

Estimação Bayesiana de Modelos Lineares com Aplicações em Economia

Schmidt, A.V.¹; Pozzobon, F.²

¹Graduanda, Departamento de Ciências Econômicas da Universidade do Estado de Santa Catarina

²Professor, Departamento de Ciências Econômicas da Universidade do Estado de Santa Catarina

Contato: aishameriane.schmidt@posgrad.ufsc.br

Introdução e Objetivos

Com o advento da computação, a econometria bayesiana tem cada vez mais espaço nos trabalhos acadêmicos, tanto por sua flexibilidade na incorporação de informações a priori como também pela possibilidade de gerar estimativas confiáveis mesmo na ausência ou escassez de eventos observáveis. Neste sentido, o presente trabalho se propõe a apresentar o modelo normal de regressão linear múltipla comparando os resultados obtidos pela estimação de mínimos quadrados ordinários (MQO) com as estimativas bayesianas utilizando uma priori conjugada normal-gama. Os principais teoremas das estimativas bayesianas para o modelo foram demonstrados analiticamente.

Adicionalmente, realizaram-se simulações onde foram comparados os resultados obtidos pela abordagem clássica e bayesiana, considerando diferentes especificações dos valores de parâmetros da densidade a priori e processo gerador dos dados. Nas simulações, constatou-se os resultados teóricos desenvolvidos, tanto com relação à convergência das estimativas bayesianas para o estimador de MQO com o aumento do tamanho da amostra, como para os efeitos de variação nas informações a priori (estes resultados estão somente no texto do trabalho). Ao final, utilizando o ferramental bayesiano apresentado, adaptou-se o estudo de Bonini e Pozzobon (2016) sobre discriminação de gênero e prêmio salarial no mercado de Tecnologia da Informação, utilizando os dados da RAIS de 2011 para trabalhadores de TI do Estado de Santa Catarina. As estimativas obtidas, utilizando como priori as estimativas reportadas no estudo original, foram consistentes com as de Bonini e Pozzobon (2016) e o desempenho da estimativa bayesiana na previsão de dados para fora da amostra foi superior do que o desempenho das estimativas de MQO.

Parte 1. Resultados teóricos do MNRL

O modelo normal de regressão linear múltipla pode ser escrito na forma matricial $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, em que \mathbf{Y} é um vetor $N \times 1$ contendo a variável endógena ou dependente; \mathbf{X} é uma matriz $N \times k$ cuja primeira coluna é toda de 1 e as demais colunas representam as variáveis exógenas (ou independentes); $\boldsymbol{\beta}$ é um vetor com k parâmetros e $\boldsymbol{\varepsilon}$ é um vetor de erros com tamanho N . Estamos assumindo que s erros tem distribuição normal multivariada com média zero e matriz de variância e covariância homocedástica, com $Var[\boldsymbol{\varepsilon}_i] = \sigma^2 = h^{-1}$ (o parâmetro h é chamado de *precisão*). Os elementos de \mathbf{X} são fixos ou, no caso de serem v.a., são independentes do termo de erro e tem f.d.p. dada por $p(\mathbf{X}|\boldsymbol{\lambda})$, em que $\boldsymbol{\lambda}$ não depende nem de $\boldsymbol{\beta}$ nem de h .

Na abordagem bayesiana é necessário definir uma distribuição à priori para os parâmetros do modelo, que neste caso são $\boldsymbol{\beta}$ e h . Neste trabalho optamos por utilizar uma priori conjugada natural (normal-gamma), que, por definição, leva a uma posteriori da mesma família de distribuições.

Ao utilizarmos uma priori conjugada natural temos o benefício de obter a posteriori em fórmula fechada, o que facilita o processo de atualização à medida que novos dados estão disponíveis. Ao mesmo tempo, é possível analisar a influência de cada componente da priori e da verossimilhança nos resultados do modelo. Estes benefícios não são gratuitos: nesta formulação não é possível ser informativo a respeito $\boldsymbol{\beta}$ ao mesmo tempo que se é não informativo a respeito de h , uma vez que a distribuição do primeiro é condicional ao segundo.

Como a posteriori é uma densidade conhecida, para fazer inferência neste modelo iremos precisar apenas das quantidades de MQO e dos hiperparâmetros da priori. Com isto podemos calcular momentos, intervalos de credibilidade e fazer comparação entre modelos.

