

# ECON220B Discussion Section 3

## Linear Regression and Bayesian Inference

---

Lapo Bini

# Roadmap

---

1. Where is  $\hat{\beta}^{\text{OLS}}$  going?
2. Linear Projection
3. Ridge Regression
4. Bayesian Inference

# Understanding The Assumptions

---

**Linear regression:**  $y_i = x_i^T \beta + u_i$  with  $\{u_1, \dots, u_n\}$  iid,  $E[u_i] = 0$ .

1. If  $u_i|x_i \sim \mathcal{N}(0, \sigma^2)$  then  $\hat{\beta}^{OLS}$  is BLUE by Markov-Gauss theorem,  
 $\hat{\beta}^{OLS} = \hat{\beta}^{MLE}$ . We are estimating a causal effect  $x \rightarrow y$ , i.e.

$$\frac{\partial}{\partial x_i} E[y_i|x_i] = \beta$$

2. If  $E[u_i|x_i] = 0$  and,  $E[u_i^2|x_i] = \sigma^2$ , then  $\hat{\beta}^{OLS}$  is BLUE by Markov-Gauss theorem. We are estimating a causal effect.
3. If  $E[u_i|x_i] \neq 0$  but  $E[u_i x_i] = 0$  still holds then  $\hat{\beta}^{OLS} \xrightarrow{P} \beta$  but we are estimating correlation between  $x$  and  $y$ , no partial effects.

## Example (1/2)

**Model:**  $y_i = \beta x_i + u_i$   $u_i = x_i^2 + \eta_i$  with true parameter  $\beta = 3$ , and

$x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ ,  $x_i \perp \eta_i$ .

$$EL[x_i m_i] = 0 \Rightarrow Cov(x_i m_i) = 0$$

(1) Suppose we estimate the model by OLS, can we apply Markov-Gauss theorem?

(2) Is  $\hat{\beta}^{OLS}$  consistent for the true  $\beta$ ?

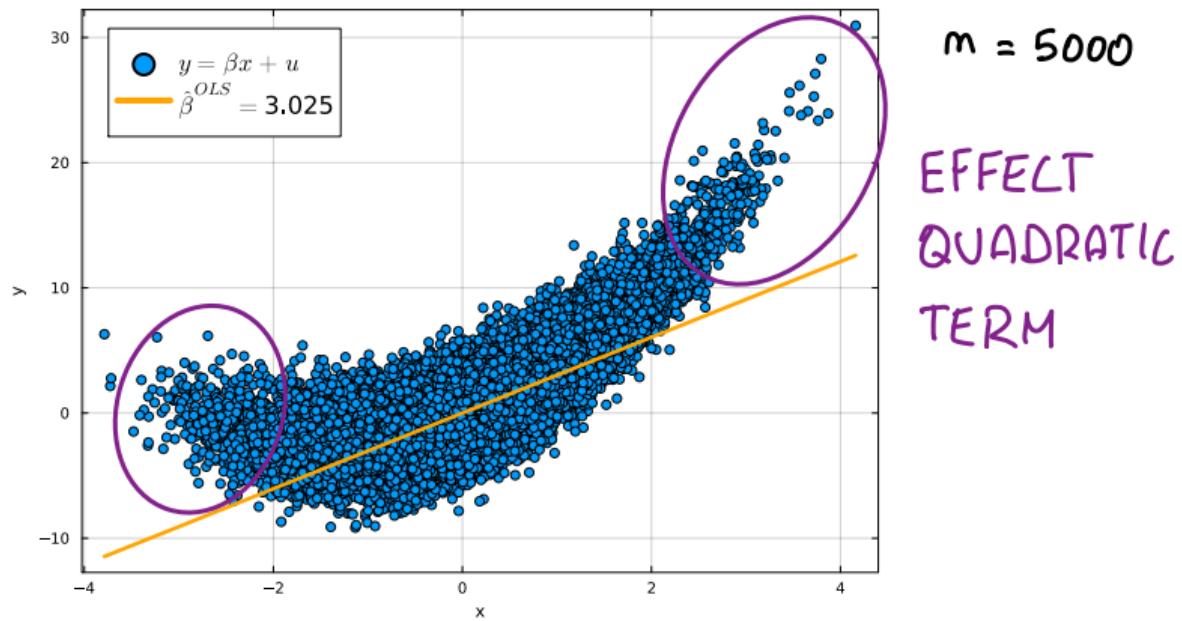
(1)  $ELu_i | x_i] = x_i^2$  NO WE CAN'T. NOTE :  $\frac{\partial}{\partial \beta} ELy_i | x_i] = \overbrace{\beta}^{\text{APE}} + 2x_i \overbrace{\eta_i}^{\text{HETEROG. EFFECT}}$

