

Mathematical Statistics II

STA2212H S LEC9101

Week 9

March 17 2021

Start recording!



Christopher Mims
@mims



This is one of the coolest and most creative visualizations of a phenomenon I've ever seen

It's also an important account of one of the biggest potential climate tipping points we may be rapidly pushing ourselves over

nytimes.com/interactive/2021/climate/tipping-points/

2021-03-04, 6:18 AM

link

- overview of Bayesian inference
- posterior predictive distribution
- Bayesian computation
- Example 11.9 AoS
- Bayesian hierarchical models

$\times 10.8$



$$p(x_{n+1} | x^n)$$

$$= \underbrace{\int f(x_{n+1} | \theta) \pi(\theta | x^n) d\theta}$$

\uparrow ass'g $x_{n+1} \perp\!\!\!\perp x_n$

$$p(x_{n+1}) = \int f(x_{n+1} | \theta) \pi(\theta) d\theta$$

\uparrow sometimes used for planning slides

- overview of Bayesian inference
- posterior predictive distribution
- Bayesian computation
- Example 11.9 AoS
- Bayesian hierarchical models

Addendum: If X_1, \dots, X_n are i.i.d. $N(\mu, \sigma^2)$, both parameters unknown, then the conjugate priors for μ and σ^2 are

$$\mu \sim N(\mu_0, \tau^2); \quad \sigma^2 \sim IG\left(\frac{\nu_0}{2}, \frac{\nu_0 \sigma_0^2}{2}\right)$$

Inverse Gamma distribution $IG(\alpha, \beta)$ has density $f(t; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \frac{1}{x^{\alpha+1}} e^{-\beta/x}$

$$\mathcal{G} \quad x^{\alpha-1} e^{-\beta x}$$

1. Notes from March 12 (Friday) still to be posted
2. Friday March 19 – Stein's paradox B-H proof
3. empirical Bayes
4. introduction to decision theory

- Mar 22 3.00 – 4.00 pm EDT
[Data Science ARES](#)
Jesse Cisewski Kehe
University of Wisconsin-Madison
“Astrostatistics: From Exoplanets to
the Large-scale Structure of the Universe”



Bayesian hierarchical model, recap

SM Ex. 11.25

$$\begin{aligned}
 & \bullet x_i | \theta_i \sim N(\theta_i, v_i), \text{ independent}; \quad \theta_i | \mu \sim N(\mu, \sigma^2), i = 1, \dots, n \text{ i.i.d.} \\
 & \hat{\theta}_i := x_i \quad \text{hyper-par} \quad \mu \sim N(\mu_0, \tau^2) \\
 & E(\theta_i | \underline{x}) = x_i \left(\frac{\sigma^2}{\sigma^2 + v_i} \right) + \mu \left(\frac{1}{\sigma^2 + v_i} \right) \quad v_i, \sigma^2, \mu_0, \tau^2 \\
 & \text{est. } \hat{\theta}_i = x_i \left(\frac{\sigma^2}{\sigma^2 + v_i} \right) + \tilde{\mu} \left(\frac{1}{\sigma^2 + v_i} \right) \quad \text{all known} \\
 & \cancel{E(\mu | \underline{x})} = \frac{\sum x_i / (\sigma^2 + v_i)}{\sum 1 / (\sigma^2 + v_i)} = \tilde{\mu} \quad \text{shrinkage estimators}
 \end{aligned}$$

If σ^2 unknown, no explicit sol'n for $\hat{\theta}_i$ & $\tilde{\mu}$: need use MCMC

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i \mid \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d; $\mu \sim N(\mu_0, \tau^2)$
- $v_i, i = 1, \dots, n, \sigma^2, \mu_0, \tau^2$ all known

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i \mid \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d; $\mu \sim N(\mu_0, \tau^2)$
- $v_i, i = 1, \dots, n, \sigma^2, \mu_0, \tau^2$ all known
- $$E(\theta_i \mid \mathbf{x}) = x_i \frac{\sigma^2}{\sigma^2 + v_i} + E(\mu \mid \mathbf{x}) \left(1 - \frac{\sigma^2}{\sigma^2 + v_i}\right)$$

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i \mid \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d; $\mu \sim N(\mu_0, \tau^2)$
- $v_i, i = 1, \dots, n, \sigma^2, \mu_0, \tau^2$ all known