Parte 2. Aplicação empírica

Na parte empírica do trabalho, atualizamos uma parte do estudo de Bonini e Pozzobon (2016) que investigaram as diferenças salariais entre homens e mulheres trabalhadores de setores específicos dos três estados da região Sul do Brasil no ano de 2011. Neste trabalho optamos por investigar apenas os trabalhadores da área de tecnologia da informação (TI) de Santa Catarina para verificar se as diferenças salariais devidas ao gênero se modificaram quando analisados os dados de 2016. As priors dos parâmetros são dadas por $\boldsymbol{\beta}|h \sim N(\underline{\boldsymbol{\beta}}, h^{-1}\underline{\mathbf{V}})$ e $h \sim \text{Gamma}(\underline{s}^{-2}, \underline{\mathbf{v}})$. Os resultados de Bonini e Pozzobon (2016) foram utilizados como hiperparâmetros e para acomodar a incerteza temporal, adotamos um valor alto da variância (10^2).

Uma vez que os logaritmos dos salários estão quase totalmente restritos ao intervalo 3.0-5.5 e por hipótese os erros são normalmente distribuídos, utilizamos $s^2 = 1$, o que implica que a priori de h é igual a 1. Como este valor é bastante arbitrário, utilizamos um valor baixo para os g.l. da priori: $\underline{\mathbf{v}} = 44$.

Aplicação empírica (dados)

Utilizamos dados da RAIS para o estado de Santa Catarina de todos trabalhadores registrados em ocupações na área de TI no ano de 2016, totalizando 10,919 pessoas. As estatísticas descritivas do banco de dados estão nas tabelas e gráficos abaixo.

	ln(sal) (R\$)	T. Emprego (meses)	Idade (anos)
Mínimo	2.067	0.00	17.00
1º quartil	4.112	14.80	27.00
Mediana	4.578	36.90	31.00
Média	4.580	55.41	32.64
3º quartil	5.015	76.50	37.00
Máximo	7.707	507.80	69.00
Desvio Padrão	0.656	62.87	8.11

Tabela 1: Descriptives of the quantitative variables entered in the model.

	n	%
Homens	8792	80.59
Ens. Superior Incompl.	1157	10.61
Superior completo	8339	76.44
Pós Graduação	184	1.69

Tabela 2: Frequências das categorias na amostra. N corresponde ao total de indivíduos na categoria. Na escolaridade o valor representa o total de pessoas cuja escolaridade máxima é a linha correspondente.

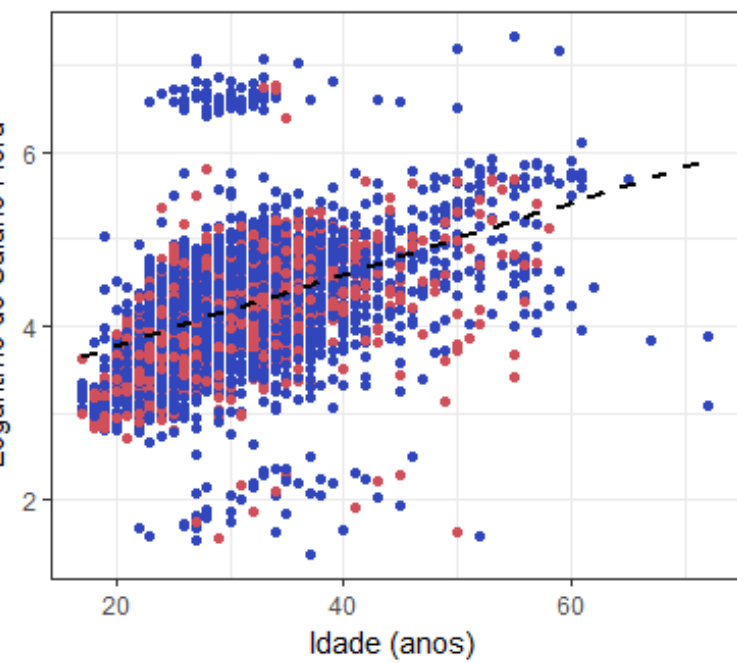


Gráfico 1: Distribuição do salário e idade, por gênero.