(2)  $\hat{\beta} = \left( \frac{1}{m} \sum x_i x_i^T \right)^{-1} \left( \frac{1}{m} \sum x_i y_i \right) = \left( \frac{1}{m} \sum x_i x_i^T \right)^{-1} \left( \frac{1}{m} \sum x_i (x_i^T \beta + x_i^2 + \eta_i) \right)$   
 $= \left( \frac{1}{m} \sum x_i x_i^T \right)^{-1} \left( \frac{1}{m} \sum x_i x_i^T \beta \right) + \left( \frac{1}{m} \sum x_i x_i^T \right)^{-1} \underbrace{\left( \frac{1}{m} \sum x_i x_i^2 \right)}_{\rightarrow ELx^3] \text{ but } x \sim N(0,1)} + \left( \frac{1}{m} \sum x_i x_i^T \right)^{-1} \underbrace{\left( \frac{1}{m} \sum x_i \eta_i \right)}_{\rightarrow ELx_i \eta_i] \text{ skew}(x) = 0 \therefore ELx^3] = 0}$   
 $= \beta + 0 + 0$

## Example (1/2)

---

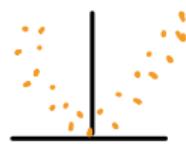
**Model:**  $y_i = \beta x_i + u_i$ ,  $u_i = x_i^2 + \eta_i$  with true parameter  $\beta = 3$ , and  $x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ ,  $x_i \perp \eta_i$ .



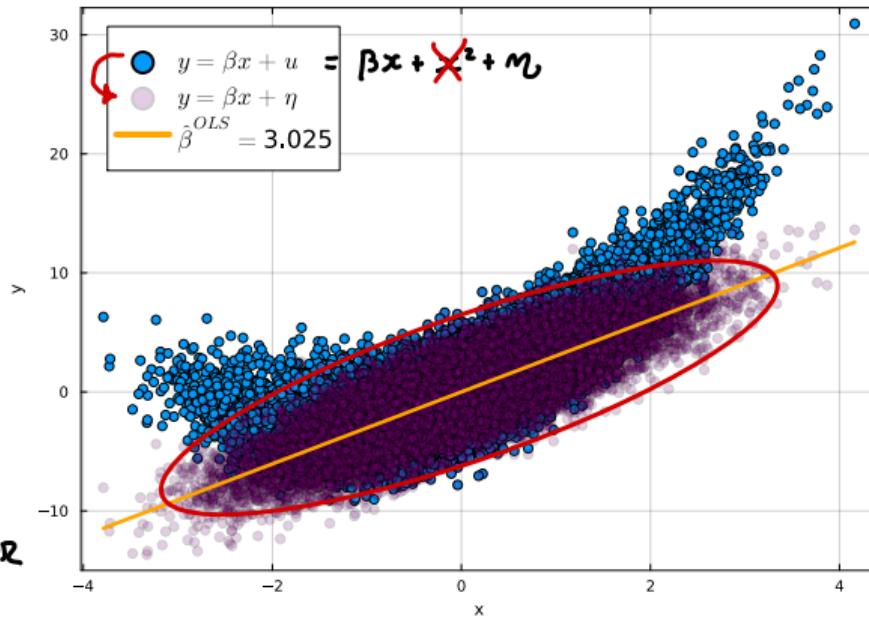
## Example (1/2)

**Model:**  $y_i = \beta x_i + u_i$   $u_i = x_i^2 + \eta_i$  with true parameter  $\beta = 3$ , and  $x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ ,  $x_i \perp \eta_i$ .

$\hat{\beta}^{OLS}$  CORR.



