{\text{loss}} \right] \rightarrow \text{sq. error loss}$$

but for N different targets

$$= \begin{pmatrix} m(X_1) \\ \vdots \\ m(X_N) \end{pmatrix} = \begin{pmatrix} \mathbb{E}[Y_1] \\ \vdots \\ \mathbb{E}[Y_N] \end{pmatrix}$$

$$= \sum_{i=1}^N \mathbb{E}[(\hat{m}(X_i) - m(X_i))^2]$$

* 11/4 Lecture

* Set-up to the K-Normal mean

Training sample : (X_i, Y_i) $i=1, \dots, N$, (X_i, Y_i) iid $F_{X,Y}$

treat $\mathbf{X} = (X_1, \dots, X_N)'$ as non-stochastic ("fixed design matrix")

$m_i = E[Y_i]$ \because note $E[Y_i] \neq E[Y_{1n}]$ in general since $X_i \neq X_{1n}$ in general

sometimes in an abuse of notation; $m_i = E[Y | X = X_i] = E[Y_i]$

\mathbf{Y} vector of outcome, $\hat{\mathbf{m}} = E[\mathbf{Y}]$

• Prediction Goal: observe $X = X_{19}$, predict Y

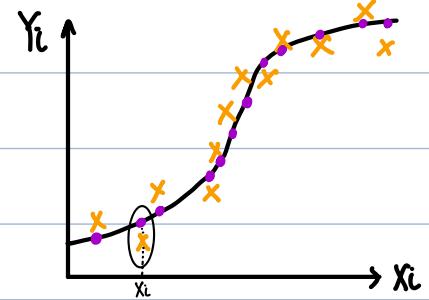
• Loss function:

$$\|\hat{\mathbf{m}} - \mathbf{m}\|^2 = \sum_{i=1}^N (\hat{m}_i - m_i)^2$$

* Side-note:
 $\begin{cases} Y_i \stackrel{!}{=} Y_j & \therefore m_i = E[Y_i] = E[Y_j] \\ Y_i \neq Y_j & \text{for } i \neq j \therefore X_i \neq X_j \end{cases}$

$E[\|\hat{\mathbf{m}} - \mathbf{m}\|^2]$: Risk | MSE
 apparent error

\rightarrow we want something that works well on avg!

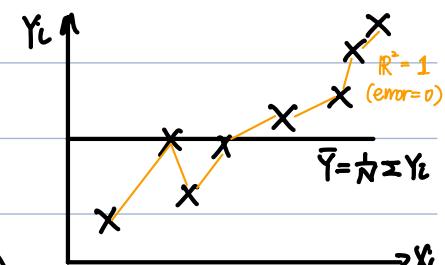


[OLS \neq BLUE under joint random sampling \rightarrow EXAM!!]

$$\begin{aligned} E[\|\mathbf{Y} - \hat{\mathbf{m}}\|^2] &= E[\|(\mathbf{Y} - \mathbf{m}) + (\mathbf{m} - \hat{\mathbf{m}})\|^2] \\ &= E[\|\mathbf{Y} - \mathbf{m}\|^2] + E[\|\hat{\mathbf{m}} - \mathbf{m}\|^2] \quad \text{① Risk func} \\ &\quad + 2E[(\mathbf{Y} - \mathbf{m})'(\mathbf{m} - \hat{\mathbf{m}})] \quad \text{② cross-product func} \end{aligned}$$

(since $Y_i = m(X_i) + \sigma u_i$, $u_i \sim N(0, 1)$,
 some mean + noise $u_i | X \sim N(0, 1)$)

$\rightarrow Y_i \sim N(m(X_i), \sigma^2)$: classical regression model



$$\textcircled{1} E[\|\mathbf{Y} - \mathbf{m}\|^2] = E\left[\sum_{i=1}^N (Y_i - m_i)^2\right] = \sum_{i=1}^N E[(Y_i - m_i)^2] = N \sigma^2$$

$$\textcircled{3} E[(\mathbf{Y} - \mathbf{m})'(\mathbf{m} - \hat{\mathbf{m}})] = E[(\mathbf{Y} - \mathbf{m})' \mathbf{m}] - E[(\mathbf{Y} - \mathbf{m})' \hat{\mathbf{m}}]$$

$$= E[\underbrace{(\mathbf{Y} - \mathbf{m})' \mathbf{m}}_{(NX1)}] - E[(\mathbf{Y} - \mathbf{m})' \hat{\mathbf{m}}] \quad \textcircled{4}$$

