

Caterpillar Regression Example: Conjugate Priors, Conditional & Marginal Posteriors, Predictive Distribution, Variable Selection

Prof. Nicholas Zabaras

Center for Informatics and Computational Science

<https://cics.nd.edu/>

*University of Notre Dame
Notre Dame, Indiana, USA*

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

September 24, 2018



Contents

- The Caterpillar Regression Problem, Linear Regression Model, MLE Estimator, Likelihood, MLE Estimates for our example, Conjugate Priors, Conditional and Marginal Posteriors, Ridge regression, Predictive Distribution, Implementation, Influence of the Conjugate Prior
- Zellner's G Prior, Marginal Posterior Mean and Variance, Predictive Modeling, Credible Intervals
- Jeffreys' non-informative Prior, Credible Intervals, Zellner's G Prior Marginal Distribution of y , Point Null Hypothesis and Bayes Factors
- Variable Selection, Model Competition, Variable Selection-Prior, Stochastic Search for the Most Probable Model, Gibb's Sampling for Variable Selection, Implementation

C. P. Robert, *The Bayesian Core*, Springer, 2nd edition, [chapter 3](#) (full text available)

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)



Regression

- Regression refers to statistical analysis that deals with the representation of dependencies between several variables.
- In particular, we want to find a representation of the distribution $f(y|\theta, \mathbf{x})$ of an observable variable y given a vector of observables \mathbf{x} , using samples of $(\mathbf{x}_i, y_i), i = 1, \dots, n$.
- Here, we will consider a particular example of modeling dependencies in pine processionary caterpillar colony size.



Linear Regression Models

- In linear regression, we analyze the linear influence of some variables on others.
- In our particular example of pine processionary caterpillar colonies,
 - **Response Variable:** y is the number of processionary caterpillar colonies
 - **Explanatory Variables:** Covariates $x = (x_1, x_2, \dots, x_k)$ as defined next (in general can be continuous, discrete or mixed type)



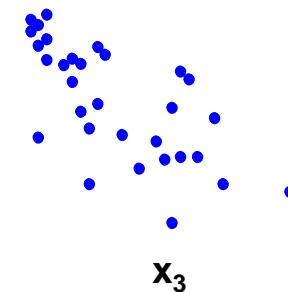
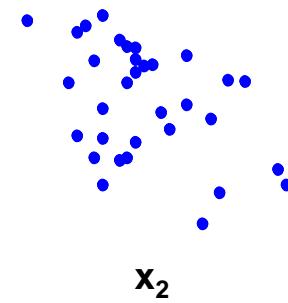
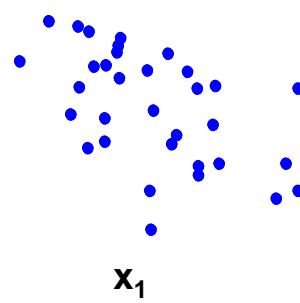
Caterpillar Regression Problem

The pine processionary caterpillar colony size (y =number of nests) is influenced by the following explanatory variables:

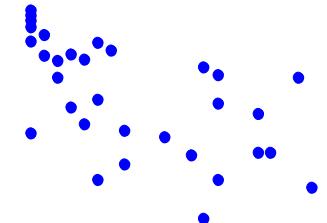
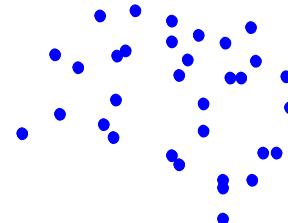
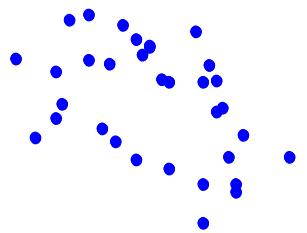
- x_1 is the altitude (in meters)
- x_2 is the slope (in degrees)
- x_3 is the number of pines in the square
- x_4 is the height (in meters) of the tree sampled at the center of the square
- x_5 is the diameter of the tree sampled at the center of the square
- x_6 is the index of the settlement density
- x_7 is the orientation of the square (from 1 if southbound to 2 otherwise)
- x_8 is the height (in meters) of the dominant tree
- x_9 is the number of vegetation strata
- x_{10} is the mix settlement index (from 1 if not mixed to 2 if mixed).



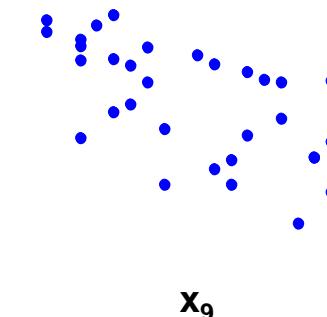
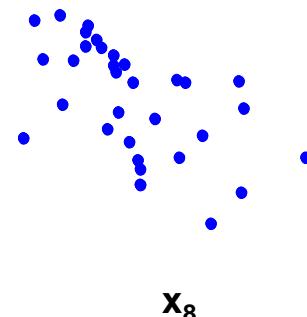
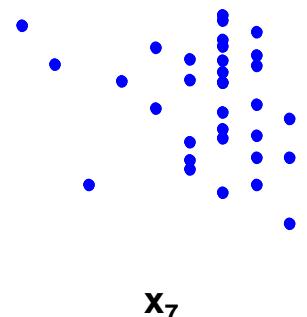
Caterpillar Regression Problem



Semilog- y plot
of the data
 $(x_i, y), i = 1, \dots, 9$



An implementation
is available
[MatLab](#), [C++](#)



Regression

- The distribution of y given x is considered in the context of a set of experimental data, $i = 1, \dots, n$, on which both y , and x_{i1}, \dots, x_{ik} are measured.
- The dataset is made from the outcomes $y = (y_1, \dots, y_n)$ and of the $n \times (k + 1)$ matrix of **explanatory variables**

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ 1 & x_{31} & x_{32} & \dots & x_{3k} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}$$



Linear Regression

- The most common linear regression model is of the form:

$$y | \beta, \sigma^2, X \sim \mathcal{N}_n(X\beta, \sigma^2 I_n)$$

- From this, we conclude that:

$$\mathbb{E}[y_i | \beta, X] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

$$Var[y_i | \sigma^2, X] = \sigma^2$$

- Note the difference between finite valued regressors x_i (like in our caterpillar problem) and categorical variables which also take finite number of values but with range that has no numerical meaning.
- Makes no sense to involve x directly in the regression: replace the single regressor x [belonging in $\{1, \dots, m\}$, say] with m indicator (or dummy) variables.

Categorical Variables

- Note the difference between finite valued regressors x_i (like in our caterpillar problem) and categorical variables which also take finite number of values but with range that has no numerical meaning.
- Makes no sense to use x directly in the regression: replace the single regressor x [in $\{1, \dots, m\}$, say] with m indicator variables

$$x_1 = \mathbb{I}_1(x), x_2 = \mathbb{I}_2(x), \dots, x_m = \mathbb{I}_m(x)$$

- Use different β_i for each class categorical variable value:

$$\mathbb{E}[y_i | \boldsymbol{\beta}, X] = \dots + \beta_1 \mathbb{I}_1(x) + \dots + \beta_m \mathbb{I}_m(x) + \dots$$

- *Identifiability* requires eliminating one of the classes (e.g. $\beta_1 = 0$) since

$$\sum_i \mathbb{I}_i(x) = 1$$



Linear Regression

- Returning back to our linear regression model:

$$y | \boldsymbol{\beta}, \sigma^2, X \sim \mathcal{N}_n(X\boldsymbol{\beta}, \sigma^2 I_n)$$

- From this, we conclude that:

$$\mathbb{E}[y_i | \boldsymbol{\beta}, X] = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

$$Var[y_i | \sigma^2, X] = \sigma^2$$

- Assume that $k + 1 < n$ and that X is of full rank: $rank(X) = k + 1$

- X is of full rank if and only if $X^T X$ is invertible

- The likelihood is then:

$$\ell(\boldsymbol{\beta}, \sigma^2 | y, X) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}(y - X\boldsymbol{\beta})^T(y - X\boldsymbol{\beta})\right)$$



MLE Estimator

- The MLE of β is the solution of the least squares minimization problem

$$\min_{\beta} (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta) = \min_{\beta} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_k x_{ik})^2 \Rightarrow \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Since $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ is a linear transform of $\mathbf{y} \sim \mathcal{N}_n(\mathbf{X}\beta, \sigma^2 \mathbf{I}_n)$

$$\hat{\beta} \sim \mathcal{N}((\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X}\beta, (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \sigma^2 \mathbf{I}_n \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1}) = \mathcal{N}(\beta, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

- Thus $\hat{\beta}$ is an unbiased estimator and $Var(\hat{\beta} | \sigma^2, \mathbf{X}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$
- Also $\hat{\beta}$ is the best linear unbiased estimator of β :

$\alpha \in \mathbb{R}^{k+1}, Var(\alpha^T \hat{\beta} | \sigma^2, \mathbf{X}) \leq Var(\alpha^T \tilde{\beta} | \sigma^2, \mathbf{X})$ where $\tilde{\beta}$ is any unbiased linear estimator of β

MLE Estimator

- The MLE of σ^2 is the solution of:

$$\min_{\sigma^2} \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \frac{n}{2} \ln \sigma |_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} = 0 \Rightarrow$$

$$\sigma_{MLE}^2 = \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})}{n}$$

- Let $\mathbf{M} = \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ be the projection $n \times n$ matrix of the data \mathbf{y} on the column space of \mathbf{X}

$$\begin{aligned} \mathbb{E}(\sigma_{MLE}^2) &= \frac{\mathbb{E}[(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})]}{n} = \frac{\mathbb{E}[(\mathbf{y} - \mathbf{M}\mathbf{y})^T (\mathbf{y} - \mathbf{M}\mathbf{y})]}{n} = \\ &= \frac{\mathbb{E}[\mathbf{y}^T (\mathbf{I} - \mathbf{M})^T (\mathbf{I} - \mathbf{M}) \mathbf{y}]}{n} = \frac{\mathbb{E}[\mathbf{y}^T (\mathbf{I} - \mathbf{M}) \mathbf{y}]}{n} \end{aligned}$$

- You can easily show that (see this slide):

For \mathbf{Y} a n -dim random vector with $\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}$, $Cov[\mathbf{Y}] = \mathbf{V}$,
and \mathbf{A} a $n \times n$ matrix then : $\mathbb{E}[\mathbf{Y}^T \mathbf{A} \mathbf{Y}] = \text{tr}(\mathbf{A} \mathbf{V}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$

MLE Estimator

$$\mathbb{E}(\sigma_{MLE}^2) = \frac{\mathbb{E}[y^T(I - M)y]}{n}$$

- You can easily show that ([see this slide](#)):

For \mathbf{Y} a n -dim random vector with $\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}$, $Cov[\mathbf{Y}] = \mathbf{V}$, and \mathbf{A} a $n \times n$ matrix then : $\mathbb{E}[Y^T A Y] = \text{tr}(A\mathbf{V}) + \boldsymbol{\mu}^T A \boldsymbol{\mu}$

- Applying this we have:

$$\begin{aligned}\mathbb{E}[y^T (I - M)y] &= \text{tr} \left((I - M) \overbrace{\text{Cov}(y)}^{\sigma^2 I_n} \right) + \mathbb{E}(y)^T (I - M) \mathbb{E}(y) = \\ &= \sigma^2 \text{rank}(I - M) + \boldsymbol{\beta}^T X^T \underbrace{(I - M)X \boldsymbol{\beta}}_0 = \sigma^2 (n - k - 1)\end{aligned}$$

Note : $I - M$ has 1 as its eigenvalue with multiplicity $n - k - 1$ and

$$\text{tr}(I - M) = \text{tr}(I) - \text{tr}(M) = n - \text{tr}\left(X(X^T X)^{-1} X^T\right) = n - \text{tr}\left((X^T X)(X^T X)^{-1}\right) = n - \text{tr}(I_{k+1}) = n - k - 1$$

MLE Estimator

- The expectation of the MLE of σ^2 is then given as:

$$\mathbb{E}(\sigma_{MLE}^2) = \frac{\sigma^2}{n}(n-k-1)$$

- The unbiased estimator of σ^2 is thus given:

$$\hat{\sigma}^2 = \frac{1}{n-k-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \frac{s^2}{n-k-1}, \text{ where: } s^2 \triangleq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

- Indeed using the result from the earlier slide we have:

$$\mathbb{E}(\hat{\sigma}^2) = \frac{1}{n-k-1} \mathbb{E} \left[(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right] = \frac{\sigma^2(n-k-1)}{n-k-1} = \sigma^2$$

- To approximate the covariance of $\hat{\boldsymbol{\beta}}$, in $\hat{\boldsymbol{\beta}} = \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$

use:

$$Var(\hat{\boldsymbol{\beta}} | \sigma^2, \mathbf{X}) = \hat{\sigma}^2(\mathbf{X}^T \mathbf{X})^{-1}$$



Appendix

For \mathbf{Y} a n -dim random vector with $\mathbb{E}[\mathbf{Y}] = \boldsymbol{\mu}$, $\text{Cov}[\mathbf{Y}] = \mathbf{V}$, and \mathbf{A} a $n \times n$ matrix then : $\mathbb{E}[\mathbf{Y}^T \mathbf{A} \mathbf{Y}] = \text{tr}(\mathbf{A} \mathbf{V}) + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$

- We can prove this theorem quite easily as follows:

$$\mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{Y} - \boldsymbol{\mu})] = \mathbb{E}[\mathbf{Y}^T \mathbf{A} \mathbf{Y} + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{A} \mathbf{Y} - \mathbf{Y}^T \mathbf{A} \boldsymbol{\mu}] \Rightarrow$$

$$\mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{Y} - \boldsymbol{\mu})] = \mathbb{E}[\mathbf{Y}^T \mathbf{A} \mathbf{Y}] + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} - \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu} = \mathbb{E}[\mathbf{Y}^T \mathbf{A} \mathbf{Y}] - \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$$

- Using the identity $\mathbb{E}[\text{tr}(\mathbf{W})] = \text{tr}[\mathbb{E}(\mathbf{W})]$ for any random square matrix \mathbf{W} , we can write the l.h.s. of the equation above as:

$$\begin{aligned}\mathbb{E}[(\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{Y} - \boldsymbol{\mu})] &= \mathbb{E}[\text{tr}((\mathbf{Y} - \boldsymbol{\mu})^T \mathbf{A} (\mathbf{Y} - \boldsymbol{\mu}))] = \mathbb{E}[\text{tr}(\mathbf{A} (\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^T)] = \\ &\text{tr}[\mathbb{E}(\mathbf{A} (\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^T)] = \text{tr}[\mathbf{A} \mathbb{E}((\mathbf{Y} - \boldsymbol{\mu})(\mathbf{Y} - \boldsymbol{\mu})^T)] = \text{tr}(\mathbf{A} \mathbf{V})\end{aligned}$$

This completed the proof of the theorem.