- OFFSET
  - NON LINEAR
  - $\text{Corr}(x, u) = 0$
- LINEAR ESTIMATOR



## Example (2/2)

---

Now we have:  $y_i = \beta x_i + \eta_i$  with true parameter  $\beta = 3$ , and  $x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ .

$$\hookrightarrow x_i = \frac{1}{\beta} y_i - \frac{1}{\beta} \eta_i$$

- (1) Suppose that instead of running a regression of  $y_i$  on  $x_i$ , you run the regression of  $x_i$  and  $y_i$ , that is you switch the dependent and independent variables:

$$x_i = \phi y_i + \nu_i$$

What is  $\hat{\phi}^{OLS}$  estimating?

$$\begin{aligned} E[\nu_i | y_i] &= -\frac{1}{\beta} E[\eta_i | y_i] \\ &= -\frac{1}{\beta} E[\eta_i | \beta x_i + \eta_i] = -\frac{\beta}{\beta} E[\eta_i | x_i] - \frac{1}{\beta} E[\eta_i^2] = -\frac{\sigma^2}{\beta} \neq 0 \end{aligned}$$

NOT  
CONSISTENT

## Example (2/2)

Now we have:  $y_i = \beta x_i + \eta_i$  with true parameter  $\beta = 3$ , and  $x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ .

- (1) Suppose that instead of running a regression of  $y_i$  on  $x_i$ , you run the regression of  $x_i$  and  $y_i$ , that is you switch the dependent and independent variables:

$$x_i = \phi y_i + \nu_i$$

$$\sigma^2 = 4$$

$$\phi = 1/\beta = 0.33$$

What is  $\hat{\phi}^{OLS}$  estimating?  $\hat{\phi} = (\frac{1}{m} \sum y_i y_i^T)^{-1} (\frac{1}{m} \sum y_i x_i)$

$$(\frac{1}{m} \sum y_i y_i^T)^{-1} (\frac{1}{m} \sum y_i (\phi y_i + \nu_i)) = \phi + (\frac{1}{m} \sum y_i^2)^{-1} (\frac{1}{m} \sum \nu_i y_i) = \phi + \frac{(-\sigma^2/\beta)}{\text{Var}(y_i) + \text{Var}(\nu_i)}$$

$\xrightarrow{P} \text{EL}[y_i^2] = \text{Var}(y_i) + \text{EL}[y_i]^2$

$$\therefore \hat{\phi}^{OLS} \xrightarrow{P} \phi + \Delta$$

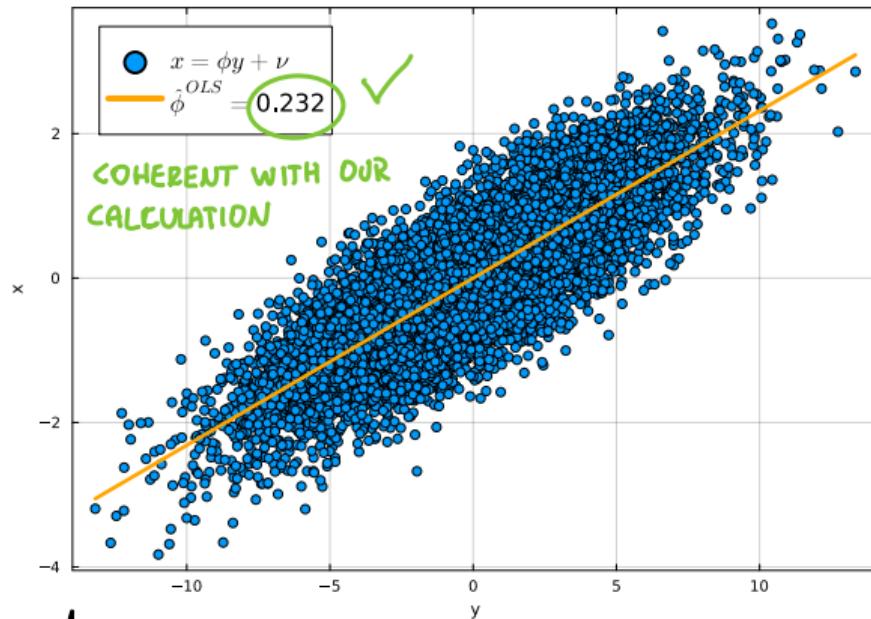
$$\hat{\phi}^{OLS} \xrightarrow{P} 0.33 + \frac{(-4/3)}{13} = 0.33 - 0.10$$