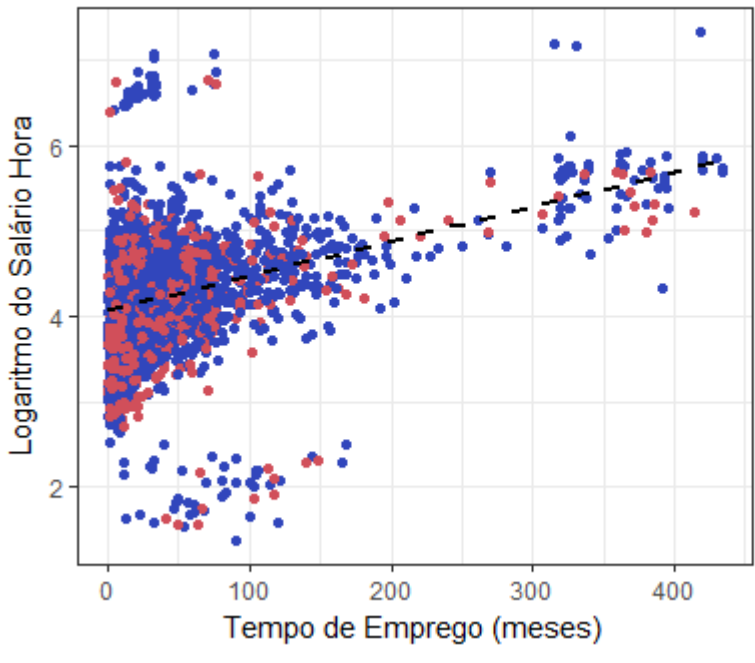


Gráfico 2: Distribuição dos salários por tempo de emprego e gênero.

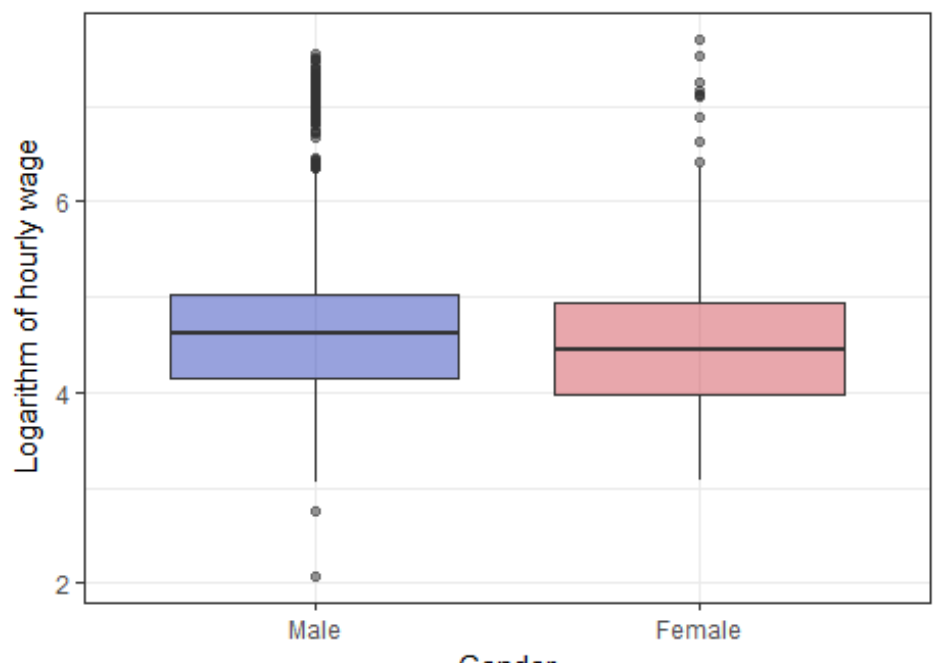


Gráfico 3: Boxplot dos salários por gênero.

Motivação

A respeito do uso de métodos bayesianos em econometria, Geweke (2001) e Zellner (1985) citam autores da metade do século XX como precursores no desenvolvimento do assunto, sendo que um dos primeiros livros textos de econometria bayesiana é a obra de Zellner (1971), intitulada “An introduction to bayesian inference in econometrics”. Daquela época até os dias atuais, com o advento da computação, métodos bayesianos se tornaram cada vez mais acessíveis e inúmeros trabalhos foram desenvolvidos na área.