T-Statistic

- We can define the standard t –statistic using the sampling distribution of the MLE as follows:

$$T_i = \frac{\hat{\beta}_i - \beta_i}{\sqrt{\hat{\sigma}^2 \omega_{ii}}} \sim \mathcal{T}(n - k - 1, 0, 1), \text{ where: } \omega_{(i,i)} = (\mathbf{X}^T \mathbf{X})^{-1}|_{(i,i)}$$

- This statistic can be used for hypothesis testing, e.g.

accept $H_0 : \beta_i = 0$ versus $H_1 : \beta_i \neq 0$ at the level α if

$$\frac{|\hat{\beta}_i|}{\sqrt{\hat{\sigma}^2 \omega_{ii}}} \leq F_{n-k-1}^{-1}(1 - \alpha/2), \text{ the } (1 - \alpha/2)\text{nd quantile of } \mathcal{T}_{n-k-1}$$

- We will see later in this lecture how the same results can be obtained using Jeffrey's uninformative prior.



T-Statistic & Frequentist Marginal Confidence Interval

- The frequentist argument in using this bound is that there is significant evidence against H_0 if the p -value is smaller than α .

$$p_i = P_{H_0} \left(|T_i| > |t_i| = \frac{|\hat{\beta}_i|}{\sqrt{\hat{\sigma}^2 \omega_{ii}}} \right) =$$
$$P_{H_0}(T_i < -|t_i|) + P_{H_0}(T_i > |t_i|) = F_{n-k-1}(-|t_i|) + (1 - F_{n-k-1}(|t_i|)) < \alpha$$

Reject H_0 if: $F_{n-k-1}(|t_i|) \geq 1 - \frac{\alpha}{2}$ or $|t_i| = \frac{|\hat{\beta}_i|}{\sqrt{\hat{\sigma}^2 \omega_{ii}}} \geq F_{n-k-1}^{-1}\left(1 - \frac{\alpha}{2}\right)$

- Finally the statistic T_i can be used to derive the frequentist marginal confidence interval is:

$$\left\{ \beta_i: |\beta_i - \hat{\beta}_i| \leq \sqrt{\hat{\sigma}^2 \omega_{ii}} F_{n-k-1}^{-1}\left(1 - \frac{\alpha}{2}\right) \right\}$$



MLE (Least Squares) Estimates

$$\hat{\beta}_i \quad \sqrt{\hat{\sigma}^2 \omega_{ii}} \quad t_i = \frac{\hat{\beta}_i}{\sqrt{\hat{\sigma}^2 \omega_{ii}}} \quad P_{H_0}\left(|T_i| > \frac{|\hat{\beta}_i|}{\sqrt{\hat{\sigma}^2 \omega_{ii}}}\right) \quad p_i =$$

Coefficients	Estimate	Std.Error	t-value	Pr(> t)
intercept	10.998412	3.060892	3.594	0.00161 **
XV1	-0.004431	0.001557	-2.846	0.00939 **
XV2	-0.053830	0.021904	-2.458	0.02232 *
XV3	0.067939	0.099492	0.683	0.50174
XV4	-1.293636	0.563925	-2.294	0.03168 *
XV5	0.231637	0.104399	2.219	0.03709 *
XV6	-0.356800	1.566782	-0.228	0.82193
XV7	-0.237469	1.006210	-0.236	0.81558
XV8	0.181060	0.236772	0.765	0.45248
XV9	-1.285316	0.865023	-1.486	0.15142
XV10	-0.433106	0.735018	-0.589	0.56162

An implementation is available [MatLab](#), [C++](#)



Likelihood

$$\ell(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

□ Note that the likelihood above can now be written in the following very useful form:

$$\ell(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) = \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} \left(\mathbf{y} - \mathbf{X}((\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \hat{\boldsymbol{\beta}})\right)^T \left(\mathbf{y} - \mathbf{X}((\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) + \hat{\boldsymbol{\beta}})\right)\right) \Rightarrow$$

$$\begin{aligned} \ell(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) - \frac{1}{2\sigma^2} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})\right) = \\ &= \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2} s^2 - \frac{1}{2\sigma^2} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})\right) \end{aligned}$$

where: $s^2 \triangleq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$,
 $\hat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$



Conjugate Priors

- Observing the form of the likelihood suggests the following conjugate prior:

$$\begin{aligned}\beta | \sigma^2, X \sim \mathcal{N}_{k+1}(\tilde{\beta}, \sigma^2 M^{-1}), M \text{ a } (k+1) \times (k+1) \text{ pos. def. symm. matrix} \\ \sigma^2 | X \sim \text{InvGamma}(a, b), a, b > 0\end{aligned}$$

- The posterior is then

$$\begin{aligned}\pi(\beta, \sigma^2 | \hat{\beta}, s^2, X) &\propto \frac{1}{(2\pi\sigma^2)^{n/2}} \exp\left(-\frac{1}{2\sigma^2}s^2 - \frac{1}{2\sigma^2}(\beta - \hat{\beta})^T X^T X (\beta - \hat{\beta})\right) \\ &\times \frac{1}{(2\pi\sigma^2)^{(k+1)/2}} \exp\left(-\frac{1}{2\sigma^2}(\beta - \tilde{\beta})^T M (\beta - \tilde{\beta})\right) \times \left(\frac{1}{\sigma^2}\right)^{a+1} \exp\left(-\frac{1}{\sigma^2}b\right) = \\ &= \sigma^{-k-1-2a-2-n} \exp\left(-\frac{1}{2\sigma^2}\{s^2 + 2b + (\beta - \hat{\beta})^T X^T X (\beta - \hat{\beta}) + (\beta - \tilde{\beta})^T M (\beta - \tilde{\beta})\}\right) = \\ &= \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2}\{s^2 + 2b + \hat{\beta}^T (M + X^T X)\beta - 2\beta^T (M\tilde{\beta} + X^T X\hat{\beta}) + \tilde{\beta}^T M\tilde{\beta} + \hat{\beta}^T (X^T X)\hat{\beta}\}\right)\end{aligned}$$

Conjugate Priors: Posterior of β given σ^2

$$\pi(\beta, \sigma^2 | \hat{\beta}, s^2, X) \propto \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2}\{s^2 + 2b + \beta^T(M + X^T X)\beta - 2\beta^T(M\tilde{\beta} + X^T X\hat{\beta}) + \tilde{\beta}^T M\tilde{\beta} + \hat{\beta}^T(X^T X)\hat{\beta}\}\right) =$$
$$\sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2}\left\{s^2 + 2b + \left(\beta - (M + X^T X)^{-1}(M\tilde{\beta} + X^T X\hat{\beta})\right)^T(M + X^T X)\left(\beta - (M + X^T X)^{-1}(M\tilde{\beta} + X^T X\hat{\beta})\right) + \tilde{\beta}^T M\tilde{\beta} + \hat{\beta}^T X^T X\hat{\beta} - (M\tilde{\beta} + X^T X\hat{\beta})^T(M + X^T X)^{-T}(M + X^T X)(M + X^T X)^{-1}(M\tilde{\beta} + X^T X\hat{\beta})\right\}\right)$$

□ Based on this, the posterior of β given σ^2 is:

$$\pi(\beta | \sigma^2, y, X) \propto \exp\left(-\frac{1}{2\sigma^2}(\beta - \mathbb{E}[\beta | y, X])^T(M + X^T X)(\beta - \mathbb{E}[\beta | y, X])\right)$$

$$\beta | \sigma^2, y, X \sim \mathcal{N}_{k+1}\left((M + X^T X)^{-1}\{X^T X\hat{\beta} + M\tilde{\beta}\}, \sigma^2(M + X^T X)^{-1}\right)$$

where

$$\mathbb{E}[\beta | \sigma^2, y, X] = (M + X^T X)^{-1}(M\tilde{\beta} + X^T X\hat{\beta})$$

and

$$Var[\beta | \sigma^2, y, X] = \sigma^2(M + X^T X)^{-1}$$



Conjugate Priors: Marginal posterior of σ^2

$$\pi(\boldsymbol{\beta}, \sigma^2 | \widehat{\boldsymbol{\beta}}, s^2, \mathbf{X}) \propto \\ \sigma^{-k-n-2a-3} \exp\left(-\frac{1}{2\sigma^2} \left\{ s^2 + 2b + (\boldsymbol{\beta} - \mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}])^T (\mathbf{M} + \mathbf{X}^T \mathbf{X}) (\boldsymbol{\beta} - \mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]) \right. \right. \\ \left. \left. + \widetilde{\boldsymbol{\beta}}^T \mathbf{M} \widetilde{\boldsymbol{\beta}} + \widetilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}} \right. \right. \\ \left. \left. - (\mathbf{M} \widetilde{\boldsymbol{\beta}} + \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}})^T (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{M} \widetilde{\boldsymbol{\beta}} + \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}}) \right\} \right)$$

□ Integrating out $\boldsymbol{\beta}$ gives the following marginal for σ^2 :

$$\pi(\sigma^2 | \widehat{\boldsymbol{\beta}}, s^2, \mathbf{X}) \\ \propto \sigma^{-n-2a-2} \exp\left(-\frac{1}{2\sigma^2} \left\{ \widetilde{\boldsymbol{\beta}}^T \mathbf{M} \widetilde{\boldsymbol{\beta}} + \widetilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}} + s^2 \right. \right. \\ \left. \left. + 2b - (\mathbf{M} \widetilde{\boldsymbol{\beta}} + \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}})^T (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{M} \widetilde{\boldsymbol{\beta}} + \mathbf{X}^T \mathbf{X} \widetilde{\boldsymbol{\beta}}) \right\} \right)$$

□ To simplify, we will use the following two identities

$$A: (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} = (\mathbf{X}^T \mathbf{X})^{-1} - (\mathbf{X}^T \mathbf{X})^{-1} \left(\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \right)^{-1} (\mathbf{X}^T \mathbf{X})^{-1}$$

$$B: \mathbf{X}^T \mathbf{X} (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \mathbf{M} = \left(\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1} \right)^{-1}$$

Conjugate Priors: Marginal Posterior of σ^2

$$\pi(\sigma^2 | \hat{\beta}, s^2, X) \propto \sigma^{-n-2a-2} \exp\left(-\frac{1}{2\sigma^2} \left\{ \tilde{\beta}^T M \tilde{\beta} + \hat{\beta}^T X^T X \hat{\beta} + s^2 + 2b - (M \tilde{\beta} + X^T X \hat{\beta})^T (M + X^T X)^{-1} (M \tilde{\beta} + X^T X \hat{\beta}) \right\}\right)$$

□ We can simplify the last term above using identities A & B:

$$(M \tilde{\beta} + X^T X \hat{\beta})^T (M + X^T X)^{-1} (M \tilde{\beta} + X^T X \hat{\beta}) = (\text{expand})$$

$$2\hat{\beta}^T X^T X (M + X^T X)^{-1} M \tilde{\beta} + \tilde{\beta}^T M^T (M + X^T X)^{-1} M \tilde{\beta} + \hat{\beta}^T X^T X (M + X^T X)^{-1} X^T X \hat{\beta} =$$

$$\overbrace{(X^T X)^{-1} - (X^T X)^{-1} (M^{-1} + (X^T X)^{-1})^{-1} (X^T X)^{-1}}$$

$$2\hat{\beta}^T X^T X (M + X^T X)^{-1} M \tilde{\beta} + \tilde{\beta}^T M^T (M + X^T X)^{-1} M \tilde{\beta} + \hat{\beta}^T X^T X ((X^T X)^{-1} - (X^T X)^{-1} (M^{-1} + (X^T X)^{-1})^{-1} (X^T X)^{-1}) X^T X \hat{\beta} =$$

$$2\hat{\beta}^T \underbrace{X^T X (M + X^T X)^{-1} M \tilde{\beta}}_{(M^{-1} + (X^T X)^{-1})^{-1}} + \tilde{\beta}^T \underbrace{M^T (M + X^T X)^{-1} M \tilde{\beta}}_{M - X^T X (M + X^T X)^{-1} M} + \hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}^T (M^{-1} + (X^T X)^{-1})^{-1} \hat{\beta} =$$

$$M - X^T X (M + X^T X)^{-1} M = M - (M^{-1} + (X^T X)^{-1})^{-1}$$

$$2\hat{\beta}^T (M^{-1} + (X^T X)^{-1})^{-1} \tilde{\beta} + \tilde{\beta}^T M \tilde{\beta} - \tilde{\beta}^T (M^{-1} + (X^T X)^{-1})^{-1} \tilde{\beta} + \hat{\beta}^T X^T X \hat{\beta} - \hat{\beta}^T (M^{-1} + (X^T X)^{-1})^{-1} \hat{\beta} =$$
$$\hat{\beta}^T M \tilde{\beta} + \hat{\beta}^T X^T X \hat{\beta} - (\tilde{\beta} - \hat{\beta})^T (M^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})$$