•

$$E(\theta_i \mid \mathbf{x}) = x_i \frac{\sigma^2}{\sigma^2 + v_i} + E(\mu \mid \mathbf{x}) \left(1 - \frac{\sigma^2}{\sigma^2 + v_i}\right)$$

•

$$E(\mu \mid \mathbf{x}) = \frac{\mu_0/\tau^2 + \sum x_i / (\sigma^2 + v_i)}{1/\tau^2 + \sum 1 / (\sigma^2 + v_i)}$$

- $x_i \mid \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i \mid \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d; $\mu \sim N(\mu_0, \tau^2)$
- $v_i, i = 1, \dots, n, \sigma^2, \mu_0, \tau^2$ all known
- $E(\theta_i \mid \mathbf{x}) = x_i \frac{\sigma^2}{\sigma^2 + v_i} + E(\mu \mid \mathbf{x})(1 - \frac{\sigma^2}{\sigma^2 + v_i})$
- $E(\mu \mid \mathbf{x}) = \frac{\mu_0/\tau^2 + \sum x_i / (\sigma^2 + v_i)}{1/\tau^2 + \sum 1 / (\sigma^2 + v_i)}$
- If σ^2 unknown, then need to sample from the posterior, no closed form available

- $x_i | \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i | \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d;
- $v_i, i = 1, \dots, n, \sigma^2, \mu_0, \tau^2$ all known

$$\mu \sim N(\mu_0, \tau^2)$$

-

$$E(\theta_i | \mathbf{x}) = x_i \frac{\sigma^2}{\sigma^2 + v_i} + E(\mu | \mathbf{x}) \left(1 - \frac{\sigma^2}{\sigma^2 + v_i}\right)$$

-

$$E(\mu | \mathbf{x}) = \frac{\mu_0 / \tau^2 + \sum x_i / (\sigma^2 + v_i)}{1 / \tau^2 + \sum 1 / (\sigma^2 + v_i)}$$

- If σ^2 unknown, then need to sample from the posterior, no closed form available

$\hat{\mu}$: arg max
 $E(\mu | \mathbf{x}^n)$

- Figure 11.11 applies similar ideas, plus sampling from the posterior, in logistic regression

$$m(\tilde{\mathbf{x}}^n)$$

622

11 · Bayesian Models

each hosp sep.

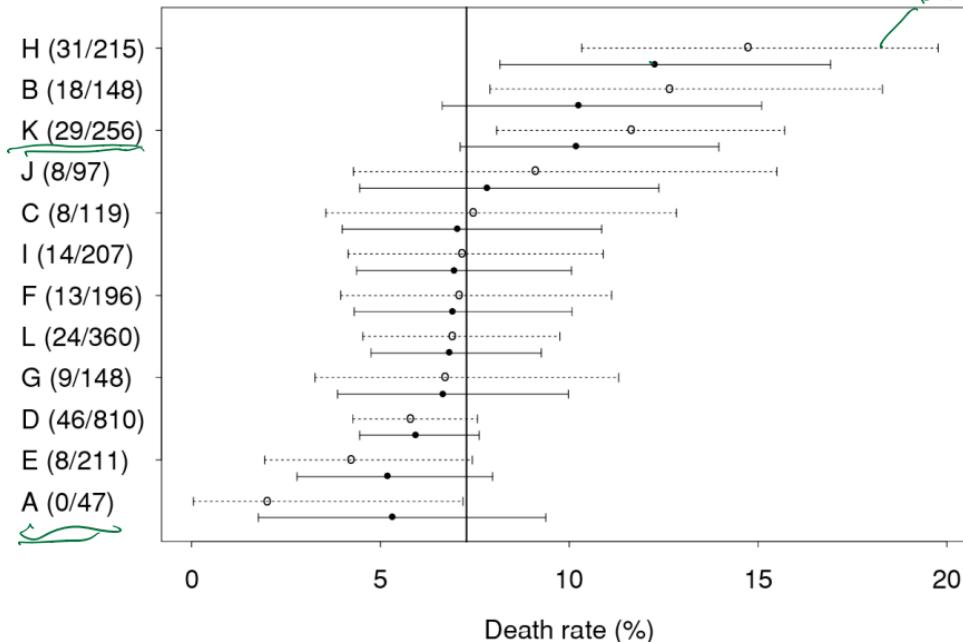


Figure 11.11 Posterior summaries for mortality rates for cardiac surgery data. Posterior means and 0.95 equitailed credible intervals for separate analyses for each hospital are shown by hollow circles and dotted lines, while blobs and solid lines show the corresponding quantities for a hierarchical model. Note the shrinkage of the estimates for the hierarchical model towards the overall posterior mean rate, shown as the solid vertical line; the hierarchical intervals are slightly shorter than those for the simpler model.