$(E[\mathbf{Y}] - \mathbf{m})' \mathbf{m} = 0$

$$\textcircled{4} E[(\mathbf{Y} - \mathbf{m})' \hat{\mathbf{m}}] = E\left[\sum (Y_i - m_i) \hat{m}_i\right] = \sum_{i=1}^N E[(Y_i - m_i) \hat{m}_i] \quad \textcircled{5}$$

$$\left(= \begin{pmatrix} Y_1 - m_1 \\ \vdots \\ Y_N - m_N \end{pmatrix}' \begin{pmatrix} \hat{m}_1 \\ \vdots \\ \hat{m}_N \end{pmatrix} \right)$$

$$= \sum_{i=1}^N \text{Cov}(Y_i, \hat{m}_i)$$

$$\left(\begin{aligned} \textcircled{5} \quad \mathbb{E}[(Y_i - m_i) \hat{m}_i] &= \mathbb{E}[Y_i \hat{m}_i] - \cancel{m_i} \mathbb{E}[\hat{m}_i] \\ &= \mathbb{E}[Y_i \hat{m}_i] - \mathbb{E}[Y_i] \mathbb{E}[\hat{m}_i] \\ &= \text{Cov}(Y_i, \hat{m}_i) \end{aligned} \right)$$

$$\Rightarrow \mathbb{E}[\|Y - \hat{m}\|^2] = \underset{\text{① noise floor}}{N\sigma^2} + \underset{\text{② risk / MSE}}{\mathbb{E}[\|\hat{m} - m\|^2]} - \underset{\text{③ degrees-of-freedom (model complexity)}}{2\sigma^2 \cdot df(\hat{m})}$$

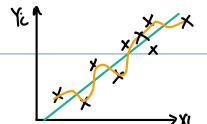
apparent error
(aka. expected in sample fit)

① noise floor ② risk / MSE ③ degrees-of-freedom (model complexity)

$$\Rightarrow \mathbb{E}[\|\hat{m} - m\|^2] = -N\sigma^2 + \mathbb{E}[\|Y - \hat{m}\|^2] + 2\sigma^2 \cdot df(\hat{m})$$

Trade-off btw "in sample fit" and "model complexity"

* Connect-the-dots est: $\mathbb{E}[\|Y - \hat{m}\|^2] = 0$ but complexity ↑↑



* Sample mean: complexity ↓↓ but $\mathbb{E}[\|Y - \hat{m}\|^2]$ huge

* K-normal means

• Set-up

$$Y_i = m(X_i) + \sigma u_i, \quad u_i | X \sim N(0,1) \quad \text{or} \quad u_i \sim N(0,1)$$

$i = 1, \dots, N, \quad X_i \in \mathbb{X}, \quad x \in \mathbb{X}$

assumption on m : $m(x) \approx \sum_{k=1}^K \alpha_k g_k(x)$ w/ $g_k(x)$ known basis functions,
 α_k unknown coefficient

smoothing assumption

$$N=1000, \quad K=100 \quad // \quad (\text{ex}) \quad g_k(x) = x^{k-1} \quad \text{where} \quad k=1, \dots, K : \text{series basis}$$

• Gram-Schmidt Orthogonalization

Euclidean space

$$\langle f, g \rangle = \frac{1}{N} \sum_{i=1}^N f(x_i) g(x_i) \quad \text{where} \quad f = \begin{pmatrix} f(x_1) \\ \vdots \\ f(x_N) \end{pmatrix}, \quad g = \begin{pmatrix} g_1(x_1) \\ \vdots \\ g_K(x_N) \end{pmatrix}$$

$$\|f\| = \sqrt{\langle f, f \rangle} = \sqrt{\frac{1}{N} \sum_{i=1}^N f(x_i)^2}$$

$$\text{proj}_f(g) = \pi(g|f) = \frac{\langle f, g \rangle}{\langle f, f \rangle} \cdot f \quad \leftarrow \text{OLS}$$

$$u = g - \text{proj}_f(g) \quad \text{: projection error}$$

$$[\text{PF}] \langle u, f \rangle = \langle g - \text{proj}_f(g), f \rangle$$

$$= \langle g - \frac{\langle f, g \rangle}{\langle f, f \rangle} f, f \rangle$$

$$= \langle g, f \rangle - \frac{\langle f, g \rangle \langle f, f \rangle}{\langle f, f \rangle} \quad (\text{bilinearity})$$