Atualmente, existem diversos modelos macroeconômicos onde os métodos bayesianos são amplamente utilizados, como por exemplo os modelos de vetores aleatórios (VAR), desenvolvidos por Christopher Sims no início da década de 80. Doan, Litterman e Sims (1984) estimaram pela primeira vez um VAR bayesiano (BVAR) enquanto que a generalização do modelo considerando parâmetros variando no tempo (TVP-VAR), desenvolvida por Cogley e Sargent em 2001, já foi estimada utilizando o amostrador de Gibbs, um método bayesiano. Mais tarde, em 2005, os mesmos autores incluíram volatilidade estocástica multivariada no TVP-VAR e novamente estimaram o modelo utilizando o amostrador de Gibbs. De acordo com Geweke, Koop e Van Dijk (2011), os modelos dinâmicos estocásticos de equilíbrio geral (DSGE), utilizados amplamente em Bancos Centrais, são predominantemente estimados utilizando inferência bayesiana. De fato, desde 2011 o Banco Central do Brasil (BACEN) utiliza o SAMBA (modelo analítico estocástico com abordagem bayesiana) para auxiliar na condução da política macroeconômica no país. Caldeira, Moura e Santos (2015) fazem uma comparação de modelos, incluindo um TVP-VAR bayesiano com uso de priori de Minnesota para previsão de dados macroeconômicos do Brasil.

Sob as hipóteses descritas anteriormente, a verossimilhança do modelo será dada por:

$$p(\mathbf{Y}|\boldsymbol{\beta}, h) = \frac{h^{N/2}}{(2\pi)^{N/2}} \exp \left[-\frac{h}{2} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta})' (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta}) \right]. \quad (1)$$

Podemos combinar (1) com uma priori normal para $\boldsymbol{\beta}|h$ vezes uma densidade gama para h de forma a obter a seguinte expressão a posteriori:

$$p(\boldsymbol{\beta}, h|\mathbf{Y}) \propto h^{\frac{\bar{\mathbf{v}}+k-2}{2}} \exp \left[-\frac{h}{2} \left(\bar{\mathbf{v}}\bar{s}^2 + (\boldsymbol{\beta} - \bar{\boldsymbol{\beta}})' \bar{\mathbf{V}}^{-1} (\boldsymbol{\beta} - \bar{\boldsymbol{\beta}}) \right) \right]. \quad (2)$$

O lado direito de (2) é o núcleo de uma densidade Normal-Gama com parâmetros:

- $\bar{\mathbf{V}} = (\underline{\mathbf{V}}^{-1} + \mathbf{X}'\mathbf{X})^{-1}$;
- $\bar{\boldsymbol{\beta}} = \bar{\mathbf{V}} (\underline{\mathbf{V}}^{-1}\underline{\boldsymbol{\beta}} + \mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}})$, em que $\hat{\boldsymbol{\beta}}$ é o estimador de MQO, enquanto que $\underline{\mathbf{V}}$ e $\underline{\boldsymbol{\beta}}$ são hiperparâmetros da priori. Logo, a expressão nos diz que o parâmetro $\bar{\boldsymbol{\beta}}$ é uma média ponderada da média à priori com o estimador de MQO. Se a precisão $\underline{\mathbf{V}}^{-1}$ for baixa, a contribuição da priori é pequena;
- $\bar{\mathbf{v}} = \underline{\mathbf{v}} + N$, ou seja, os graus de liberdade da priori mais o tamanho amostral resultam nos g.l. da posteriori;
- $\bar{\mathbf{v}}\bar{s}^2 = \underline{\mathbf{v}}s^2 + \mathbf{v}s^2 + (\hat{\boldsymbol{\beta}} - \underline{\boldsymbol{\beta}})' [\underline{\mathbf{V}} + (\mathbf{X}'\mathbf{X})^{-1}]^{-1} (\hat{\boldsymbol{\beta}} - \underline{\boldsymbol{\beta}})$, na qual $\underline{\mathbf{v}}s^2$ é a SQR da priori, $\mathbf{v}s^2$ é a SQR da verossimilhança e o último termo serve de penalização para diferenças entre a média à priori e o estimador de MQO.