As above, $x_i | \theta_i \sim N(\theta_i, v_i)$, independent; $\theta_i | \mu \sim N(\mu, \sigma^2), i = 1, \dots, n$ i.i.d;

$$\theta_i | x_i \sim N\left(x_i \frac{\sigma^2}{\sigma^2 + v_i} + \mu \frac{v_i}{\sigma^2 + v_i}, \frac{\sigma^2 v_i}{\sigma^2 + v_i}\right)$$

$$\pi(\theta_i | x_i) = f(x_i | \theta_i) \pi(\theta_i) / \underbrace{\int f(x_i | \theta_i) \pi(\theta_i) d\theta_i}_{\text{marginal}} = \frac{f(x_i)}{f(x_i)}$$

$$\prod_{i=1}^n f_{x_i}(\underline{x}_i | \mu) = L(\mu, \underline{x}^n)$$

$$\text{(not bc ind't)} \quad \hat{\mu} = \hat{\mu}(\underline{x}^n) = \frac{\sum x_i / (v_i + \sigma^2)}{\sum 1 / (v_i + \sigma^2)}$$

SM "x_i marginally ind't, $N[\mu, v_i / (\sigma^2 + v_i)]$ "

$$\tilde{\theta}_i = \underbrace{\left(x_i \frac{\sigma^2}{\sigma^2 + v_i} \right)}_{\text{weight}} + \hat{\mu} \underbrace{\left(\frac{v_i}{\sigma^2 + v_i} \right)}_{\text{weight}}$$

$$= \hat{\mu} + (1 - \tilde{\gamma}_i)(x_i - \hat{\mu})$$

$$\tilde{\gamma}_i = \frac{v_i}{\sigma^2 + v_i} \quad \text{weight} \quad \begin{array}{l} \sigma^2 \rightarrow \infty \\ v_i \rightarrow \infty \end{array}$$

Using weight density for x^n
to estimate parameters in the prior
is called empirical Bayes

"large sets of parallel situations carry their own prior information"

$$\begin{aligned} E \tilde{\theta}_i &\neq \theta_i = \theta_i \frac{\sigma^2}{\sigma^2 + v_i} + \mu \frac{v_i}{\sigma^2 + v_i} \\ \text{bias} &\neq 0 \quad \text{perhaps } \text{mse}(\tilde{\theta}_i) \\ &= E (\tilde{\theta}_i - \theta_i)^2 < \text{Var}(x_i) \end{aligned}$$

$$x_i | \theta_i \sim N(\theta_i, v_i) \quad \uparrow$$

Alternative to ridge regression

$$\begin{aligned} \sum z_i &= 0 \\ \sum z_i^2 &= 1 \quad \uparrow \\ \leftarrow \frac{z_i - \bar{z}}{s_z} \quad \hat{\beta}_{\text{ridge}} &= (\mathbf{Z}^T \mathbf{Z} + \lambda \mathbf{I})^{-1} \mathbf{Z}^T \mathbf{y} \quad \leftarrow \text{shrinkage est.}'s \end{aligned}$$

$$\min_{\beta} \left\{ \sum_{i=1}^n (x_i - z_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

β_j = change in j
pointing
change in z_j

76

Empirical Bayes

Table 6.1 Counts y_x of number of claims x made in a single year by 9461 automobile insurance policy holders. Robbins' formula (6.7) estimates the number of claims expected in a succeeding year, for instance 0.168 for a customer in the $x = 0$ category. Parametric maximum likelihood analysis based on a gamma prior gives less noisy estimates.