$$= 0$$

Recursive Basis function

$$f_1(x) = g_1(x) : \text{generally a constant}$$

$$f_2(x) = g_2(x) - \frac{\langle g_2, f_1 \rangle}{\langle f_1, f_1 \rangle} f_1(x)$$

$$\left(\begin{array}{l} \text{assume } g_k(x) = x^{k-1}; 1, x, x^2, \dots \Rightarrow f_1(x) = 1 \\ f_2(x) = x - \bar{x} \quad (\bar{x} = \frac{1}{N} \sum x_i) \end{array} \right)$$

$$f_3(x) = g_3(x) - \frac{\langle g_3, f_2 \rangle}{\langle f_2, f_2 \rangle} f_2(x) - \frac{\langle g_3, f_1 \rangle}{\langle f_1, f_1 \rangle} f_1(x)$$

:

:

$$f_k(x) = g_k(x) - \sum_{j=1}^{k-1} \frac{\langle g_k, f_j \rangle}{\langle f_j, f_j \rangle} f_j(x)$$

$$\langle f_1, f_2 \rangle = 0, \langle f_1, f_3 \rangle = 0, \dots \langle f_k, f_j \rangle = 0 \quad k \neq j, k, j < K$$

$$\Rightarrow \underbrace{g_1, \dots, g_k}_{m(x)} \Rightarrow \underbrace{f_1, \dots, f_k}_{\text{orthogonal to each other}}$$

$$\Rightarrow \text{orthogonality} : \phi_1(x) = \frac{f_1(x)}{\|f_1(x)\|}, \dots \phi_k(x) = \frac{f_k(x)}{\|f_k(x)\|}$$

$$\|\phi_k\| = 1$$

• Maximum Likelihood Estimator (MLE)

$$m(x) = \sum_{k=1}^K \theta_k \phi_k(x)$$

$$(Y_i = m(x_i) + \sigma u_i, u_i \sim N(0, 1), \sigma^2 \text{ known})$$

$$W(X) = \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_k(x) \end{pmatrix}, \theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_K \end{pmatrix}$$

log likelihood func:

$$l(Y | X; \theta) = -\frac{N}{2} \ln 2\pi - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - \underbrace{w(X_i)' \theta}_{w_i})^2$$

$$\hat{\theta}_{ML} = \left[\frac{1}{N} \Sigma w_i w_i' \right]^{-1} \left[\frac{1}{N} \Sigma w_i Y_i \right]$$

$$\hat{m}(x) = w(x)' \hat{\theta}_{ML}$$

$$\boxed{\frac{1}{N} \sum_{i=1}^N w_i w_i' = I_k}$$

>Show this: $\frac{1}{N} \sum_{i=1}^N \begin{bmatrix} \phi_1(x_i) \\ \phi_2(x_i) \\ \vdots \\ \phi_N(x_i) \end{bmatrix} \begin{bmatrix} \phi_1(x_i) \\ \phi_1(x_i)\phi_2(x_i) \\ \vdots \\ \phi_N(x_i)\phi_1(x_i) \end{bmatrix}' = I_k$

$$\hat{\theta}_{ML} = \bar{Z} = \frac{1}{N} \Sigma w_i Y_i$$

\Rightarrow note that: $\bar{Z} \sim N(\theta, \frac{\sigma^2}{N} I_k)$

k-normal mean

* 11/6 Lecture

- Set-up

$$Y_i = m(x_i) + \sigma u_i, \quad u_i | X \sim N(0, 1) \quad i=1, \dots, N, \quad X \text{ non-stochastic, fixed}$$

σ^2 known

Assume $m(x) = \sum_{k=1}^K \alpha_k \phi_k(x)$

Gram-Schmidt orthonormalization : $m(x) = \sum_{k=1}^K \theta_k \phi_k(x)$ orthonormal basis

$$\hat{\theta}_{ML} = \left[\frac{1}{N} \sum_{i=1}^N w_i w_i' \right]^{-1} \cdot \left[\frac{1}{N} \sum_{i=1}^N w_i Y_i \right] \quad \text{where } \frac{1}{N} \sum_{i=1}^N w_i w_i' = I_N, \quad w_i = w(x_i)$$

$$:= Z = \frac{1}{N} \sum_{i=1}^N w_i Y_i \quad \text{Gaussian}$$

$Z \sim N(\theta, \frac{\sigma^2}{N} I_K)$

Gaussian

$$w(x) = \begin{pmatrix} \phi_1(x) \\ \vdots \\ \phi_K(x) \end{pmatrix}, \quad \theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_K \end{pmatrix}$$