□ Finally:

$$\pi(\sigma^2 | \hat{\beta}, s^2, X) \propto \sigma^{-n-2a-2} \exp\left\{-\frac{1}{2\sigma^2} \left((\tilde{\beta} - \hat{\beta})^T (M^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta}) + s^2 + 2b \right)\right\}$$

Conjugate Priors: Marginal Posterior of σ^2

$$\pi(\sigma^2 | \tilde{\beta}, s^2, X) \propto \sigma^{-n-2a-2} \exp \left\{ -\frac{1}{2\sigma^2} \left((\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta}) + s^2 + 2b \right) \right\}$$

□ The marginal posterior is:

$$\sigma^2 | y, X \sim \text{InvGamma} \left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})}{2} \right)$$

□ The marginal posterior mean for $n \geq 2$ is then:

$$\mathbb{E}^\pi[\sigma^2 | y, X] = \frac{2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})}{n + 2a - 2}$$

Inverse-gamma	$\theta \sim \text{Inv-gamma}(\alpha, \beta)$ $p(\theta) = \text{Inv-gamma}(\theta \alpha, \beta)$	shape $\alpha > 0$ scale $\beta > 0$
---------------	---	---

$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$	$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$ $\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$ $\text{mode}(\theta) = \frac{\beta}{\alpha+1}$
--	--

Conjugate Priors: MLE, Posterior Mean and Ridge Estimator

- Note that setting $\mathbf{M} = \mathbf{I}_{k+1}/c$, $c > 0$ and $\tilde{\boldsymbol{\beta}} = \mathbf{0}_{k+1}$ in the conditional posterior mean

$$\mathbb{E}[\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}] = (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{M} \tilde{\boldsymbol{\beta}} + \mathbf{X}^T \mathbf{y})$$

we obtain the classical Ridge Regression estimate:

$$\mathbb{E}[\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}] = \left(\frac{1}{c} \mathbf{I}_{k+1} + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y} = \left(\frac{1}{c} \mathbf{I}_{k+1} + \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \mathbf{y}$$

- The general estimator $\mathbb{E}[\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}]$ can be seen as a weighted average of the prior mean and the MLE.

Conjugate Priors: Marginal Posterior of β

$$\pi(\beta|y, X) = \int \pi(\beta|\sigma^2, y, X) \pi(\sigma^2|y, X) d\sigma^2$$

$$\pi(\beta|\sigma^2, y, X) \propto \mathcal{N}_{k+1} \left((\mathbf{M} + X^T X)^{-1} \{ X^T X \hat{\beta} + \mathbf{M} \tilde{\beta} \}, \sigma^2 (\mathbf{M} + X^T X)^{-1} \right)$$

$$\pi(\sigma^2|y, X) \propto \mathcal{IG} \left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})}{2} \right)$$

□ To integrate, we use the following transformations:

$$\tau = 1/\sigma^2, \quad d\tau = -\tau^2 d\sigma^2, \quad \hat{\mu} = \{\mathbf{M} + X^T X\}^{-1} [X^T X \hat{\beta} + \mathbf{M} \tilde{\beta}]$$

$$Z = \tau \frac{1}{2} \left[(\beta - \hat{\mu})^T (\mathbf{M} + X^T X) (\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta}) \right]$$

$$d\tau = \frac{2}{G} dZ$$

Inverse-gamma	$\theta \sim \text{Inv-gamma}(\alpha, \beta)$ $p(\theta) = \text{Inv-gamma}(\theta \alpha, \beta)$	shape $\alpha > 0$ scale $\beta > 0$
---------------	---	---

$\tau = 1/\sigma^2, \quad d\tau = -\tau^2 d\sigma^2, \quad \hat{\mu} = \{\mathbf{M} + X^T X\}^{-1} [X^T X \hat{\beta} + \mathbf{M} \tilde{\beta}]$ $Z = \tau \frac{1}{2} \left[(\beta - \hat{\mu})^T (\mathbf{M} + X^T X) (\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta}) \right]$ $d\tau = \frac{2}{G} dZ$

Conjugate Priors: Marginal Posterior of β

$$\tau = 1/\sigma^2, \quad d\tau = -\tau^2 d\sigma^2, \quad \hat{\mu} = (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} [\mathbf{X}^T \mathbf{X} \hat{\beta} + \mathbf{M} \tilde{\beta}]$$

$$Z = \tau \frac{1}{2} [(\beta - \hat{\mu})^T (\mathbf{M} + \mathbf{X}^T \mathbf{X}) (\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\beta} - \hat{\beta})]$$

$$d\tau = \frac{2}{G} dZ$$

□ The marginal now takes the form:

$$\begin{aligned}\pi(\beta | y, X) &\propto \int_0^\infty e^{-Z^2} \tau^{\left(\frac{k+1}{2}\right) + \left(\frac{n}{2} + a + 1\right)} d\sigma^2 \sim \int_0^\infty e^{-Z^2} \tau^{\frac{k}{2} + \frac{n}{2} + a - \frac{1}{2}} d\tau \\ &\sim G^{-\frac{k}{2} - \frac{n}{2} - a - \frac{1}{2}} \int_0^\infty e^{-Z^2} Z^{\frac{k}{2} + \frac{n}{2} + a - \frac{1}{2}} dZ \sim G^{-\frac{k}{2} - \frac{n}{2} - a - \frac{1}{2}}\end{aligned}$$

where

$$G = (\beta - \hat{\mu})^T (\mathbf{M} + \mathbf{X}^T \mathbf{X}) (\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\beta} - \hat{\beta})$$

□ The last integral above is a scalar quantity.

Conjugate Priors: Marginal Posterior of β

- Thus integrating out σ^2 leads to a multivariate Student's T marginal posterior on β :

$$\pi(\beta|y, X) \propto \left[(\beta - \hat{\mu})^T (M + X^T X)(\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (M^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta}) \right]^{-\frac{k}{2} - \frac{n}{2} - a - \frac{1}{2}}$$

where:

$$\hat{\mu} = \{M + X^T X\}^{-1} [X^T X \hat{\beta} + M \tilde{\beta}]$$

- Recall that the density of the multivariate $\mathcal{T}_p(\nu, \theta, \Sigma)$ is:

$$\mathcal{T}_p(\nu, \theta, \Sigma) = \frac{\Gamma((\nu + p)/2)}{\sqrt{\det(\Sigma)\nu\pi}} \left[1 + \frac{(t - \theta)^T \Sigma^{-1} (t - \theta)}{\nu} \right]^{-\frac{\nu + p}{2}}$$

Conjugate Priors: Marginal Posterior of β

- We can now see that our marginal

$$\pi(\beta|y, X) \propto \left[(\beta - \hat{\mu})^T (\mathbf{M} + X^T X)(\beta - \hat{\mu}) + 2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1}(\tilde{\beta} - \hat{\beta}) \right]^{-\frac{k}{2} - \frac{n}{2} - a - \frac{1}{2}}$$

can be written (can check with substitution of the Eqs below) in the form of

$$\mathcal{T}_p(\nu, \theta, \Sigma) = \frac{\Gamma((\nu + p)/2) / \Gamma(\nu/2)}{\sqrt{\det(\Sigma)\nu\pi}} \left[1 + \frac{(\theta - \theta)^T \Sigma^{-1} (\theta - \theta)}{\nu} \right]^{-\frac{\nu+p}{2}}$$

as follows ($p = k + 1, \nu = n + 2a$):

$$\beta|y, X \sim \mathcal{T}_{k+1}(n + 2a, \hat{\mu}, \hat{\Sigma})$$

$$\hat{\mu} = (\mathbf{M} + X^T X)^{-1} ((X^T X)\hat{\beta} + \mathbf{M}\tilde{\beta})$$

$$\hat{\Sigma} = \frac{2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (\mathbf{M}^{-1} + (X^T X)^{-1})^{-1}(\tilde{\beta} - \hat{\beta})}{n + 2a} (\mathbf{M} + X^T X)^{-1}$$



Multivariate Student's- \mathcal{T} Distribution

- Recall the properties of the multivariate Student's- \mathcal{T} distribution:*

Multivariate
Student- t

$$\begin{aligned}\theta &\sim t_\nu(\mu, \Sigma) \\ p(\theta) &= t_\nu(\theta|\mu, \Sigma) \\ (\text{implicit dimension } d)\end{aligned}$$

degrees of freedom $\nu > 0$
location $\mu = (\mu_1, \dots, \mu_d)$
symmetric, pos. definite
 $d \times d$ scale matrix Σ

$$p(\theta) = \frac{\Gamma((\nu+d)/2)}{\Gamma(\nu/2)\nu^{d/2}\pi^{d/2}} |\Sigma|^{-1/2} \times (1 + \frac{1}{\nu}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu))^{-(\nu+d)/2}$$

$$\begin{aligned}\mathbb{E}(\theta) &= \mu, \text{ for } \nu > 1 \\ \text{var}(\theta) &= \frac{\nu}{\nu-2} \Sigma, \text{ for } \nu > 2 \\ \text{mode}(\theta) &= \mu\end{aligned}$$

*Note that in the notation of the tables above the degrees of freedom of the distribution rather than the dimensionality are shown as subscripts, e.g., t_ν vs t_d in the earlier slide. Hopefully the notation will be clear from the discussion.

Bayesian Data Analysis, [A. Gelman, J. Carlin, H. Stern and D. Rubin](#), 2004



Conjugate Priors: Predictive Distribution

- For a given $(m, k + 1)$ explanatory matrix \tilde{X} , the outcome \tilde{y} can be inferred through the predictive distribution*

$$\pi(\tilde{y}|\sigma^2, y, X, \tilde{X})$$

- Since $\pi(\tilde{y}|\tilde{X}, \beta, \sigma^2) \propto \mathcal{N}(\tilde{X}\beta, \sigma^2 I_m)$,

and since the posterior of β conditional on σ^2 is given as

$$\beta|\sigma^2, y, X \sim \mathcal{N}_{k+1} \left((\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{X} \hat{\beta} + \mathbf{M} \tilde{\beta} \}, \sigma^2 (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \right)$$

we can see that:

$$\pi(\tilde{y}|\sigma^2, y, X, \tilde{X}) \propto \int \pi(\tilde{y}|\beta, \sigma^2, y, X, \tilde{X}) \pi(\beta|\sigma^2, y, X) d\beta$$

is a Gaussian.

* We will later on integrate σ^2 out to compute: $\pi(\tilde{y}|y, X, \tilde{X})$



Appendix

- We will make use next of the following conditional expectation equalities:

$$\mathbb{E}[X|Z] = \mathbb{E}[\mathbb{E}[X|Y,Z]|Z]$$

$$Var[X|Z] = Var[\mathbb{E}[X|Y,Z]|Z] + \mathbb{E}[Var[X|Y,Z]|Z]$$

- The proof of these is quite straightforward, e.g.

$$\mathbb{E}[X|Z] = \int_x x p(x|z) dx = \int_x x \int_y p(x,y|z) dy dx$$

$$\int_x x \int_y p(x|y,z) p(y|z) dy dx = \int_y \left(\int_x x p(x|y,z) dx \right) p(y|z) dy = \mathbb{E}[\mathbb{E}[X|Y,Z]|Z]$$

Conjugate Priors: Predictive Distribution

- We can compute the mean and the variance of this predictive distribution as follows:

$$\begin{aligned}\mathbb{E}_\pi[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] &= \mathbb{E}_\pi(\mathbb{E}_\pi(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}})|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) = \\ \mathbb{E}_\pi(\tilde{\mathbf{X}}\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &= \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1}(\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{M} \tilde{\boldsymbol{\beta}})\end{aligned}$$

and

$$\begin{aligned}Var_\pi[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] &= \\ \mathbb{E}_\pi(Var_\pi(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}})|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) + Var_\pi(\mathbb{E}_\pi(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}})|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &= \\ \mathbb{E}_\pi(\sigma^2 \mathbf{I}_m |\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) + Var_\pi(\tilde{\mathbf{X}}\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &= \\ \sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}} \sigma^2 (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T &= \sigma^2 (\mathbf{I}_m + \tilde{\mathbf{X}} (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)\end{aligned}$$

- In conclusion:

$$\begin{aligned}\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}} &\sim \mathcal{N}_m\left(\tilde{\mathbf{X}} \mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}], \sigma^2 (\mathbf{I}_m + \tilde{\mathbf{X}} (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)\right) \\ \mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}] &= (\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \{ \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{M} \tilde{\boldsymbol{\beta}} \}\end{aligned}$$

Conjugate Priors: Predictive Distribution

$$\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}} \sim \mathcal{N}(\tilde{\mathbf{X}}\mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}], \sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}}Var[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]\tilde{\mathbf{X}}^T)$$

$$\mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] = \tilde{\mathbf{X}}\mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}] = \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1}\{\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{M} \tilde{\boldsymbol{\beta}}\},$$

$$Var[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] = \sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}}Var[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]\tilde{\mathbf{X}}^T = \sigma^2 \mathbf{I}_m + \sigma^2 \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1}\tilde{\mathbf{X}}^T$$