Claims x	0	1	2	3	4	5	6	7
Counts y_x	7840	1317	239	42	14	4	4	1
Formula (6.7)	.168	.363	.527	1.33	1.43	6.00	1.75	
Gamma MLE	.164	.398	.633	.87	1.10	1.34	1.57	

??

The handwritten annotations include:

- A green arrow points from the label "Counts y_x " to the first column of the table.
- Three green arrows point from the label "Formula (6.7)" to the second column of the table.
- Two green arrows point from the label "Gamma MLE" to the third column of the table.
- Handwritten numbers are circled in green: ".168" (under Formula), ".363" (under Formula), ".398" (under Gamma MLE), "6.00" (under Formula), "1.75" (under Formula), and "1.57" (under Gamma MLE).
- A large green circle encloses the entire row of expected values (Formula).
- A large green circle encloses the entire row of maximum likelihood estimates (MLE).
- A green arrow points from the circled "1.75" to the circled "1.57".
- A green arrow points from the circled "1.57" to the circled "6.00".
- A green arrow points from the circled "6.00" to the circled "1.75".
- A green arrow points from the circled "1.75" back to the circled "1.57".
- A question mark "?" is written below the circled "6.00".

$$f_k(x; \theta_k) = \frac{\theta_k^x e^{-\theta_k}}{x!} \quad x=0, 1, 2, \dots$$

$$\theta_k^0 e^{-\theta_k} / \theta_k! = e^{-\theta_k}$$

$\nwarrow P_n$

k - customer k
had x claims
in the year

$$\theta_k \sim \Gamma(\nu, 1/\sigma)$$

density $\theta_k^{\nu-1} e^{-\theta_k/\sigma} / P(x) \sigma^\nu$

marginal density $f(x) = \prod_{k=1}^K \int f(x|\theta_k) \pi(\theta_k) d\theta_k$

$$= f(x | \nu, \sigma)$$

now use ML to estimate γ, σ

$$\ln p(\underline{x} | \gamma, \sigma) \propto f(\underline{x} | \gamma, \sigma)$$

estimates in the table are

$$E(\theta_i | \underline{x}) \quad \text{using } \hat{\sigma}, \hat{\gamma} \text{ in the prior}$$

Other option is the prior is nonparametric

Nonparametric empirical Bayes

$$f(x) = \int_0^{\infty} \frac{e^{-\theta} \theta^x}{x!} \pi(\theta) d\theta$$

$f(x)$
|||

$$E(\theta|x) = \frac{\int \underline{\theta} f(\theta|x) \pi(\theta) d\theta}{\int f(\theta|x) \pi(\theta) d\theta}$$

|||

- Loss function

$$L(\theta; \hat{\theta}) : \mathbb{H} \times \mathbb{H} \rightarrow \mathbb{R}$$

θ unk. par. in model
 $\hat{\theta}$ estimate

- Risk of an estimator

$L(\theta; \hat{\theta}) = (\theta - \hat{\theta})^2$ sq. err

$\frac{1}{2} |\theta - \hat{\theta}|$ abs. err

- Mean-squared error

$\text{class } \approx$

expected loss

$\begin{cases} \theta & , \hat{\theta} = \theta \\ 1 & , \hat{\theta} \neq \theta \end{cases}$ ($\hat{\theta} = \theta$ loss)

$|\theta - \hat{\theta}|^p$ L_p loss

- Bayes risk

Kullback-Leibler loss

"decision" is

a choice of estimator

$$\hat{\theta} = \hat{\theta}(x).$$

$$L(\theta, \hat{\theta}) = \int \left(\log \frac{f(x; \theta)}{f(x; \hat{\theta})} \right) f(x; \theta) dx$$

- Loss function

$$L(\theta; \hat{\theta})$$

unk. choose

$$f(x; \theta)$$

- Risk of an estimator

$$R(\theta; \hat{\theta}) = \int L(\theta; \hat{\theta}) f(x; \theta) dx = E_{\theta} \{ L(\theta, \hat{\theta}(x)) \}$$

est. $\theta \leftarrow$ undermodel

- Mean-squared error

↑ under sq'd error loss

expected loss

- Bayes risk

$$R(\theta; \hat{\theta}) = \int \{ \theta - \hat{\theta}(x) \}^2 f(x; \theta) dx$$