• Unbiasedness of $\hat{\theta}_{ML}$ (remember X is non-stochastic) if X is stochastic, mean of θ doesn't always exist

$$\begin{aligned} - E[Z | X] &= \frac{1}{N} \sum_{i=1}^N w_i E[Y_i | X] \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N w_i w_i' \theta}_{I_N} + \frac{\sigma}{N} \sum_{i=1}^N w_i E[u_i | X] \end{aligned}$$

(why? $Y_i = w_i' \theta + \sigma u_i$)

= θ : OLS is unbiased when X is non-stochastic

Covariance matrix

$$\begin{aligned} - V[Z_k | X] &= V\left[\frac{1}{N} \sum_{i=1}^N \phi_k(x_i) Y_i | X\right] \\ &\quad \text{diagonal elements} \quad \text{WKL} \\ &= \frac{1}{N^2} \sum_{i=1}^N V(\phi_k(x_i) Y_i | X) \\ &= \frac{\sigma^2}{N} \frac{1}{N} \sum_{i=1}^N \phi_k(x_i)^2 \\ &\quad \underbrace{\|\phi_k\|^2}_{} = 1 \\ &= \frac{\sigma^2}{N} \quad \text{diagonal elements } (k=1, \dots, K) \end{aligned}$$

$\because Y_i$'s independent : $V(\Sigma) = \Sigma V$

$$V(aX+bY) = a^2 V(X) + b^2 V(Y) + 2ab C(X, Y)$$

$$\begin{aligned} - C(z_j, z_k | X) &= C\left[\frac{1}{N} \sum_{i=1}^N \phi_j(x_i) Y_i, \frac{1}{N} \sum_{i=1}^N \phi_k(x_i) Y_i | X\right] \\ &\quad \text{off diagonal elements} = 0 \\ &= \frac{\sigma^2}{N} \cdot \underbrace{\frac{1}{N} \sum_{i=1}^N \phi_j(x_i) \cdot \phi_k(x_i)}_{}, \quad j \neq k \end{aligned}$$

$\langle \phi_j, \phi_k \rangle = 0$ j -th k -th elements of the orthogonal basis : orthogonal by construction
 \rightarrow inner product = 0

$\Rightarrow \hat{\theta}_{ML} = Z \sim N(\theta, \frac{\sigma^2}{N} I_K)$

$\boxed{x \times x \times x \times x} \hat{m}(x) = w(x)' \cdot \hat{\theta}_{ML}$

• Getting back to "Risk" (motivation: forecasting problem)

$$\underbrace{\|\hat{m} - m\|^2}_{\text{squared error loss}} = \frac{1}{N} \sum_{i=1}^N \underbrace{(w_i' \theta_{ML})}_{m(x_i)} - w_i' \theta$$

$$\hat{\theta}_{ML} = \mathbf{z}$$

(to emphasize $\hat{\theta}_{ML}$ is r.v.)

$$= \frac{1}{N} \sum_{i=1}^N w_i' (\mathbf{z} - \theta)$$

$$= \frac{1}{N} \sum_{i=1}^N \left(\sum_{k=1}^K \phi_k(x_i) \cdot (z_k - \theta_k) \right)^2$$

orthogonal to one another

$$= \sum_{k=1}^K (z_k - \theta_k)^2$$

$$w_i = \begin{pmatrix} \phi_1(x_i) \\ \vdots \\ \phi_K(x_i) \end{pmatrix}$$

Loss function: $L(\mathbf{z}, \theta) = \|\mathbf{z} - \theta\|^2 = \sum_{k=1}^K (z_k - \theta_k)^2$

Risk function: $R(\mathbf{z}, \theta) = \mathbb{E}[L(\mathbf{z}, \theta)]$

$$\mathbb{E}\left[\sum_{k=1}^K (z_k - \theta_k)^2\right] = \sum_{k=1}^K \mathbb{E}[(z_k - \theta_k)^2] = \frac{K\sigma^2}{N}$$

\downarrow multivariate normal, identity cov. matrix

The Risk of MLE varies with:

- (i) model complexity, $K \uparrow$ not a choice!
- (ii) amount of data, $N \downarrow$
- (iii) noise floor, $\sigma^2 \uparrow$

* James-Stein Type Estimators

• Loss & Risk function

Loss: $\mathcal{L} = \{G\mathbf{z} : G = \text{diag}\{c_1, \dots, c_K\}, c_k \in [0, 1], k=1, \dots, K\}$

Risk: $\mathbb{E}\left[\sum_{k=1}^K (c_k z_k - \theta_k)^2\right] = \mathbb{E}\left[\sum_{k=1}^K (c_k(z_k - \theta_k) - (1-c_k)\theta_k)^2\right]$
 $= \underbrace{\frac{\sigma^2}{N} \sum_{k=1}^K c_k^2}_{V(C, K)} + \underbrace{\sum_{k=1}^K (1-c_k)^2 \theta_k^2}_{\text{Bias}^2}$

• "Oracle"

(Thought experiment: think of a new θ est.
 \rightarrow which C would minimize risk?)

$$c_k^* = \frac{\theta_k}{\frac{1}{N} + \theta_k^2}, \quad k=1, \dots, K \quad \Rightarrow \quad C^* \cdot \mathbf{z} = \hat{\theta}_0$$

if θ big \rightarrow deviation = lots of bias
small \rightarrow not much bias, more variance \rightarrow scale a lot
 \rightarrow don't scale a lot

\hookrightarrow imagine it's like a constant

(plug c_k^* back into the risk func)

$$\Rightarrow \inf_{\theta \in \Omega} R(\hat{\theta}_0, \theta) = \frac{\sigma^2}{N} \sum_{k=1}^K \left(\frac{\theta_k^2}{\frac{\sigma^2}{N} + \theta_k^2} \right) < \frac{K\sigma^2}{N} \quad (= R_{ML}) \quad : \text{"bound"}$$

• Stein - Unbiased - Risk - Estimate (SURE)

\checkmark model selection
 \checkmark cross validation
 \checkmark cross fitting .

$$Z \sim N(\theta, \sigma^2 I_K), \quad \hat{\theta} = \hat{\theta}(Z), \quad g(Z) = \hat{\theta}(Z) - Z$$

$$\hat{R}_{\text{SURE}}(Z_i) = -K\sigma^2 + 2\sigma^2 \cdot \sum_{k=1}^K \frac{\partial g_k(Z)}{\partial z_k} + \sum_{k=1}^K (\hat{\theta}_k - z_k)^2$$

: unbiased risk estimator

$$E[\hat{R}_{\text{SURE}}(Z_i)] = E[\|\hat{\theta} - \theta\|^2] \quad : \text{UNBIASED}$$

• James-Stein Type Estimators

$$\hat{\theta}_{JS}(Z) = \left(1 - \frac{(K-2)}{Z'Z} \frac{\sigma^2}{N} \right) Z$$

$$g_{JS}(Z) = -\left(\frac{(K-2)\sigma^2/N}{Z'Z} \right) Z \quad \Rightarrow \quad \sum_{k=1}^K \frac{\partial g_k(Z)}{\partial z_k} = -\frac{(K-2)^2 \cdot \sigma^2}{Z'Z}$$

$$\hat{R}_{\text{SURE}}(Z_i) = \underbrace{\frac{K}{N}\sigma^2}_{\hat{\theta}_{JS}} - \frac{(K-2)^2}{Z'Z} \cdot \frac{\sigma^4}{N^2}$$

$$R(\hat{\theta}_{JS}, \theta) = E[\hat{R}_{\text{SURE}} | Z]$$

$$= \underbrace{\frac{K}{N}\sigma^2}_{R(\hat{\theta}_{ML}, \theta)} - \underbrace{(K-2)^2 \cdot \frac{\sigma^4}{N^2} \cdot E\left[\frac{1}{Z'Z}\right]}_{\textcircled{1}} \quad (\text{for } K \geq 3) \quad < R(\hat{\theta}_{ML}, \theta)$$

$$Z_k^2 = \frac{\sigma^2}{N} \left(\frac{\theta_k}{\sigma/\sqrt{N}} + u \right)^2, \quad u \sim N(0, 1), \quad Z_k \sim N(\theta_k, \frac{\sigma^2}{N})$$