□ We now integrate σ^2 against the posterior distribution

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \mathcal{IG}\left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{2}\right)$$

to obtain:

$$\begin{aligned} \pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &\propto \int_0^\infty \mathcal{N}(\tilde{\mathbf{X}}\mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}], \sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}}Var[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]\tilde{\mathbf{X}}^T) \\ &\times \mathcal{IG}\left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{2}\right) d\sigma^2 \end{aligned}$$

Conjugate Priors: Predictive Distribution

$$\begin{aligned}\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &\propto \int_0^\infty \mathcal{N}(\tilde{\mathbf{X}}\mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}], \sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}}Var[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]\tilde{\mathbf{X}}^T) \\ &\times \mathcal{IG}\left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{2}\right) d\sigma^2\end{aligned}$$

□ We substitute:

$$\sigma^2 \mathbf{I}_m + \tilde{\mathbf{X}}Var[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}]\tilde{\mathbf{X}}^T = \sigma^2 (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)$$

$$\begin{aligned}\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &\propto \int_0^\infty \mathcal{N}(\tilde{\mathbf{X}}\mathbb{E}[\boldsymbol{\beta}|\sigma^2, \mathbf{y}, \mathbf{X}], \sigma^2 (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)) \\ &\times \mathcal{IG}\left(\frac{n}{2} + a, b + \frac{s^2}{2} + \frac{(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})}{2}\right) d\sigma^2 \\ &\propto \int_0^\infty \sigma^{-m-n-2a-2} \exp\left(-\frac{1}{2\sigma^2} \{2b + s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \right. \\ &\quad \left. (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])^T (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)^{-1} (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])\}\right) d\sigma^2\end{aligned}$$

Conjugate Priors: Predictive Distribution

$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto$$

$$\propto \int_0^\infty \sigma^{-m-n-2a-2} \exp\left(-\frac{1}{2\sigma^2} \{2b + s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])^T (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)^{-1} (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}]) \} \right) d\sigma^2$$

□ If we call the term inside $\{.\}$ as G , then:

$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto \int_0^\infty \sigma^{-m-n-2a-2} e^{-\frac{1}{2\sigma^2}G} d\sigma^2 \Rightarrow$$
$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto \int_0^\infty (\sigma^2)^{-\frac{m+n+2a+2}{2}} e^{-\frac{1}{2\sigma^2}G} d\sigma^2 \Rightarrow \pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto G^{-\frac{m+n+2a+2}{2}} \int_0^\infty Z^{\frac{m+n+2a+2}{2}} e^{-Z} \left(-\frac{1}{2} G \frac{1}{Z^2} dZ \right) \sim G^{-\frac{m+n+2a+2}{2}}$$
$$Z = \frac{1}{2\sigma^2} G$$

$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto \{2b + s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])^T (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)^{-1} (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}]) \}^{-(m+n+2a)/2}$$



Conjugate Priors: Predictive Distribution

$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto \{2b + s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \\ (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])^T (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)^{-1} (\tilde{\mathbf{y}} - \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}])\}^{-(m+n+2a)/2}$$

□ This corresponds to a Student's \mathcal{T} distribution:

$$\mathcal{T}_p(\nu, \theta, \Sigma) = \frac{\Gamma((\nu + p)/2)}{\sqrt{\det(\Sigma)\nu\pi}} \left[1 + \frac{(\mathbf{t} - \theta)^T \Sigma^{-1} (\mathbf{t} - \theta)}{\nu} \right]^{-\frac{\nu+p}{2}}$$

with

$$\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}} \sim \mathcal{T}_m(n + 2a, \hat{\boldsymbol{\mu}}, \hat{\Sigma})$$

$$\hat{\boldsymbol{\mu}} = \mathbb{E}[\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] = \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{M} \tilde{\boldsymbol{\beta}})$$

$$\hat{\Sigma} = \frac{\{2b + s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T (\mathbf{M}^{-1} + (\mathbf{X}^T \mathbf{X})^{-1})^{-1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})\} (\mathbf{I}_m + \tilde{\mathbf{X}}(\mathbf{M} + \mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T)}{n + 2a}$$



Implementation of the Conjugate Prior

Our Prior:

$$\beta | \sigma^2, X \sim \mathcal{N}_{k+1}(\tilde{\beta}, \sigma^2 M^{-1}), M \text{ a } (k+1) \times (k+1) \text{ pos. def. symm. matrix}$$
$$\sigma^2 | X \sim IG(a, b), a, b > 0$$

- We assume that there is no precise information available about $\tilde{\beta}, M, a, b$.
- Let us choose as $M = I_{k+1}/c$ and $\tilde{\beta} = \mathbf{0}_{k+1}$.

$$\beta | \sigma^2, X \sim \mathcal{N}_{k+1}(0_{k+1}, c\sigma^2 I_{k+1})$$

- Let us take $a = 2.1$ and $b = 2$, i.e. prior mean and prior variance of σ^2 equal to 1.82 and 33.06.

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \theta > 0$$

$$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$$

$$\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$$

$$\text{mode}(\theta) = \frac{\beta}{\alpha+1}$$

- The intuition is that if c is large, the prior on β should be more diffuse and have less bearing on the outcome.
- However, this turns out not to be the case. There is lasting influence of c on the posterior means of σ^2 and β_0 .



Influence of the Prior Scale c on Bayesian Estimates of β

***** Conjugate Priors *****

$\mathbb{E}_\pi(\sigma^2 y, X)$	$\mathbb{E}_\pi(\beta_0 y, X)$	$V_\pi(\beta_0 y, X)$	
c	$\mathbb{E}[\text{sigma}^2 y, X]$	$\mathbb{E}[\text{beta}_0 y, X]$	$V[\text{beta}_0 y, X]$
<hr/>			
0.1	1.0044	0.1251	0.0988
1.0	0.8541	0.9031	0.7733
10.0	0.6976	4.7299	3.8991
100.0	0.5746	9.6626	6.8355
1000.0	0.5470	10.8476	7.3419

An implementation is available [MatLab](#), [C++](#)

$$\mathbb{E}^\pi[\sigma^2 | y, X] = \frac{2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (M^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})}{n + 2a - 2}$$

$$\mathbb{E}^\pi[\beta | y, X] = (M + X^T X)^{-1} ((X^T X) \hat{\beta} + M \tilde{\beta})$$

$$\text{Var}[\beta | y, X] = \frac{2b + s^2 + (\tilde{\beta} - \hat{\beta})^T (M^{-1} + (X^T X)^{-1})^{-1} (\tilde{\beta} - \hat{\beta})}{n + 2a} (M + X^T X)^{-1}$$

Bayesian Estimates of β for $c = 100$

***** Conjugate Priors *****

Bayes estimates of beta for $c = 100$

	$\mathbb{E}_\pi(\beta_i y, X)$	$V_\pi(\beta_i y, X)$
beta_i	$\mathbb{E}[\text{beta_i} y, X]$	$V[\text{beta_i} y, X]$
-----	-----	-----
beta_0	9.6626	6.8355
beta_1	-0.0040	0.0000
beta_2	-0.0516	0.0004
beta_3	0.0418	0.0077
beta_4	-1.2633	0.2615
beta_5	0.2307	0.0090
beta_6	-0.0832	1.9310
beta_7	-0.1917	0.8254
beta_8	0.1608	0.0462
beta_9	-1.2069	0.6127
beta_10	-0.2567	0.4267

An implementation
is available
[MatLab](#), [C++](#)



Influence of the Conjugate Prior

- The value of c thus has a significant influence on the estimators and even more on the posterior variance.
- The Bayes estimates stabilize for very large values of c .

Thus the prior associated with a particular c should not be considered as a weak or pseudo-noninformative prior but, on the opposite, associated with a specific proper prior information.

- The dependence on (a, b) is equally strong.

Considering these limitations of conjugate priors on at least the posterior variance, a more sophisticated noninformative strategy is needed.

- We first look a middle-ground perspective which settles the problem of the choice of M .



Zellner's Informative G-Prior

- We start with a middle ground solution by introducing only information about the location parameter of the regression but bypassing the selection of the prior correlation structure.

$$\boldsymbol{\beta} | \sigma^2, \mathbf{X} \sim \mathcal{N}_{k+1}(\tilde{\boldsymbol{\beta}}, c\sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$$

$\sigma^2 \sim \pi(\sigma^2 | \mathbf{X}) \propto \sigma^{-2}$ improper Jeffreys prior

- Zellner, A. (1971). [An Introduction to Bayesian Econometrics](#). John Wiley, New York.
- Zellner, A. (1984). [Basic Issues in Econometrics](#). University of Chicago Press, Chicago.
- Carlin, B. and Louis, T. (1996). [Bayes and Empirical Bayes Methods for Data Analysis](#). Chapman and Hall, New York.



Zellner's Informative G-Prior

- We start with a middle ground solution by introducing only information about the location parameter of the regression but bypassing the selection of the prior correlation structure.

$$\beta | \sigma^2, \mathbf{X} \sim \mathcal{N}_{k+1}(\tilde{\beta}, c\sigma^2(\mathbf{X}^T \mathbf{X})^{-1})$$

$\sigma^2 \sim \pi(\sigma^2 | \mathbf{X}) \propto \sigma^{-2}$ *improper Jeffreys prior*

- The prior determination is restricted to the choices of $\tilde{\beta}$ and constant c . c can be interpreted as a measure of the information available in the prior relative to the sample.
 - E.g., we will see that setting $1/c = 0.5$ gives the prior the same weight as 50% of the sample.
- There is still strong influence of c .



Zellner's Informative G-Prior: Conditional Posterior of β

- The joint posterior now takes the form (note $X^T X$ is used in both likelihood and prior):

$$\begin{aligned}\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto (\sigma^2)^{-(n/2)} \exp\left(-\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) - \frac{1}{2\sigma^2} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})\right) (\sigma^2)^{-(k+1)/2} \\ &\quad \times \exp\left(-\frac{1}{2c\sigma^2} (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})\right) (\sigma^2)^{-1}\end{aligned}$$

- We can now compute the following:

$$\begin{aligned}\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} &\sim \exp\left(-\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}) - \frac{1}{2c\sigma^2} (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \tilde{\boldsymbol{\beta}})\right) \\ &\sim \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right)^T \frac{c+1}{c\sigma^2} \mathbf{X}^T \mathbf{X} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right) \right) \Rightarrow\end{aligned}$$

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim \mathcal{N}_{k+1} \left(\frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right), \frac{c\sigma^2}{c+1} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

Zellner's Informative G-Prior: Posterior Marginal of σ^2

□ Starting again with the posterior

$$\begin{aligned}\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &\propto (\sigma^2)^{-(n+k+3/2)} \exp\left(-\frac{1}{2\sigma^2} s^2 - \frac{1}{2\sigma^2} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})\right) \\ &\quad \times \exp\left(-\frac{1}{2c\sigma^2} (\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \widetilde{\boldsymbol{\beta}})\right) = \\ &\sim (\sigma^2)^{-(k+1)/2} \exp\left(-\frac{1}{2} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\widetilde{\boldsymbol{\beta}}}{c} + \widehat{\boldsymbol{\beta}}\right)\right)^T \frac{c+1}{c\sigma^2} \mathbf{X}^T \mathbf{X} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\widetilde{\boldsymbol{\beta}}}{c} + \widehat{\boldsymbol{\beta}}\right)\right)\right) \\ &\quad \times (\sigma^2)^{-(n/2+1)} \exp\left(-\frac{1}{2\sigma^2} \left(s^2 + \frac{1}{c+1} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})\right)\right)\end{aligned}$$

□ We can now compute the following posterior marginal with integration in $\boldsymbol{\beta}$:

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \text{InvGamma}\left(\frac{n}{2}, \frac{s^2}{2} + \frac{1}{2(c+1)} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})\right)$$

$$\begin{aligned}p(\theta) &= \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0 \\ E(\theta) &= \frac{\beta}{\alpha-1}, \quad \text{for } \alpha > 1 \\ \text{var}(\theta) &= \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \quad \alpha > 2 \\ \text{mode}(\theta) &= \frac{\beta}{\alpha+1}\end{aligned}$$

$$\mathbb{E}[\sigma^2 | \mathbf{y}, \mathbf{X}] = \frac{s^2 + \frac{1}{(c+1)} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\widetilde{\boldsymbol{\beta}} - \widehat{\boldsymbol{\beta}})}{n-2}$$

Zellner's Informative G-Prior: Posterior Marginal of β

□ Starting again with the posterior

$$\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-(n+k+3/2)} \exp\left\{-\frac{1}{2\sigma^2} \left[\left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right)^T \frac{c+1}{c} \mathbf{X}^T \mathbf{X} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right) + s^2 + \frac{1}{c+1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) \right] \right\}$$

□ We can now compute the following posterior marginal with integration in σ^2 :

$$\boldsymbol{\beta} | \mathbf{y}, \mathbf{X} \sim \left\{ \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right)^T \frac{c+1}{c} \mathbf{X}^T \mathbf{X} \left(\boldsymbol{\beta} - \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right) \right) + s^2 + \frac{1}{c+1} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) \hat{\boldsymbol{\beta}} \right\}^{-(n+k+1/2)}$$

$$\boldsymbol{\beta} | \mathbf{y}, \mathbf{X} \sim \mathcal{T}_{k+1} \left(n, \frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right), \frac{c \left(s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) / (c+1) \right)}{n(c+1)} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