- maximum risk

$$\begin{aligned} &= E \{ \hat{\theta}(x) - \theta \}^2 = E_{\theta} \{ \hat{\theta}(x) - \underline{\theta} \\ &\quad + \underline{\theta} - \theta \}^2 \\ &= \text{var} \{ \hat{\theta}(x) \} + \text{bias}^2 \{ \hat{\theta}(x) \} \end{aligned}$$

$$\bar{X} \sim N(\theta, 1)$$

$$\hat{\theta}_1 = \bar{X}$$

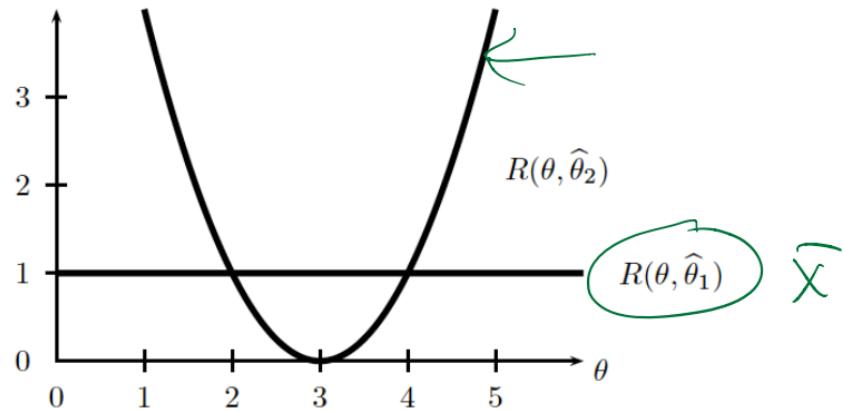
$$\hat{\theta}_2 = 3$$

$$R_1(\theta; \hat{\theta}) = E\{\theta - \hat{\theta}(x)\}^2$$

$$1 = E(\theta - \bar{X})^2 = \frac{1}{n}$$

$$2 = E(\theta - 3)^2 = (\theta - 3)^2$$

12.2 Comparing Risk Functions 195

FIGURE 12.1. Comparing two risk functions. Neither risk function dominates the other at all values of θ .

often can't find an estimator
for which $R(\theta; \hat{\theta})$ is smallest $\forall \theta \in \Theta$

Examples

AoS 12.2, 12.3

12.3

$$X \sim \text{Bin}(p_1, p)$$

$$\hat{p}_1 = \bar{x} = \frac{x}{n}$$

$$R(p_1, \hat{p}_1) = \text{var } \bar{x} = \frac{p(1-p)}{n}$$

$$\hat{p}_2 = \frac{x + \alpha}{n + \alpha + \beta} \quad \begin{array}{l} \text{post. mean} \\ \text{if } B(\alpha, \beta) \end{array}$$

$$R_{\hat{p}_2}(p, \hat{p}_2) = \text{var } \hat{p}_2 + \left(E(\hat{p}_2) \right)^2$$

$$= \text{see Eq. below Fig 12.2}$$

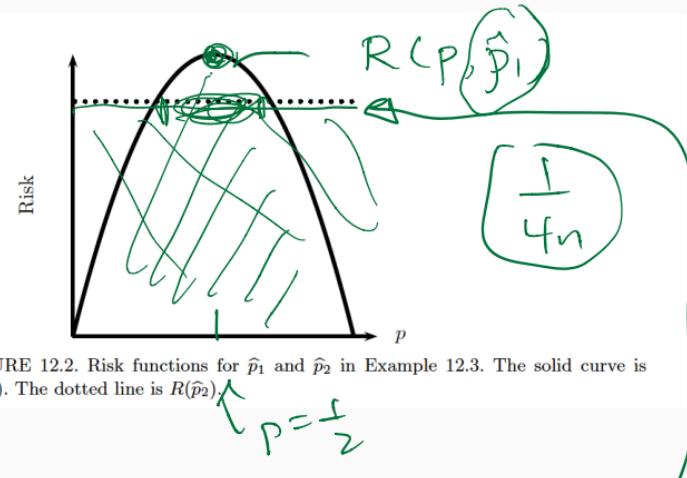


FIGURE 12.2. Risk functions for \hat{p}_1 and \hat{p}_2 in Example 12.3. The solid curve is $R(\hat{p}_1)$. The dotted line is $R(\hat{p}_2)$.