$$Z'Z \sim \frac{\sigma^2}{N} V \quad \text{where } V \text{ is a non-central } \chi^2 \text{ r.v. w.r.t. } K \text{ degrees of freedom}$$

and non-centrality parameter $\rho = \sum_{k=1}^K N(\theta_k^2 / \sigma^2)$

$$V \sim \chi^2_{K+2W} \quad \text{where } W \sim \text{Poisson}(\frac{\rho}{2}), \quad E[W] = \frac{\rho}{2} \quad : \text{fact}$$

$$E\left[\frac{1}{Z'Z}\right] = \frac{N}{\sigma^2} E\left[\frac{1}{V}\right]$$

$$(L.I.E.) = \frac{N}{\sigma^2} E\left[E\left[\frac{1}{V} | W\right]\right]$$

$\because K \geq 2, \quad S \sim \chi^2_k, \quad E\left[\frac{1}{S}\right] = \frac{1}{k-2}$

conditioned on W , becomes central χ^2 r.v.

$$= \frac{N}{\sigma^2} E\left[\frac{1}{K-2+2W}\right]$$

$$(Jensen's \text{ Inequality}) \geq \frac{N}{\sigma^2} \frac{1}{K-2+N \sum_{k=1}^K (\theta_k^2 / \sigma^2)} \quad \left(= \frac{N}{\sigma^2} \cdot \frac{1}{K-2+2E[W]} \right)$$

$$= \frac{1}{(K-2) \cdot \frac{\sigma^2}{N} + \|\theta\|^2}$$

$$R(\hat{\theta}_{JS}, \theta) \leq \frac{k}{N} \sigma^2 - \frac{(k-2)^2 \frac{\sigma^4}{N^2}}{(k-2) \frac{\sigma^2}{N} + \|\theta\|^2} \leq R(\hat{\theta}_{ML}, \theta)$$

$$= \frac{2}{N} \sigma^2 + \frac{(k-2) \frac{\sigma^2}{N} \|\hat{\theta}\|^2}{(k-2) \frac{\sigma^2}{N} + \|\theta\|^2}$$

\Rightarrow When $k \geq 3$, $R(\hat{\theta}_{JS}, \theta) \leq R(\hat{\theta}_{ML}, \theta)$ for $\forall \theta \in \Theta$

\therefore MLE is inadmissible? ($\exists \hat{\theta}_{JS}$ that dominates $\hat{\theta}_{ML}$)

X. Ensemble Risk

vs

Component Risk

$$\sum_{k=1}^K (\hat{\theta}_{jk} - \theta_{jk})^2$$

$$(\hat{\theta}_{jk} - \theta_{jk})^2$$

\rightarrow shrink

The James–Stein estimator made this point dramatically in 1961, and made it in the context of just a few unknown parameters, not hundreds or thousands. It begins the story of *shrinkage estimation*, in which deliberate biases are introduced to improve overall performance, at a possible danger to individual estimates. Chapters XXX and XXX will carry on the story in its modern implementations.

* 11/17 Lecture

• Finishing up K-Normal means (& JS est)

• Set-up

$$\mathbf{Z} \sim N(\boldsymbol{\theta}, \frac{\sigma^2}{N} \mathbf{I}_K), \quad Y = m(X) + \sigma u, \quad u \sim N(0, \sigma^2)$$

$$m(X) = \sum_{k=1}^K \varphi_k(x) \theta_k \quad (K=200, N=1000)$$

large enough s.t.
we don't have to worry about
whether K is suff. large

$$\|\hat{m} - m\|^2 = \sum_{k=1}^K (\hat{\theta}_k - \theta_k)^2 : \text{sq. error loss}$$

$$\hat{\theta}_{ML} = \mathbf{Z} = \underbrace{W'Y}_{(N \times K)(N \times 1)} \quad \text{where } (W'W) = \mathbf{I}_K$$

→ MLE is inadmissible

→ showed this by demonstrating that the risk function of $\hat{\theta}_{JS}$
lies below that of $\hat{\theta}_{ML}$ for $\forall \theta \in \Theta$ ($K \geq 3$)

• JS Estimator

$$\hat{\theta}_{JS}(\mathbf{Z}) = \left(1 - \frac{(K-2)}{\mathbf{Z}'\mathbf{Z}} \frac{\sigma^2}{N} \right) \mathbf{Z}$$

"shrinkage estimator"

$$\mathcal{L} = \{ \mathbf{CZ} : \mathbf{C} = \text{diag}\{c_1, \dots, c_K\}, \quad c_k \in [0, 1], \quad k=1, \dots, K \}$$

multiply each coefficient by 0~1 $\begin{cases} x_1: \text{MLE (unbiased)} \\ x_0: \text{zero} \end{cases}$ ↗ bias ↘ retain all variance

Recall: we consider "oracle" choice for C.