Multivariate Student-t	$\theta \sim t_\nu(\mu, \Sigma)$ $p(\theta) = t_\nu(\theta \mu, \Sigma)$ (implicit dimension d)	degrees of freedom $\nu > 0$ location $\mu = (\mu_1, \dots, \mu_d)$ symmetric, pos. definite $d \times d$ scale matrix Σ	$p(\theta) = \frac{\Gamma((\nu+d)/2)}{\Gamma(\nu/2)\nu^{d/2}\pi^{d/2}} \Sigma ^{-1/2} \times (1 + \frac{1}{\nu}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu))^{-(\nu+d)/2}$	$E(\theta) = \mu$, for $\nu > 1$ $\text{var}(\theta) = \frac{\nu}{\nu-2} \Sigma$, for $\nu > 2$ mode(θ) = μ
------------------------	--	---	---	--



Zellner's Informative G-Prior: Bayes Estimates

- The Bayes estimate of β can be derived from:

$$\beta | y, X \sim \mathcal{T}_{k+1} \left(n, \frac{c}{c+1} \left(\frac{\tilde{\beta}}{c} + \hat{\beta} \right), \frac{c \left(s^2 + (\tilde{\beta} - \hat{\beta})^T X^T X (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n(c+1)} (X^T X)^{-1} \right)$$

$$\mathbb{E}[\beta | y, X] = \frac{1}{c+1} (\tilde{\beta} + c\hat{\beta})$$

- The posterior variance of β is also given from above as:

$$V_\pi[\beta | y, X] = \frac{c}{c+1} \frac{\left(s^2 + (\tilde{\beta} - \hat{\beta})^T X^T X (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n-2} (X^T X)^{-1}$$

$$p(\theta) = \frac{\Gamma((\nu+d)/2)}{\Gamma(\nu/2)\nu^{d/2}\pi^{d/2}} |\Sigma|^{-1/2} \times (1 + \frac{1}{\nu}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu))^{-(\nu+d)/2}$$
$$\begin{aligned} E(\theta) &= \mu, \text{ for } \nu > 1 \\ \text{var}(\theta) &= \frac{\nu}{\nu-2} \Sigma, \text{ for } \nu > 2 \\ \text{mode}(\theta) &= \mu \end{aligned}$$



Zellner's Informative G-Prior: Bayes Estimates

$$\mathbb{E}[\beta | y, X] = \frac{1}{c+1} (\tilde{\beta} + c\hat{\beta})$$

- If $c = 1$, you can see from the above equation that it is like putting the same weight on the prior information and on the sample:

$$\mathbb{E}[\beta | y, X] = \frac{1}{2} (\tilde{\beta} + \hat{\beta})$$

which is the average of the prior mean and the MLE estimator.

- If $c = 100$, the prior weights 1% of the sample.



Zellner's Informative G-Prior: Bayes Estimates

- The Bayes estimate of σ^2 (use the mean of InvGamma) can be derived from

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \mathcal{IG} \left(\frac{n}{2}, \frac{s^2}{2} + \frac{1}{2(c+1)} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) \right) \Rightarrow$$

as

$$\mathbb{E}[\sigma^2 | \mathbf{y}, \mathbf{X}] = \frac{s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) / (c+1)}{n-2}$$

- Note that only when c goes to infinity the influence of the prior vanishes.

$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$	$E(\theta) = \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1$
	$\text{var}(\theta) = \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2$
	$\text{mode}(\theta) = \frac{\beta}{\alpha+1}$



Zellner's Informative G-Prior: Marginal Posterior Mean and Variance of β

***** Zellner's G-Prior *****

$$\mathbb{E}[\boldsymbol{\beta}|y, X] = \frac{1}{c+1}(\tilde{\boldsymbol{\beta}} + c\hat{\boldsymbol{\beta}})$$

Posterior mean and variance of beta for $c = 100$

$$\mathbb{E}_\pi(\beta_i|y, X)$$

$$V_\pi(\beta_i|y, X)$$

$$\log_{10}(BF)$$

$$\frac{c}{c+1} \frac{\left(s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T X^T X (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})/(c+1)\right)}{n-2} (X^T X)^{-1}$$

beta_i	$\mathbb{E}[\text{beta_i} y, X]$	$V[\text{beta_i} y, X]$	$\log_{10}(BF)$
(Intercept)	10.8895	6.8229	2.1873 (****)
X1	-0.0044	0.0000	1.1571 (***)
X2	-0.0533	0.0003	0.6667 (**)
X3	0.0673	0.0072	-0.8585
X4	-1.2808	0.2316	0.4726 (*)
X5	0.2293	0.0079	0.3861 (*)
X6	-0.3533	1.7877	-0.9860
X7	-0.2351	0.7373	-0.9849
X8	0.1793	0.0408	-0.8225
X9	-1.2726	0.5449	-0.3461
X10	-0.4288	0.3934	-0.8949

[See here for BF calculation](#)

(Intercept) 10.8895

X1 -0.0044

X2 -0.0533

X3 0.0673

X4 -1.2808

X5 0.2293

X6 -0.3533

X7 -0.2351

X8 0.1793

X9 -1.2726

X10 -0.4288

6.8229

0.0000

0.0003

0.0072

0.2316

0.0079

1.7877

0.7373

0.0408

0.5449

0.3934

An implementation
is available

[MatLab](#), [C++](#)

Evidence against H_0 :

(****) decisive

(***) strong

(**) substantial

(*) poor



Zellner's Informative G-Prior: Marginal Posterior Mean and Variance of β

***** Zellner's G-Prior *****

$$\mathbb{E}[\beta|y, X] = \frac{1}{c+1}(\tilde{\beta} + c\hat{\beta})$$

Posterior mean and variance of beta for $c = 1000$

$\mathbb{E}_\pi(\beta_i y, X)$	$V_\pi(\beta_i y, X)$	$\log_{10}(BF)$	$V_\pi[\beta y, X] = \frac{c}{c+1} \frac{(s^2 + (\tilde{\beta} - \hat{\beta})^T X^T X (\tilde{\beta} - \hat{\beta})/(c+1))}{n-2} (X^T X)^{-1}$
--------------------------------	-----------------------	-----------------	--

beta_i	$\mathbb{E}[\text{beta_i} y, X]$	$V[\text{beta_i} y, X]$	$\log10(BF)$
--------	-----------------------------------	--------------------------	--------------

[See here for BF calculation](#)

(Intercept)	10.9874	6.6644	1.7973 (***)
X1	-0.0044	0.0000	0.7375 (**)
X2	-0.0538	0.0003	0.2313 (*)
X3	0.0679	0.0070	-1.3506
X4	-1.2923	0.2262	0.0307
X5	0.2314	0.0078	-0.0588
X6	-0.3564	1.7461	-1.4834
X7	-0.2372	0.7202	-1.4822
X8	0.1809	0.0399	-1.3130
X9	-1.2840	0.5323	-0.8177
X10	-0.4327	0.3843	-1.3885

An implementation
is available

[MatLab](#), [C++](#)

Evidence against H_0 :

(****) decisive

(***) strong

(**) substantial

(*) poor



Zellner's Informative G-Prior: Predictive Modeling

- We want to predict ($m \geq 1$) future observations when the explanatory variables \tilde{X} but not the outcome variables

$$\tilde{\mathbf{y}} \sim \mathcal{N}_m(\tilde{\mathbf{X}}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_m)$$

have been observed.

- Predictive distribution on $\tilde{\mathbf{y}}$ defined as marginal of the joint posterior distribution on $(\tilde{\mathbf{y}}, \boldsymbol{\beta}, \sigma^2)$. Can be computed analytically by

$$\pi(\tilde{\mathbf{y}}|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \propto \int \pi(\tilde{\mathbf{y}}|\sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) \pi(\sigma^2|\mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) d\sigma^2$$

Zellner's Informative G-Prior: Predictive Modeling

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim \mathcal{N}_{k+1} \left(\frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right), \frac{c\sigma^2}{c+1} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- Conditional on σ^2 the future vector of observations has a Gaussian distribution with

$$\begin{aligned}\mathbb{E}_\pi(\tilde{\mathbf{y}} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &= \mathbb{E}_\pi[\mathbb{E}_\pi(\tilde{\mathbf{y}} | \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] \\ &= \mathbb{E}_\pi[\tilde{\mathbf{X}} \boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] = \tilde{\mathbf{X}} \mathbb{E}_\pi[\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] \Rightarrow\end{aligned}$$

$$\mathbb{E}_\pi(\tilde{\mathbf{y}} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) = \tilde{\mathbf{X}} \frac{\tilde{\boldsymbol{\beta}} + c\hat{\boldsymbol{\beta}}}{c+1} \text{ (independent of } \sigma^2\text{)}$$

- This representation is quite intuitive, being the product of the matrix of explanatory variables $\tilde{\mathbf{X}}$ by the Bayes estimate of $\boldsymbol{\beta}$.

Zellner's Informative G-Prior: Predictive Modeling

$$\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X} \sim \mathcal{N}_{k+1} \left(\frac{c}{c+1} \left(\frac{\tilde{\boldsymbol{\beta}}}{c} + \hat{\boldsymbol{\beta}} \right), \frac{c\sigma^2}{c+1} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

- Similarly, we can compute:

$$\begin{aligned} V_\pi(\tilde{\mathbf{y}} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) &= \mathbb{E}_\pi [V(\tilde{\mathbf{y}} | \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] \\ &\quad + V_\pi [\mathbb{E}_\pi (\tilde{\mathbf{y}} | \boldsymbol{\beta}, \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] = \\ &= \mathbb{E}_\pi [\sigma^2 \mathbf{I}_m | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] + V_\pi [\tilde{\mathbf{X}} \boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}] \end{aligned}$$

$$V_\pi(\tilde{\mathbf{y}} | \sigma^2, \mathbf{y}, \mathbf{X}, \tilde{\mathbf{X}}) = \sigma^2 \left(\mathbf{I}_m + \frac{c}{c+1} \tilde{\mathbf{X}} (\mathbf{X}^T \mathbf{X})^{-1} \tilde{\mathbf{X}}^T \right)$$

- Here, we are interested on the highest posterior density (HPD) regions on subvectors of the parameter β derived from the marginal posterior distribution of β .
- For a single parameter,

$$\beta_i | \mathbf{y}, \mathbf{X} \sim \mathcal{T}_1 \left(n, \frac{c}{c+1} \left(\frac{\tilde{\beta}_i}{c} + \hat{\beta}_i \right), \frac{c \left(s^2 + (\tilde{\beta} - \hat{\beta})^T \mathbf{X}^T \mathbf{X} (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n(c+1)} \omega_{ii} \right)$$

where ω_{ii} is the (i, i) element of $(\mathbf{X}^T \mathbf{X})^{-1}$

- Let us define

$$\tau = \frac{\tilde{\beta} + c\hat{\beta}}{c + 1}$$

- Also

$$\mathbf{K} = \frac{c \left(s^2 + (\tilde{\beta} - \hat{\beta})^T X^T X (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n(c+1)} (\mathbf{X}^T \mathbf{X})^{-1} = (K_{ij})$$

- The variable

$$\zeta_i = \frac{\beta_i - \tau_i}{\sqrt{K_{ii}}}$$

has a \mathcal{T} -distribution with n degrees of freedom.

- A $1 - \alpha$ HPD interval on β_i is thus given by ([see also here](#))

$$\left[\tau_i - \sqrt{K_{ii}} F_n^{-1}(1 - \alpha/2), \tau_i + \sqrt{K_{ii}} F_n^{-1}(1 - \alpha/2) \right]$$

where F_n is the CDF of $\mathcal{T}(\nu = n)$.

- Note that these HPD are different from the frequentist confidence intervals defined earlier as:

$$\beta_i | \mathbf{y}, \mathbf{X} \sim \mathcal{T}(n - k - 1, \hat{\beta}_i, \omega_{(i,i)} s^2 / (n - k - 1))$$

$$\frac{\beta_i - \hat{\beta}_i}{\omega_{(i,i)} s^2 / (n - k - 1)} \sim \mathcal{T}(\nu = n - k - 1, 0, 1), \text{ where: } \omega_{(i,i)} = (\mathbf{X}^T \mathbf{X})^{-1}|_{(i,i)}$$

$$\left\{ \beta_i : |\beta_i - \hat{\beta}_i| \leq F_{n-k-1}^{-1}(1 - \alpha/2) \sqrt{\omega_{(i,i)} s^2 / (n - k - 1)} \right\}$$

$$s^2 \triangleq (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^T (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

beta_i	HPD Interval
beta_0	[4.6518, 17.3450]
beta_1	[-0.0077, -0.0012]
beta_2	[-0.0992, -0.0084]
beta_3	[-0.1384, 0.2742]
beta_4	[-2.4629, -0.1244]
beta_5	[0.0152, 0.4481]
beta_6	[-3.6054, 2.8918]
beta_7	[-2.3238, 1.8489]
beta_8	[-0.3099, 0.6720]
beta_9	[-3.0789, 0.5083]
beta_10	[-1.9571, 1.0909]

$c = 100$

An implementation
is available
[MatLab](#), [C++](#)

$$\left\{ \beta_i : |\beta_i - \hat{\beta}_i| \leq F_{n-k-1}^{-1}(1 - a/2) \sqrt{\omega_{(i,i)} s^2 / (n - k - 1)} \right\}$$

Zellner's Informative G-Prior: 90% High Posterior Density (HPD) Intervals for β_i 's

beta_i	HPD Interval	$\left\{ \beta_i : \beta_i - \hat{\beta}_i \right.$
beta_0	[5.7435, 16.2533]	
beta_1	[-0.0071, -0.0018]	
beta_2	[-0.0914, -0.0162]	
beta_3	[-0.1029, 0.2387]	
beta_4	[-2.2618, -0.3255]	
beta_5	[0.0524, 0.4109]	
beta_6	[-3.0466, 2.3330]	
beta_7	[-1.9649, 1.4900]	
beta_8	[-0.2254, 0.5875]	
beta_9	[-2.7704, 0.1998]	
beta_10	[-1.6950, 0.8288]	

$c = 100$

An implementation
is available
[MatLab](#), [C++](#)

Note: The results given in ``C. P. Robert, [The Bayesian Core](#), Springer, 2nd edition, [chapter 3](#)'' refer to 90% HPD intervals rather than 95% as posted. These results agree with what is shown above.