$\rightarrow \text{Var}(\beta x_i) + \text{Var}(\eta_i) = 9 \cdot 1 + 4 = 13$	$\rightarrow (\text{EL}[y_i])^2 = (\text{EL}[\beta x_i] + \text{EL}[\eta_i])^2 = 0^2$
---	---

## Example (2/2)

---

Now we have:  $y_i = \beta x_i + \eta_i$  with true parameter  $\beta = 3$ , and  $x_i \sim \mathcal{N}(0, 1)$ ,  $\eta_i \sim \mathcal{N}(0, 4)$ .



$$\hat{\phi}_{OLS} \xrightarrow{\rho} \phi + \Delta := \delta$$

# Linear Projection

---

- If  $E[u_i x_i] \neq 0$  then  $\hat{\beta}^{OLS} \xrightarrow{P} \delta \equiv \beta + \Delta$  it converges to the coefficient of the **linear projection**.
- The linear projection  $y_i = x_i^T \delta + u_i$  is also called the **minimum mean square linear predictor** since  $\delta$  solves the following problem:

$$\min_{\mathbf{d} \in \mathbb{R}^k} E[(y_i - x_i^T \mathbf{d})^2]$$

- The linear projection **always** satisfies  $E[x_i u_i] = 0$  and  $E[u_i] = 0$ .

$$\hat{u}_i = y_i - x_i^T \hat{\beta} = x_i^T \beta + u_i - x_i^T \hat{\beta} = u_i + x_i^T (\beta - \hat{\beta}) \quad \text{TAKEN SUM AND}$$

$$\text{DIVIDE BOTH TERMS BY } m: \quad \frac{1}{m} \sum \hat{u}_i = \frac{1}{m} \sum u_i + \underbrace{\frac{1}{m} \sum x_i^T (\beta - \hat{\beta})}_{\text{CAN WE APPLY WLLN? IS THIS ZERO?}}$$

$$\hat{u}_i = y_i - x_i^T \hat{\beta} = x_i^T \delta + u_i - x_i^T \hat{\beta}$$

$$\hat{u}_i = u_i + x_i^T (\delta - \hat{\beta})$$

$$\textcircled{a} \quad \frac{1}{m} \sum \hat{u}_i = \frac{1}{m} \sum u_i + \frac{1}{m} \sum x_i^T (\delta - \hat{\beta})$$

$\hat{\beta} \xrightarrow{P} \delta$

$$\textcircled{b} \quad x_i \hat{u}_i = x_i u_i + x_i x_i^T (\delta + \hat{\beta})$$

$$\frac{1}{m} \sum x_i \hat{u}_i = \frac{1}{m} \sum x_i u_i + \underbrace{\frac{1}{m} \sum x_i x_i^T (\delta + \beta)}_{Op(1)}$$

$$\underbrace{\frac{1}{m} \sum x_i \hat{u}_i}_{=0} \xrightarrow{P} E[x_i u_i]$$

## When Does $E[u_i x_i] = 0$ Fail?

---

- **Omitted variable bias:** consider the following linear regression model  
 $y_i = \beta_0 + \beta_1 x_i + \beta_2 z_i + u_i$  where  $y_i, x_i, z_i, u_i$  are all scalars and  
 $E[u_i x_i] = E[u_i z_i] = 0$
- Suppose we regress  $y_i$  on  $x_i$  only: what is the probability limit of  $\hat{\beta}_1^{OLS}$ ?  
When does the limit coincide with the true parameter  $\beta_1$ ?

WE ARE GOING TO  
START FROM AN  
EXAMPLE BEFORE  
FORMAL ANALYSIS

---

## Two Religions: Frequentists vs Bayesians

Given  $\{y_1, \dots, y_n\}$  iid sample with  $y_i \sim \mathcal{N}(\mu, \sigma^2)$  we are interested in the population mean  $\mu$ . We already know that MLE estimator is  $\hat{\mu}^{MLE} = n^{-1} \sum_{i=1}^n y_i \sim \mathcal{N}(\mu, \sigma^2/n)$ . Two different approaches:

1. **Frequentist**: the data is the result of sampling from a random process. Frequentists see the data as varying and the parameter  $\mu$  of this random process that generates the data as being fixed.  $\mathcal{N}(\mu, \sigma^2/n)$  describes a distribution across different samples.
2. **Bayesian**:  $\mu$  treated as a random variable. Bayesians have prior beliefs about  $\mu$  (**prior distribution**), which is updated after observing the data (**likelihood function**) using **Bayes' Rule**. The **posterior distribution** summarises the uncertainty about credible values of  $\mu$ .