Use SURE to think about how to choose C in practice.

$$\hat{R}_{\text{SURE}}(\mathbf{Z}, \mathbf{C}) = \frac{K}{N} \sigma^2 - \frac{2}{N} \sigma^2 \sum_{k=1}^K (1-c_k) + \sum_{k=1}^K (\hat{\theta}_k - z_{ik})^2 \quad (\hat{\theta} = \mathbf{CZ})$$

$$(R_{\text{SURE}}(\mathbf{Z}) = K\sigma^2 + 2\sigma^2 \sum_{k=1}^K \frac{\partial \hat{\theta}_k(\mathbf{Z})}{\partial z_{ik}} + \sum_{k=1}^K (\theta_k - z_{ik})^2)$$

$E[\hat{R}_{\text{SURE}}(\mathbf{Z}, \mathbf{C})]$: verify our earlier calculation.

Choose $\hat{\mathbf{C}}$ to minimize an estimate of risk

$$(*) = \underbrace{\frac{\sigma^2}{N} \sum_{k=1}^K c_k^2}_{\text{variance reduction}} + \underbrace{\sum_{k=1}^K (z_k^2 - \frac{\sigma^2}{N}) (1-c_k)^2}_{\text{estimate of bias sq.}} \quad (\text{heuristic})$$

$$\left(E[z_k^2] = \underbrace{V(z_k)}_{\sigma^2/N} + \underbrace{E[z_k]^2}_{\theta_k^2} \right)$$

$$\left(\theta_{k,0} = c_k z_k \right)$$

$$V(\hat{\theta}_{k,0}) = c_k^2 V(z_k) = c_k^2 \cdot \frac{\sigma^2}{N}$$

$$(\text{FOC}) \quad \hat{c}_k = 1 - \frac{N^{-1} \sigma^2}{z_k^2} \quad k=1, \dots, K ; \quad z_k^2 = \frac{\sigma^2}{N} + \theta_k^2 + \text{noise}$$

⇒ USE: Universal Series Estimator

└ How to choose optimal $\hat{\epsilon}_k$?

─ How does it connect w/ oracle?

⇒ oracle shrinks a lot

when θ_{ik} is small relative to "noise", σ^2

• x. oracle:

$$\hat{\epsilon}_k^* = \frac{\theta_{ik}^2}{\sigma_k^2 + \theta_{ik}^2}$$

if we knew the truth, this is the optimal shrinkage coefficient (let $c \approx 0$ if true $\theta = 0$)

─ Oracle: choose $\hat{\epsilon}_{ik0}$ to minimize Risk

(evaluating risk requires omniscience)

USE: choose $\hat{\epsilon}_k$ to minimize SURE b/c we don't know the actual risk!

(= Empirical Risk Minimization)

• Summary

Setting: $\{X_i, Y_i\}_{i=1}^N$, $Y_i = m(X_i) + \sigma u_i$,

$$m(x) = \sum_{k=1}^K p_k(x_i) \hat{\theta}_k$$

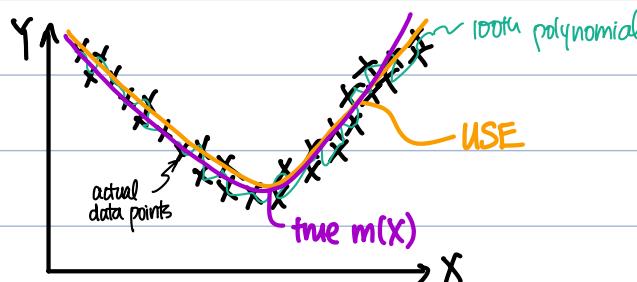
Steps: ① use OLS to compute $\hat{\theta}_{MLE} = \bar{z}$ → inadmissible (that's why we start looking into shrinking the coefficients)

$$\hat{\epsilon}_k = 1 - \frac{\sigma^2/N}{\hat{\theta}_k^2} \quad (k=1, \dots, K)$$

$$\hat{m}^{USE}(x_i) = \sum_{k=1}^K p_k(x_i) \hat{\theta}_k^{USE}$$

$$\text{where } \hat{\theta}_k^{USE} = \hat{\epsilon}_k \cdot \hat{\theta}_{MLE} = \hat{\epsilon}_k \cdot \bar{z}$$

$$\hat{m}^{OLS}(x_i) = \sum_{k=1}^K p_k(x_i) \hat{\theta}_k^{OLS}$$



* Bayesian Bootstrapping

• Intro

measuring uncertainty

want to start thinking about uncertainty w/ "sample in hand"