O modelo estimado é dado por:

$\ln(\text{sal}) = \boldsymbol{\beta}_0 + \boldsymbol{\beta}_1 \cdot \mathbf{S}_1 + \boldsymbol{\beta}_2 \cdot \mathbf{S}_2 + \boldsymbol{\beta}_3 \cdot \mathbf{S}_3 + \boldsymbol{\beta}_4 \cdot \text{Idade} + \boldsymbol{\beta}_5 \cdot \text{Idade}^2 + \boldsymbol{\beta}_6 \cdot \text{TE} + \boldsymbol{\beta}_7 \cdot \text{Gênero} + \boldsymbol{\varepsilon}$, em que:

- $\ln(\text{sal})$ é o vetor com o logaritmo natural dos salários (em reais);
- $\boldsymbol{\beta}_0$ é o intercepto do modelo que representa fatores não observáveis comuns a todos trabalhadores;
- $\mathbf{S}_1, \mathbf{S}_2$ e \mathbf{S}_3 são variáveis dummy que representam a escolaridade – superior incompleto, completo e pós graduação. Quando as três são iguais a zero, a escolaridade é Ensino médio completo;
- Idade é a idade em anos, que estamos conjecturando ter retornos marginais decrescentes;
- TE é o tempo de emprego, em meses;
- Gênero é uma variável binária que assume valor 0 para mulheres e 1 para homens;
- as hipóteses sobre os erros são as mesmas anteriores.

Aplicação empírica (resultados)

As distribuições a posteriori dos parâmetros são dadas por $\boldsymbol{\beta}|\mathbf{y} \sim t(\bar{\boldsymbol{\beta}}, \bar{s}^2\bar{\mathbf{V}}, \bar{\mathbf{v}})$, $h|\mathbf{y} \sim (\bar{s}^{-2}, \bar{\mathbf{v}})$ e sua conjunta é $\boldsymbol{\beta}, h|\mathbf{y} \sim NG(\bar{\boldsymbol{\beta}}, \bar{\mathbf{V}}, \bar{s}^2, \bar{\mathbf{v}})$, cujos parâmetros podem ser consultados em Koop (2003). As médias a posteriori, desvio padrão, e probabilidades associadas aos $\boldsymbol{\beta}_i$ foram calculados utilizando a distribuição t .

Parâmetro	Média Priori	Média posteriori	Des. Pad. Posteriori	IC95%	P($\boldsymbol{\beta}_i > 0$)
$\boldsymbol{\beta}_0$	2.07	1.8739	0.0647	[1.747;2.001]	1
$\boldsymbol{\beta}_1$	0.12	0.0949	0.0198	[0.0561;0.1338]	1
$\boldsymbol{\beta}_2$	0.45	0.4943	0.0153	[0.4642;0.5243]	1
$\boldsymbol{\beta}_3$	1.07	0.9525	0.0387	[0.8766;1.0284]	1
$\boldsymbol{\beta}_4$	0.13	0.1053	0.0037	[0.0981;0.1126]	1
$\boldsymbol{\beta}_5$	-0.00143	-0.0011	≈ 0	[-0.0012;-0.001]	0
$\boldsymbol{\beta}_6$	0.00271	0.0030	0.0001	[0.0028;0.0032]	1
$\boldsymbol{\beta}_7(\text{Gênero})$	-0.15	-0.1334	0.0117	[-0.1563;-0.1105]	0

Tabela 3: Resultados a posteriori

Conclusões

O MNRL em sua versão bayesiana confere mais flexibilidade aos pesquisadores pois permite a inclusão de informações a priori além de cálculos de probabilidades a posteriori e maior facilidade na atualização dos resultados quando novos dados estão disponíveis. A priori conjugada natural normal-gamma permite que o procedimento de inferência seja feito em fórmula fechada, o que torna o processo computacional mais simples e rápido.

Neste trabalho vimos os resultados teóricos do MNRL bem como os resultados das simulação baseadas em experimentos de Monte Carlo (estes resultados foram omitidos do poster por questões de espaço). Por fim, propusemos um exercício empírico onde o modelo teórico utilizado no trabalho de Ponini e Pozzobon (2015) foi adaptado para uma subamostra com dados mais atuais. Os resultados mostraram concordância entre os dois estudos no sentido de que os salários das mulheres são mais baixos que os salários dos homens na amostra estudada.