Zellner's Informative G-Prior: Marginal Distribution of y

- The marginal distribution of y (evidence) is a multivariate \mathcal{T} .
- Since $\beta | \sigma^2, X \sim \mathcal{N}_{k+1}(\tilde{\beta}, c\sigma^2(X^T X)^{-1})$, the linear transform of β satisfies:

$$X\beta | \sigma^2, X \sim \mathcal{N}(X\tilde{\beta}, c\sigma^2 X(X^T X)^{-1} X^T)$$

which implies

$$y | \sigma^2, X \sim \mathcal{N}_n \left(X\tilde{\beta}, \sigma^2(I_n + cX(X^T X)^{-1} X^T) \right)$$

- Integration in σ^2 with $\pi(\sigma^2) = 1/\sigma^2$ gives ([see Appendices](#)):

$$f(y|X, c) = \int_0^\infty \frac{1}{(2\pi)^{n/2}} \frac{1}{(c+1)^{(k+1)/2}} (\sigma^2)^{-n/2-1} e^{-\frac{1}{2\sigma^2}(y-X\tilde{\beta})^T (I_n + cX(X^T X)^{-1} X^T)^{-1} (y-X\tilde{\beta})} d\sigma^2$$

$$f(y|X, c) = (c+1)^{-(k+1)/2} \pi^{-n/2} \Gamma\left(\frac{n}{2}\right) \left[y^T y - \frac{c}{c+1} y^T X(X^T X)^{-1} X^T y + \frac{1}{c+1} \tilde{\beta}^T X^T X \tilde{\beta} - \frac{2}{c+1} y^T X \tilde{\beta} \right]^{-n/2}$$

Appendix

□ It can be easily shown that

$$\left(\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \right)^{-1} = \mathbf{I}_n - \frac{c}{c+1}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$$

□ Indeed:

$$\begin{aligned} & \left(\mathbf{I}_n - \frac{c}{c+1}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \right) \left(\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \right) = \mathbf{I}_n - \frac{c}{c+1}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T + c\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ & - \frac{c^2}{c+1}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T = \mathbf{I}_n + \left(c - \frac{c}{c+1} - \frac{c^2}{c+1} \right) \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T = \mathbf{I}_n \end{aligned}$$

□ This can simplify the term in the exponential in the earlier slide:

$$\begin{aligned} & (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^T \left(\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \right)^{-1} (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) = (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})^T \left(\mathbf{I}_n - \frac{c}{c+1}\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \right) (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}) = \\ & \mathbf{y}^T\mathbf{y} - \frac{c}{c+1}\mathbf{y}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} - \mathbf{y}^T\mathbf{X}\tilde{\boldsymbol{\beta}} + \frac{c}{c+1}\mathbf{y}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\tilde{\boldsymbol{\beta}} \\ & - \tilde{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{y} + \tilde{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}\tilde{\boldsymbol{\beta}} + \frac{c}{c+1}\tilde{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} - \frac{c}{c+1}\tilde{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\tilde{\boldsymbol{\beta}} = \\ & = \mathbf{y}^T\mathbf{y} - \frac{c}{c+1}\mathbf{y}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} + \frac{1}{c+1}\tilde{\boldsymbol{\beta}}^T\mathbf{X}^T\mathbf{X}\tilde{\boldsymbol{\beta}} - \frac{2}{c+1}\mathbf{y}^T\mathbf{X}\tilde{\boldsymbol{\beta}} \end{aligned}$$

Appendix

□ Let us revisit the matrix $\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

□ Note that for any $\boldsymbol{\beta} \in \mathbb{R}^{k+1}$:

$$\left[\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right] \mathbf{X} \boldsymbol{\beta} = \mathbf{X} \boldsymbol{\beta} + c\mathbf{X} \boldsymbol{\beta} = (1+c)\mathbf{X} \boldsymbol{\beta}$$

which implies that $\mathbf{X} \boldsymbol{\beta}$ is an eigenvector with eigenvalue $(1+c)$. There are $(k+1)$ of those.

□ Finally note that for any z in the null space of \mathbf{X}^T , $\mathbf{X}^T z = 0$,

$$\left[\mathbf{I}_n + c\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \right] z = z + c\mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T z = z$$

which implies that these z ($n-k-1$ in number) are eigenvectors of our matrix with 1 as the eigenvalues.

□ The determinant of the matrix is then $(1+c)^{k+1} 1^{n-k-1} = (1+c)^{k+1}$.
This explains the derivation in the earlier slide.



Zellner's Informative G-Prior: Point Null Hypothesis

- If a null hypothesis is $H_0 : R\beta = r$, R being a $q \times (k+1)$ matrix, the model under H_0 can be rewritten as

$$y | \beta^0, \sigma^2, X_0 \stackrel{H_0}{\sim} \mathcal{N}_n(X_0 \beta^0, \sigma^2 I_n)$$

where β^0 is $(k + 1 - q)$ dimensional.

- Under the prior

$$\beta^0 | \sigma^2, X_0 \sim \mathcal{N}_{k+1-q}(\tilde{\beta}^0, c_0 \sigma^2 (X_0^T X_0)^{-1})$$

- The marginal distribution of y under H_0 is:

$$f(y|X_0, H_0) = (c_0 + 1)^{-(k+1-q)/2} \pi^{-n/2} \Gamma\left(\frac{n}{2}\right)$$

$$\times \left[y^T y - \frac{c_0}{c_0 + 1} y^T X_0 (X_0^T X_0)^{-1} X_0^T y + \frac{1}{c_0 + 1} \tilde{\beta}_0^T X_0^T X_0 \tilde{\beta}_0 - \frac{2}{c_0 + 1} y^T X_0 \tilde{\beta}_0 \right]^{-n/2}$$



Zellner's Informative G-Prior: Bayes Factor

- The Bayes factor is then given in analytical form as:

$$B_{10}^{\pi} = \frac{f(\mathbf{y}|\mathbf{X})}{f(\mathbf{y}|\mathbf{X}_0, H_0)} = \frac{(c_0 + 1)^{(k+1-q)/2}}{(c + 1)^{(k+1)/2}} \times \\ \left[\frac{\mathbf{y}^T \mathbf{y} - \frac{c_0}{c_0 + 1} \mathbf{y}^T \mathbf{X}_0 (\mathbf{X}_0^T \mathbf{X}_0)^{-1} \mathbf{X}_0^T \mathbf{y} + \frac{1}{c_0 + 1} \tilde{\boldsymbol{\beta}}_0^T \mathbf{X}_0^T \mathbf{X}_0 \tilde{\boldsymbol{\beta}}_0 - \frac{2}{c_0 + 1} \mathbf{y}^T \mathbf{X}_0 \tilde{\boldsymbol{\beta}}_0}{\mathbf{y}^T \mathbf{y} - \frac{c}{c + 1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\beta}} - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X} \tilde{\boldsymbol{\beta}}} \right]^{n/2}$$

- Note that we use the same σ^2 in both models
- The Bayes factor depends on c_0 and c .
- This calculation can be used to evaluate the inclusion or not of each of the terms $\beta_i, i = 1, \dots, k + 1$ in our regression model.

Note: $\mathbf{X}_0 = \mathbf{X}(:, setdiff(1:k + 1, j))$

Noninformative Prior Analysis: Jeffreys' Prior

- Considering the robustness issues of the two priors examined earlier in the case of a complete lack of prior information, we consider now the non-informative Jeffreys' prior.
- The Jeffreys' prior is a flat prior on $(\beta, \log \sigma^2)$

$$\pi^J(\beta, \sigma^2 | X) \propto \sigma^{-2}$$

- The posterior is then given as:

$$\begin{aligned}\pi^J(\beta, \sigma^2 | y, X) &\propto (\sigma^{-2})^{n/2} \exp\left(-\frac{1}{2\sigma^2} (y - X\hat{\beta})^T (y - X\hat{\beta}) - \frac{1}{2\sigma^2} (\beta - \hat{\beta})^T X^T X (\beta - \hat{\beta})\right) \times \sigma^{-2} \\ &= (\sigma^{-2})^{(k+1)/2} \exp\left(-\frac{1}{2\sigma^2} (\beta - \hat{\beta})^T X^T X (\beta - \hat{\beta})\right) \times \\ &\quad (\sigma^{-2})^{(n-k-1)/2+1} \exp\left(-\frac{1}{2\sigma^2} s^2\right)\end{aligned}$$

Noninformative Prior Analysis: Jeffreys' Prior

$$\pi^J(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^{-2})^{(k+1)/2} \exp\left(-\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})\right) \times \\ (\sigma^{-2})^{(n-k-1)/2+1} \exp\left(-\frac{1}{2\sigma^2} s^2\right)$$

- From this joint posterior, we can immediately evaluate the following conditional and marginal posteriors:

$$\pi^J(\sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^2)^{-(n-k-1)/2-1} \exp\left(-\frac{1}{2\sigma^2} s^2\right) = \mathcal{IG}\left(\frac{(n-k-1)}{2}, \frac{s^2}{2}\right)$$

$$\pi^J(\boldsymbol{\beta} | \sigma^2, \mathbf{y}, \mathbf{X}) \propto (\sigma^{-2})^{(k+1)/2} \exp\left(-\frac{1}{2\sigma^2} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}})\right) \\ = \mathcal{N}_{k+1}(\widehat{\boldsymbol{\beta}}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$$

- From the 1st of these eqs., the Bayes' estimate of σ^2 is:

$$\mathbb{E}^\pi(\sigma^2 | \mathbf{y}, \mathbf{X}) = \frac{s^2}{n-k-3}$$

$$p(\theta) = \frac{\beta^\alpha}{\Gamma(\alpha)} \theta^{-(\alpha+1)} e^{-\beta/\theta}, \quad \theta > 0$$

$$\begin{aligned} E(\theta) &= \frac{\beta}{\alpha-1}, \text{ for } \alpha > 1 \\ \text{var}(\theta) &= \frac{\beta^2}{(\alpha-1)^2(\alpha-2)}, \alpha > 2 \\ \text{mode}(\theta) &= \frac{\beta}{\alpha+1} \end{aligned}$$

- This estimate is larger (more pessimistic) than earlier estimates



Noninformative Prior Analysis: Jeffreys' Prior

$$\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto (\sigma^{-2})^{n/2+1} \exp\left(-\frac{s^2 + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{2\sigma^2}\right)$$

- To compute the marginal posterior of $\boldsymbol{\beta}$, we need to integrate in σ^2 . Using symbolic integrator:

$$\begin{aligned}\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) &\propto \int (\sigma^2)^{-n/2-1} \exp\left(-\frac{s^2 + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{2\sigma^2}\right) d\sigma^2 \\ &\propto \left(s^2 + (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})\right)^{-n/2} = \left(1 + \frac{(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})^T \left(\frac{s^2 (\mathbf{X}^T \mathbf{X})^{-1}}{n-k-1}\right)^{-1} (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}{n-k-1}\right)^{-n/2}\end{aligned}$$

- This is a Students' \mathcal{T} distribution:

$$\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) \propto \mathcal{T}_{k+1}(\nu = n - k - 1, \boldsymbol{\theta}, \boldsymbol{\Sigma}), \text{ with: } \boldsymbol{\theta} = \hat{\boldsymbol{\beta}}, \boldsymbol{\Sigma} = \frac{s^2 (\mathbf{X}^T \mathbf{X})^{-1}}{n - k - 1}$$

- Thus the Bayes' estimate of $\boldsymbol{\beta}$ is: $\mathbb{E}^\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \hat{\boldsymbol{\beta}}$

- The HPD intervals coincide with the frequentist confidence intervals.