{ Bayesian

Frequentist (large sample)

why Bayesian bootstrapping: (i) a way to quantify uncertainty from Bayesian perspective

{ (ii) introduction to simulation (w/ numerical integration)

(iii) Bayesian b.s. is also a weighted frequentist b.s.

• Set-up

$$Z = (X', Y')' \in \bar{Z} = \{z_1, \dots, z_J\} \quad (J \text{ is massive})$$

(Kx1) x 1
(Kx1)
(1x1)

support
is a set

(ex) Y : earnings, $Y_i \in \{0, 1, \dots, 88,000\}$

$$\left. \begin{array}{l} X_1: \text{gender} \\ X_2: \text{yrs of schooling} \\ X_3: \text{AFQT score} \end{array} \right\} \quad \left. \begin{array}{l} X_{1i} \in \{0, 1\} \\ X_{2i} \in \{0, \dots, 20\} \\ X_{3i} \in \{1, \dots, 100\} \end{array} \right\}$$

$$J = 2 \times 21 \times 100 \times 88,000 = 369,600,000$$

$$z \sim F_z$$

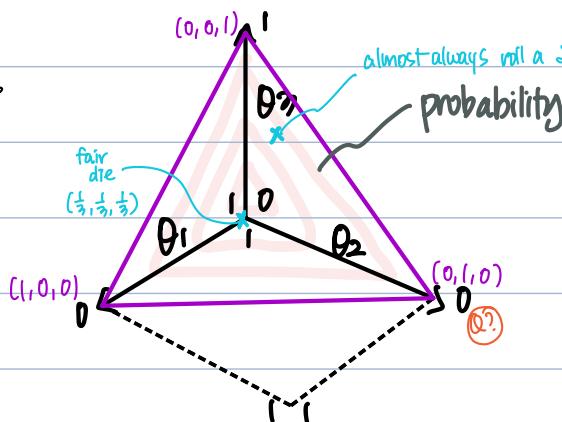
$$\Pr(Z = z_j | \underline{\theta}) = \theta_j \quad \text{where } \underline{\theta} = (\theta_1, \dots, \theta_J)'$$

$$\underline{\theta} \in \Theta = \mathbb{S}^{J-1} = \left\{ \underline{\theta} = (\theta_1, \dots, \theta_J)' \in \mathbb{R}^J, \theta_j \geq 0, \sum_{j=1}^J \theta_j = 1 \right\}$$

probability (unit)
simplex

all non-negative
sum to 1

imagine a 6-sided die
but only w/ 3 numbers
(1, 1, 2, 2, 3, 3)



probability mass dist.

\Rightarrow Each point of the simplex
is a possible population.
(state of the world)

• Likelihood (multinomial) (\because Bayesian!)

$$\text{Set-up: } \underline{\theta} \in \Theta = \mathbb{S}^{J-1}, \Pr(Z = z_j | \underline{\theta}) = \theta_j \quad (\underline{\theta}: \text{"population"})$$

$$\bar{Z} = (z_1, \dots, z_N) = z = (1, 1, 1, 3, 2, 1, 3, 2, \dots)$$

$$\text{Likelihood: } f(z | \underline{\theta}) = \frac{N!}{N_1! \times \dots \times N_J!} \cdot \prod_{j=1}^J \theta_j^{N_j}$$

capturing different ordering of draws

$$\text{Prior (Dirichlet)} : \pi(\theta_1, \dots, \theta_J) = \frac{\Gamma(\sum_{j=1}^J \alpha_j)}{\prod_{j=1}^J \Gamma(\alpha_j)} \cdot \prod_{j=1}^J \theta_j^{\alpha_j - 1} \quad [\text{gamma func}]$$

describes a "population of populations"

(from a bag of dice \rightarrow pick a die ("population") \rightarrow roll a die ("sample"))

assigns probability mass to simplex