$$\alpha = \beta = \sqrt{n}/4$$

$$R_{\hat{p}_2}(p, \hat{p}_2) = \frac{n}{4(n+\sqrt{n})^2}$$

$$R(\theta, \hat{\theta}(x)) = \int L(\theta, \hat{\theta}(x)) f(x; \theta) dx \quad \leftarrow \text{still dep}$$

$$R(\theta, d(x))$$

$$\hat{p}_1 = \bar{x};$$

$$R(p, \hat{p}_1) = p(1-p)/n$$

$$\hat{p}_2 = \frac{n\bar{X} + \sqrt{n/4}}{n + \sqrt{n}};$$

$$R(p, \hat{p}_2) = \frac{n}{4(n + \sqrt{n})^2}$$

Bin

$$\bar{R}(\hat{\theta}(x)) = \sup_{\theta \in \Theta} R\{\theta, \hat{\theta}(x)\} \quad \text{max risk}$$

minimax estimator is defined as $\inf_{\tilde{\theta} \in \tilde{\Theta}} \sup_{\theta \in \Theta} R(\theta, \tilde{\theta})$
 "minimax lower bd" "minimax decision rule"

Bayes rules

Bayes risk
of $\hat{\theta}$

$$\int R(\theta, \hat{\theta}) \pi(\theta) d\theta = s_\pi(\hat{\theta})$$
$$= r(f, \hat{\theta}) = r(\pi, \hat{\theta})$$

Bayes est. is the f^* $\hat{\theta}(x)$ that satisfies

$$\min_{\tilde{\theta} \in \Theta} r(f, \tilde{\theta}) = \min_{\tilde{\theta} \in \Theta} r_\pi(\tilde{\theta})$$

$$= \inf_{\tilde{\theta}} s_\pi(\tilde{\theta})$$

Often called Bayes rules (decisions > estimators)

~~Admissibility~~

How to find Bayes rules? (i.e. Bayes estimators)

$$\text{def} \hat{\theta} : \inf_{\tilde{\theta}} R_{\pi}(\tilde{\theta}) = \inf_{\tilde{\theta}} \int_{\Theta} R(\theta, \tilde{\theta}(x)) \pi(\theta) d\theta$$

Theorem 12.7 : note typo in stmt,

$$\theta \leftarrow \hat{\theta}$$

Exp'd loss under posterior

$$= \int L(\theta, \hat{\theta}(x)) \pi(\theta|x) d\theta \quad \text{dep. } \hat{\theta}, x$$

choose $\hat{\theta}_\pi(x)$ to $\min_{\hat{\theta}} \int L(\theta, \hat{\theta}(x)) \pi(\theta|x) d\theta$

What's the Bayes risk of $\hat{\theta}_\pi(x)$?

$$= \int R(\theta, \hat{\theta}_\pi(x)) \pi(\theta) d\theta = r_\pi(\hat{\theta}) \quad \text{dep. on } \pi$$

$$\min_{\hat{\theta}} \int R(\theta, \hat{\theta}_\pi(x)) \pi(\theta) d\theta \quad \leftarrow$$

$$= \min_{\hat{\theta}} \int \left\{ \int L(\theta, \hat{\theta}_\pi(x)) f(x|\theta) dx \right\} \pi(\theta) d\theta$$

$$f(x|\theta) \pi(\theta) = \pi(\theta|x) f(x)$$

$$= \min_{\hat{\theta}} \int \int L(\theta, \hat{\theta}_\pi(x)) \pi(\theta|x) f(x) dx d\theta$$

$$= \min_{\hat{\theta}} \int f(x) \int L(\theta, \hat{\theta}_\pi(x)) \pi(\theta|x) d\theta dx$$

$$\text{only need to } \min_{\hat{\theta}} \int L(\theta, \hat{\theta}_\pi(x)) \pi(\theta|x) d\theta \quad \leftarrow$$

Conclusion: "Bayes rule" w/ prior = best posterior estimate

$$\min_{\hat{\theta}} \int (\theta - \hat{\theta})^2 \pi(\theta|x) d\theta$$

$$= \dots = \hat{\theta} = E(\theta|x) \text{ answer}$$

admissibility \leftarrow needed w/ Friday

\cup \equiv $)$