Zellner's Non-Informative G -Prior

- The main difference with informative G -prior setup is that **we now consider c as unknown.**
- We use the same G –prior distribution with $\tilde{\beta} = \mathbf{0}_{k+1}$ conditional on c , and introduce a diffuse prior on c ,

$$\pi(c) = c^{-1} \mathbb{I}_{N^*}(c)$$

- The corresponding marginal posterior is given as:

$$\begin{aligned}\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) &= \int \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}, c) \pi(c | \mathbf{y}, \mathbf{X}) dc \\ &\propto \sum_{c=1}^{\infty} \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}, c) \frac{f(\mathbf{y} / \mathbf{X}, c) \pi(c)}{\pi(\mathbf{y} / \mathbf{X})} \\ &\propto \sum_{c=1}^{\infty} \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}, c) f(\mathbf{y} / \mathbf{X}, c) c^{-1}\end{aligned}$$

Zellner's Non-Informative G-Prior

$$\pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}) \propto \sum_{c=1}^{\infty} \pi(\boldsymbol{\beta}, \sigma^2 | \mathbf{y}, \mathbf{X}, c) f(\mathbf{y} | \mathbf{X}, c) c^{-1}$$

□ We have proved earlier (for constant c) that:

$$f(\mathbf{y} | \mathbf{X}, c) = (c + 1)^{-(k+1)/2} \pi^{-n/2} \Gamma\left(\frac{n}{2}\right) \\ \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c + 1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\beta}} - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X} \tilde{\boldsymbol{\beta}} \right]^{-n/2}$$

□ This for our problem with $\tilde{\boldsymbol{\beta}} = \mathbf{0}_{k+1}$

$$f(\mathbf{y} | \mathbf{X}, c) \propto (c + 1)^{-(k+1)/2} \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c + 1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right]^{-n/2}$$

Zellner's Non-Informative G-Prior: Posterior Mean of β

□ Recall that

$$\beta | \mathbf{y}, \mathbf{X} \sim \mathcal{T}_{k+1} \left(n, \frac{c}{c+1} \left(\frac{\tilde{\beta}}{c} + \hat{\beta} \right), \frac{c \left(s^2 + (\tilde{\beta} - \hat{\beta})^T \mathbf{X}^T \mathbf{X} (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n(c+1)} (\mathbf{X}^T \mathbf{X})^{-1} \right)$$

□ The Bayes estimates of β for $\tilde{\beta} = \mathbf{0}_{k+1}$ is now given by:

$$\mathbb{E}_\pi(\beta | \mathbf{y}, \mathbf{X}) = \mathbb{E}_\pi[\mathbb{E}_\pi(\beta | \mathbf{y}, \mathbf{X}, c) | \mathbf{y}, \mathbf{X}] = \mathbb{E}_\pi \left[\frac{c}{c+1} \hat{\beta} | \mathbf{y}, \mathbf{X} \right] = \frac{\sum_{c=1}^{\infty} \frac{c}{c+1} \pi(c | \mathbf{y}, \mathbf{X})}{\sum_{c=1}^{\infty} \pi(c | \mathbf{y}, \mathbf{X})} \hat{\beta} \Rightarrow$$

$$\mathbb{E}_\pi(\beta | \mathbf{y}, \mathbf{X}) = \frac{\sum_{c=1}^{\infty} \frac{c}{c+1} f(\mathbf{y} | \mathbf{X}, c) c^{-1}}{\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}} \hat{\beta}$$

Zellner's Non-Informative G-Prior: Posterior Mean of σ^2

□ Recall that

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \mathcal{IG}\left(\frac{n}{2}, \frac{s^2}{2} + \frac{1}{2(c+1)} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})\right) \Rightarrow$$

$$\mathbb{E}[\sigma^2 | \mathbf{y}, \mathbf{X}] = \frac{s^2 + (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})^T \mathbf{X}^T \mathbf{X} (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) / (c+1)}{n-2}$$

□ The Bayes estimate of σ^2 is given similarly by:

$$\mathbb{E}_\pi(\sigma^2 | \mathbf{y}, \mathbf{X}) = \frac{\sum_{c=1}^{\infty} \frac{s^2 + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} / (c+1)}{n-2} f(\mathbf{y} | \mathbf{X}, c) c^{-1}}{\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}}$$

□ Both Bayes estimates involve infinite summations on c .
The denominator in both cases is the normalizing constant
of the posterior $\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}$



Zellner's Non-Informative G-Prior: Posterior Variance of β

$$\beta | y, X \sim \mathcal{T}_{k+1} \left(n, \frac{c}{c+1} \left(\frac{\tilde{\beta}}{c} + \hat{\beta} \right), \frac{c \left(s^2 + (\tilde{\beta} - \hat{\beta})^T X^T X (\tilde{\beta} - \hat{\beta}) / (c+1) \right)}{n(c+1)} (X^T X)^{-1} \right)$$

$$\begin{aligned}
V_\pi(\beta | y, X) &= \mathbb{E}_\pi[V_\pi(\beta | y, X, c) | y, X] + V_\pi[\mathbb{E}_\pi(\beta | y, X, c) | y, X] = \\
&\mathbb{E}_\pi \left[\frac{c}{(n-2)(c+1)} \left(s^2 + \hat{\beta}^T X^T X \hat{\beta} / (c+1) \right) (X^T X)^{-1} \right] + V_\pi \left[\frac{c}{(c+1)} \hat{\beta} | y, X \right] = \\
&\left[\frac{\sum_{c=1}^{\infty} \frac{f(y|X,c)}{(n-2)(c+1)} \left(s^2 + \hat{\beta}^T X^T X \hat{\beta} / (c+1) \right)}{\sum_{c=1}^{\infty} f(y|X,c)c^{-1}} \right] (X^T X)^{-1} + \\
&\hat{\beta} \left[\frac{\sum_{c=1}^{\infty} \left(\frac{c}{(c+1)} - \mathbb{E}_\pi \left(\frac{c}{(c+1)} | y, X \right) \right)^2 f(y|X,c)c^{-1}}{\sum_{c=1}^{\infty} f(y|X,c)c^{-1}} \right] \hat{\beta}^T
\end{aligned}$$

See [earlier slide](#) for the term $f(y|X,c)$



Zellner's Non-Informative G-Prior: Marginal Distribution

- The marginal distribution of the dataset is available in closed form

$$f(\mathbf{y}/\mathbf{X}) = \sum_{c=1}^{\infty} f(\mathbf{y}/\mathbf{X}, c) c^{-1} \propto \\ \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(k+1)/2} \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c+1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \right]^{-n/2}$$

- The \mathcal{T} -shape means that we can also compute the normalizing constant!

Zellner's Informative G-Prior: Posterior Mean and Variance

$$\mathbb{E}_\pi(\beta_i | \mathbf{y}, \mathbf{X})$$

$$V_\pi(\beta_i | \mathbf{y}, \mathbf{X})$$

$$\mathbb{E}_\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \frac{\sum_{c=1}^{\infty} \frac{c}{c+1} f(\mathbf{y} | \mathbf{X}, c) c^{-1}}{\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}} \hat{\boldsymbol{\beta}}$$

beta_i	$\mathbb{E}[\text{beta}_i \mathbf{y}, \mathbf{X}]$	$V[\text{beta}_i \mathbf{y}, \mathbf{X}]$
--------	--	---

(Intercept)	9.2714	9.6424
X1	-0.0037	0.0000
X2	-0.0454	0.0005
X3	0.0573	0.0092
X4	-1.0905	0.3079
X5	0.1953	0.0105
X6	-0.3008	2.2750
X7	-0.2002	0.9383
X8	0.1526	0.0522
X9	-1.0835	0.7063
X10	-0.3651	0.5020

$$V_\pi(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \left[\frac{\sum_{c=1}^{\infty} \frac{f(\mathbf{y} | \mathbf{X}, c)}{(n-2)(c+1)} (s^2 + \hat{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\beta}} / (c+1))}{\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}} \right] (\mathbf{X}^T \mathbf{X})^{-1} + \hat{\boldsymbol{\beta}} \left[\frac{\sum_{c=1}^{\infty} \left(\frac{c}{(c+1)} - \mathbb{E}_\pi \left(\frac{c}{(c+1)} | \mathbf{y}, \mathbf{X} \right) \right)^2 f(\mathbf{y} | \mathbf{X}, c) c^{-1}}{\sum_{c=1}^{\infty} f(\mathbf{y} | \mathbf{X}, c) c^{-1}} \right] \hat{\boldsymbol{\beta}}^T$$

A Matlab implementation can be downloaded [here](#)

(We use $\tilde{\boldsymbol{\beta}} = \mathbf{0}_{11}$, $c = 100$)



Zellner's Non-Informative G-Prior: Point Null Hypothesis

- If a null hypothesis is $H_0 : R\beta = r$, the model under H_0 can be rewritten as

$$y | \beta^0, \sigma^2, X_0 \stackrel{H_0}{\sim} \mathcal{N}_n(X_0 \beta^0, \sigma^2 I_n)$$

where β^0 is $(k + 1 - q)$ dimensional.

- Under $\pi(c) = c^{-1}$ and the prior

$$\beta^0 | \sigma^2, X_0, c \sim \mathcal{N}_{k+1-q} \left(\mathbf{0}_{k+1-q}, c \sigma^2 (X_0^T X_0)^{-1} \right)$$

the marginal distribution of y under H_0 is:

$$f(y | X_0, H_0) \propto \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(k+1-q)/2} \left[y^T y - \frac{c}{c+1} y^T X_0 (X_0^T X_0)^{-1} X_0^T y \right]^{-n/2}$$

- The Bayes factor $B_{10}^\pi = f(y | X) / f(y | X_0, H_0)$ can now be computed



Zellner's Non-Informative G-Prior: Posterior Mean and Variance

For $H_0 : \beta_7 = \beta_8 = 0$, $\log_{10}(B_{10}^\pi) = -0.7884$ (We use $\tilde{\beta} = \mathbf{0}_{11}$)

	$\mathbb{E}_\pi(\beta_i y, X)$	$V_\pi(\beta_i y, X)$	$\log_{10}(BF)$	
beta_i	$\mathbb{E}[\text{beta_i} y, X]$	$V[\text{beta_i} y, X]$	$\log_{10}(BF)$	
(Intercept)	9.2714	9.6424	1.4205 (***)	
X1	-0.0037	0.0000	0.8503 (**)	
X2	-0.0454	0.0005	0.5664 (**)	
X3	0.0573	0.0092	-0.3609	
X4	-1.0905	0.3079	0.4520 (*)	
X5	0.1953	0.0105	0.4007 (*)	
X6	-0.3008	2.2750	-0.4411	
X7	-0.2002	0.9383	-0.4404	Evidence against H_0 :
X8	0.1526	0.0522	-0.3383	(****) decisive
X9	-1.0835	0.7063	-0.0424	(***) strong
X10	-0.3651	0.5020	-0.3838	(**) substantial (*) poor

A Matlab implementation can be downloaded [here](#)



Variable Selection

- Let us return to our regression model with one dependent random variable y and a set of k $\{x_1, x_2, \dots, x_k\}$ explanatory variables.
- Are all the x_i 's needed in the regression?
- We assume that every q –subset $\{i_1, i_2, \dots, i_q\}$, $0 \leq q \leq k$, of the explanatory variables,

$$\left\{1, x_{i_1}, x_{i_2}, \dots, x_{i_q}\right\}$$

is a proper set of explanatory variables for the regression of y (as before, the intercept is included in all models).

- We have a total of 2^k models to select from!



Variable Selection

- Following earlier notation, we denote: $X = [\mathbf{1}_n \ x_1 \ x_2 \dots x_k]$ as the matrix that contains $\mathbf{1}_n$ and the k potential predictor variables.
- Each model M_γ is associated with binary indicator vector

$$\gamma \in \Gamma = \{0,1\}^k$$

where $\gamma_i = 1$ means that the variable x_i is included in the model M_γ and $\gamma_i = 0$ that it is not.

- The number of variables included in the model M_γ is:
- The indices of the variables included in the model and not included in the model are denoted, respectively, as: $t_1(\gamma), t_0(\gamma)$



Variable Selection - Models in Competition

- For $\beta \in \mathbb{R}^{k+1}$ and X , we define β_γ as the sub-vector

$$\beta_\gamma = \left(\beta_0, (\beta_i)_{i \in t_1(\gamma)} \right)$$

- Let X_γ be the submatrix of X where only the column $\mathbf{1}_n$ and the columns in $t_1(\gamma)$ have been left.
- The model M_γ is then defined as:

$$y | \gamma, \beta_\gamma, \sigma^2, X \sim \mathcal{N}_n(X_\gamma \beta_\gamma, \sigma^2 I_n)$$

where $\beta_\gamma \in \mathbb{R}^{q_\gamma + 1}$, $\sigma^2 \in \mathbb{R}_+^*$ are the unknown parameters.

- The σ^2 is common to all models and we use the same prior for all models.



Variable Selection - Models in Competition

- We have a high number 2^k of models in competition.
- We cannot specify a prior on every M_γ in a completely subjective and autonomous manner.
- We derive all priors from a single global prior **associated with the full model** that corresponds to $\gamma = (1, \dots, 1)$.

Zellner's Informative Prior: Variable Selection

- For the full model that corresponds to $\gamma = (1, \dots, 1)$, we use the Zellner's informative G-prior:

$$\begin{aligned}\boldsymbol{\beta} | \sigma^2, \mathbf{X} &\sim \mathcal{N}_{k+1}(\tilde{\boldsymbol{\beta}}, c\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}) \\ \sigma^2 &\sim \pi(\sigma^2 | \mathbf{X}) \propto \sigma^{-2} \text{ improper Jeffreys prior}\end{aligned}$$

- For each model M_γ , the prior distribution of β_γ conditional on σ^2 is fixed as:

$$\boldsymbol{\beta}_\gamma | \gamma, \sigma^2 \sim \mathcal{N}_{q_\gamma+1}(\tilde{\boldsymbol{\beta}}_\gamma, c\sigma^2(\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1})$$

where $\tilde{\boldsymbol{\beta}}_\gamma = (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \tilde{\boldsymbol{\beta}}$ and same prior on σ^2 .

Zellner's Informative Prior: Variable Selection - Prior

- The joint prior for model M_γ is the improper prior

$$\pi(\boldsymbol{\beta}_\gamma, \sigma^2 | \mathbf{y}) \propto (\sigma^2)^{-(q_\gamma+1)/2-1} \exp \left[-\frac{1}{2(c\sigma^2)} (\boldsymbol{\beta}_\gamma - \tilde{\boldsymbol{\beta}}_\gamma)^T (\mathbf{X}_\gamma^T \mathbf{X}_\gamma) (\boldsymbol{\beta}_\gamma - \tilde{\boldsymbol{\beta}}_\gamma) \right]$$

- Infinitely many ways of defining a prior on the model index γ :

Our choice is a uniform prior $p(\gamma | \mathbf{X}) = 2^{-k}$.

- Posterior distribution of γ is central to variable selection since it is proportional to marginal density of \mathbf{y} on M_γ (or evidence of M_γ)

Zellner's Informative Prior: Variable Selection - Prior

- Posterior distribution of γ is proportional to the marginal density of y on M_γ (so it can also be used to compute Bayes factors)

$$\begin{aligned}\pi(\gamma | \mathbf{y}, \mathbf{X}) &\propto f(\mathbf{y} | \gamma, \mathbf{X}) \pi(\gamma | \mathbf{X}) \propto f(\mathbf{y} | \gamma, \mathbf{X}) \\ &= \int \left(\int f(\mathbf{y} | \gamma, \boldsymbol{\beta}, \sigma^2, \mathbf{X}) \pi(\boldsymbol{\beta} | \gamma, \sigma^2, \mathbf{X}) d\boldsymbol{\beta} \right) \pi(\sigma^2 | \mathbf{X}) d\sigma^2\end{aligned}$$

where (see earlier derivation)

$$\begin{aligned}f(\mathbf{y} | \boldsymbol{\gamma}, \sigma^2, \mathbf{X}) &= \int f(\mathbf{y} | \boldsymbol{\gamma}, \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta} | \boldsymbol{\gamma}, \sigma^2) d\boldsymbol{\beta} = \\ &= (c + 1)^{-(q_\gamma + 1)/2} (2\pi)^{-n/2} (\sigma^2)^{-n/2} \times \\ \exp \left(-\frac{1}{2\sigma^2} \mathbf{y}^T \mathbf{y} + \frac{1}{2\sigma^2(c + 1)} \left\{ c \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \mathbf{y} - \tilde{\boldsymbol{\beta}}_\gamma^T \mathbf{X}_\gamma^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma + 2 \mathbf{y}^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma \right\} \right)\end{aligned}$$

Zellners Informative Prior: Variable Selection - Prior

$$\pi(\gamma | \mathbf{y}, \mathbf{X}) \propto \int f(\mathbf{y} | \gamma, \sigma^2, \mathbf{X}) \pi(\sigma^2 | \mathbf{X}) d\sigma^2 = \int f(\mathbf{y} | \gamma, \sigma^2, \mathbf{X}) \frac{1}{\sigma^2} d\sigma^2$$

□ Posterior distribution of γ is then given as:

$$f(\gamma | \mathbf{y}, \mathbf{X}) \propto (c + 1)^{-(q_\gamma + 1)/2} \times \\ \left(\mathbf{y}^T \mathbf{y} - \frac{c}{(c + 1)} \mathbf{y}^T \mathbf{X}_\gamma (\mathbf{X}_\gamma^T \mathbf{X}_\gamma)^{-1} \mathbf{X}_\gamma^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}_\gamma^T \mathbf{X}_\gamma^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X}_\gamma \tilde{\boldsymbol{\beta}}_\gamma \right)^{-n/2}$$

□ We already have seen this distribution earlier for a fixed c :

$$f(\mathbf{y} | \mathbf{X}, c) = (c + 1)^{-(k+1)/2} \pi^{-n/2} \Gamma\left(\frac{n}{2}\right) \left[\mathbf{y}^T \mathbf{y} - \frac{c}{c + 1} \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} + \frac{1}{c + 1} \tilde{\boldsymbol{\beta}}^T \mathbf{X}^T \mathbf{X} \tilde{\boldsymbol{\beta}} - \frac{2}{c + 1} \mathbf{y}^T \mathbf{X} \tilde{\boldsymbol{\beta}} \right]^{-n/2}$$

Model Selection

- Most likely models ordered by decreasing posterior probabilities using Zellner's informative G –prior with $c = 100$.

$t_1(\gamma)$	$\pi(\gamma y, X)$
t1_gamma	pi(gamma y, X)
0 1 2 4 5	0.231543
0 1 2 4 5 9	0.037358
0 1 9	0.034435
0 1 2 4 5 10	0.032975
0 1 4 5	0.030606
0 1 2 9	0.025016
0 1 2 4 5 7	0.024144
0 1 2 4 5 8	0.023784
0 1 2 4 5 6	0.023735
0 1 2 3 4 5	0.023207
0 1 6 9	0.014587
0 1 2 3 9	0.014491
0 9	0.014281
0 1 2 6 9	0.013551
0 1 4 5 9	0.012761
0 1 3 9	0.011712
0 1 2 8	0.011477
0 1 8	0.009519
0 1 2 3 4 5 9	0.009036
0 1 2 4 5 6 9	0.009031

An implementation
is available
[MatLab](#), [C++](#)



Model Selection

- Model M_γ with the highest posterior probability is $t_1(\gamma) = (1, 2, 4, 5)$, which corresponds to the variables
 - altitude,
 - slope,
 - height of the tree sampled in the center of the area, and
 - diameter of the tree sampled in the center of the area.



Model Selection

- For the Zellner's non-informative prior with $p(c) = 1/c$, we have ($\tilde{\beta} = \mathbf{0}_{k+1}$) :

$$\pi(\gamma | y, X) = \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(q_\gamma+1)/2} \left[y^T y - \frac{c}{c+1} y^T X_\gamma (X_\gamma^T X_\gamma)^{-1} X_\gamma^T y \right]^{-n/2}$$

- Again we have seen this before as

$$f(y/X) = \sum_{c=1}^{\infty} f(y/X, c) c^{-1} \propto \\ \sum_{c=1}^{\infty} c^{-1} (c+1)^{-(k+1)/2} \left[y^T y - \frac{c}{c+1} y^T X (X^T X)^{-1} X^T y \right]^{-n/2}$$

Model Selection

- Most likely models ordered by decreasing posterior probabilities using Zellner's non-informative G –prior.

$t_1(\gamma)$	$\pi(\gamma y, X)$
t1_gamma	$\pi(\text{gamma} y, X)$
0 1 2 4 5	0.092914
0 1 2 4 5 9	0.032553
0 1 2 4 5 10	0.029512
0 1 2 4 5 7	0.023114
0 1 2 4 5 8	0.022843
0 1 2 4 5 6	0.022807
0 1 2 3 4 5	0.022409
0 1 2 3 4 5 9	0.016733
0 1 2 4 5 6 9	0.016725
0 1 2 4 5 8 9	0.013726
0 1 4 5	0.011031
0 1 2 4 5 9 10	0.009933
0 1 2 3 9	0.009698
0 1 2 9	0.009316
0 1 2 4 5 7 9	0.009253
0 1 2 6 9	0.009189
0 1 4 5 9	0.008756
0 1 2 3 4 5 10	0.007933
0 1 2 4 5 8 10	0.007901
0 1 2 4 5 7 10	0.007896

An implementation
is available
[MatLab](#), [C++](#)



Stochastic Search for the Most Likely Model

- When k is large, it becomes computationally intractable to compute the posterior probabilities of the 2^k models.
- Need of a tailored algorithm that samples from $p(\gamma|y, X)$ and selects the most likely models.
- Can be done by Gibbs sampling*, given the availability of the **full conditional posterior probabilities of the γ_i 's**. If

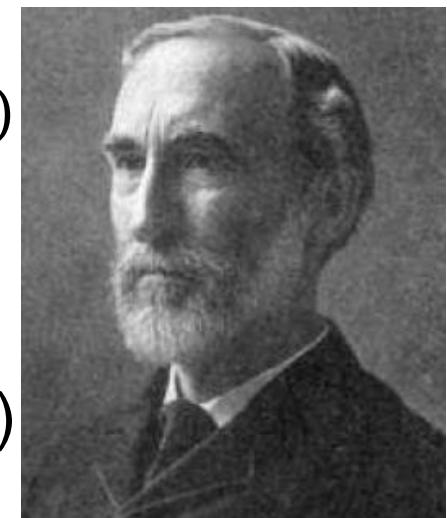
$\gamma_{-i} = (\gamma_1, \dots, \gamma_{i-1}, \gamma_{i+1}, \dots, \gamma_k)$ ($1 \leq i \leq k$)
then

$$p(\gamma_i | y, \gamma_{-i}, X) \propto p(\gamma | y, X)$$

(to be evaluated in both $\gamma_i = 0$ and $\gamma_i = 1$)

* Gibbs and other sampling algorithms will be introduced and discussed in detail in

forthcoming lectures.



Gibbs Sampling for Variable Selection

Initialization: Draw γ^0 from the uniform distribution on Γ

Iteration t: Given $(\gamma_1^{(t-1)}, \dots, \gamma_k^{(t-1)})$, generate

- $\gamma_1^{(t)}$ according to $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$
- $\gamma_2^{(t)}$ according to $\pi(\gamma_2 | \mathbf{y}, \gamma_1^{(t)}, \gamma_3^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$
- ..
- ..
- $\gamma_k^{(t)}$ according to $\pi(\gamma_k | \mathbf{y}, \gamma_1^{(t)}, \gamma_2^{(t)}, \dots, \gamma_{k-1}^{(t)}, \mathbf{X})$

Gibbs Sampling for Variable Selection

Question: How to sample γ_1^t according to $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$

1. The conditional distribution $\pi(\gamma_1 | \mathbf{y}, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)}, \mathbf{X})$ is proportional to $\pi(\gamma | \mathbf{y}, \mathbf{X})$
2. Since γ_1^t only has two possible values which are 0 and 1, we get

$$p_0 \propto \pi(\gamma_1 = 0, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)} | \mathbf{y}, \mathbf{X})$$

$$p_1 \propto \pi(\gamma_1 = 1, \gamma_2^{(t-1)}, \dots, \gamma_k^{(t-1)} | \mathbf{y}, \mathbf{X})$$

So the probability that $\gamma_1^t = 1$ is p_1 , then we can use Gibbs Sampling to approximate the distribution of $\{\gamma_i^t\}$

Gibbs Sampling: Posterior Probabilities

- After $T \gg 1$ MCMC iterations, we approximate the posterior probabilities $p(\gamma|y, X)$ by empirical averages

$$\hat{\pi}(\gamma|y, X) = \frac{1}{T - T_0 + 1} \sum_{t=T_0}^T \mathbb{I}_{\gamma^{(t)}=\gamma}$$

The T_0 first values (burn in) in the MCMC chain are eliminated.



Model Choice Comparison: Gibbs Estimates

First level Informative G-prior model with ($\tilde{\beta} = 0_{11}$, $c = 100$) compared with the Gibbs estimates of the top ten posterior probabilities

$t_1(\gamma)$	$\pi(\gamma y, X)$	$\hat{\pi}(\gamma y, X)$
t1_gamma	$\text{pi}(\text{gamma} y, X)$	$\text{pi_hat}(\text{gamma} y, X)$
0 1 2 4 5	0.231543	0.239276
0 1 2 4 5 9	0.037358	0.034397
0 1 9	0.034435	0.032397
0 1 2 4 5 10	0.032975	0.030097
0 1 4 5	0.030606	0.029397
0 1 2 9	0.025016	0.025297
0 1 2 4 5 7	0.024144	0.022498
0 1 2 4 5 8	0.023784	0.024898
0 1 2 4 5 6	0.023735	0.023598
0 1 2 3 4 5	0.023207	0.022998



A [MatLab](#) implementation is available

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2017, N. Zabaras)

Model Choice Comparison: Gibbs Estimates

Non-informative G-prior variable model choice compared with the Gibbs estimates of the top ten posterior probabilities

$t_1(\gamma)$	$\pi(\gamma y, X)$	$\hat{\pi}(\gamma y, X)$
t1_gamma	$\pi(\text{gamma} y, X)$	$\hat{\pi}(\text{gamma} y, X)$
0 1 2 4 5	0.092914	0.093391
0 1 2 4 5 9	0.032553	0.033097
0 1 2 4 5 10	0.029512	0.032597
0 1 2 4 5 7	0.023114	0.025097
0 1 2 4 5 8	0.022843	0.023098
0 1 2 4 5 6	0.022807	0.022498
0 1 2 3 4 5	0.022409	0.021698
0 1 2 3 4 5 9	0.016733	0.015998
0 1 2 4 5 6 9	0.016725	0.014899
0 1 2 4 5 8 9	0.013726	0.013399

[A MatLab](#) implementation is available



Gibbs Sampling: Probabilities of Inclusion

- An approximation of the probability to include the i-th variable:

$$\hat{P}^\pi(\gamma_i = 1 | \mathbf{y}, \mathbf{X}) = \frac{1}{T - T_0 + 1} \sum_{t=T_0}^T \mathbb{I}_{\gamma_i^{(t)}=1}$$



Probabilities of Inclusion Estimates

Informative ($\tilde{\beta} = 0_{11}$, $c = 100$) and non-informative G-prior variable inclusion estimates (based on the same Gibbs output as in the earlier two tables)

γ_i	$\hat{P}^\pi(\gamma_i = 1 \mathbf{y}, \mathbf{X})$	$\hat{P}^\pi(\gamma_i = 1 \mathbf{y}, \mathbf{X})$
gamma_i	P(gamma_i y, X)	P(gamma_i y, X)
gamma_1	0.8733	0.8806
gamma_2	0.7100	0.7789
gamma_3	0.1515	0.2958
gamma_4	0.6842	0.7422
gamma_5	0.6635	0.7234
gamma_6	0.1659	0.2992
gamma_7	0.1343	0.2812
gamma_8	0.1478	0.2740
gamma_9	0.3942	0.5015
gamma_10	0.1135	0.2556

[A MatLab implementation is available](#)