# Ridge Regression

---

- Consider the follow linear regression model  $y_i = x_i^T \beta + u_i$ ,  $u_i \sim \mathcal{N}(0, 1)$ .
- Assume that the parameters  $\beta \in \mathbb{R}^d$  follow the distribution  $\beta \sim \mathcal{N}(0, \lambda^2 I_d)$  where  $\lambda > 0$  and  $I_d$  is the  $(dx d)$  identity matrix.
- Lastly, assume that  $u_i, x_i, \beta$  are mutually independent.

- (1) Prove that  $f_\beta(\beta) = \lambda^{-d} \prod_{j=1}^d \phi(\beta_j / \lambda)$ .
- (2) Show that  $f_{\mathbf{Y}|\beta, \mathbf{X}}(y_1, \dots, y_n | \beta, \mathbf{X}) = \prod_{i=1}^n \phi(y_i - x_i^T \beta)$ .
- (3) Derive the Maximum Likelihood Estimator  $\hat{\beta}^{MLE}$ .
- (4) Find the posterior distribution  $f_{\beta|\mathbf{Y}, \mathbf{X}}(\beta | \mathbf{Y}, \mathbf{X})$  and derive the Bayes estimator defined as

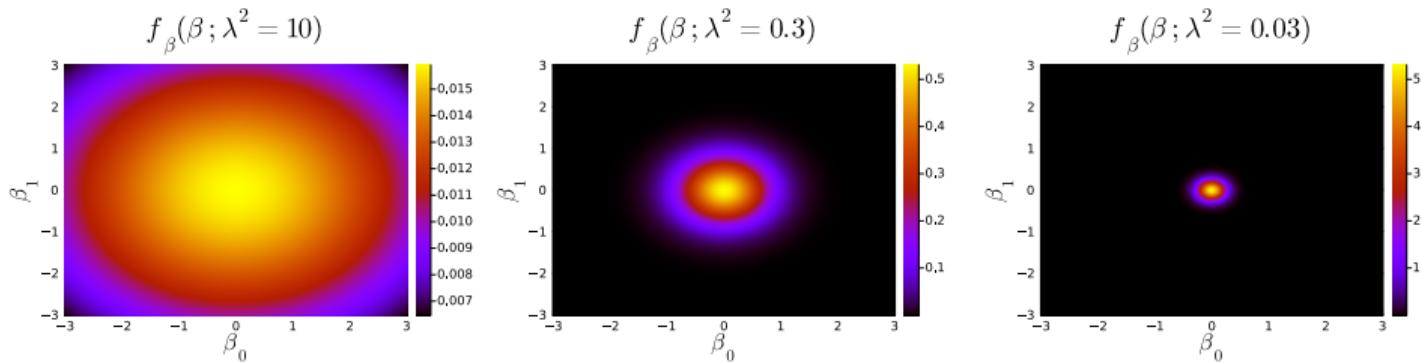
$$\hat{\beta}^{Bayes} \equiv \arg \max_{\beta} f_{\beta|\mathbf{Y}, \mathbf{X}}(\beta | \mathbf{Y}, \mathbf{X})$$

# Ridge Regression - Prior Distribution

---

Before observing the data, our **prior belief** is that the parameters are most likely to be close to zero. The parameter  $\lambda^2$  represents the uncertainty of our guess, i.e.  $\beta \sim \mathcal{N}(0, \lambda^2 I_2)$ .

**Figure:** Prior distribution for different values of  $\lambda^2$ .



$$\textcircled{1} \quad \beta \sim N(0, \lambda^2 \text{Id})$$

UNCORRELATED  $\Rightarrow$  INDEP.

$$\therefore \beta_i \sim N(0, \lambda^2)$$

$$f_{\beta} = \prod_{i=1}^m f_{\beta_i} = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\lambda} e^{\left\{-\frac{1}{2}(\beta_i/\lambda)^2\right\}}$$

$$= \frac{1}{\lambda^m} \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{\left\{-\frac{1}{2}(\beta_i/\lambda)^2\right\}}$$

$$\therefore f_{\beta}(b) = \frac{1}{\lambda^m} \prod_{i=1}^m \phi\left(\frac{\beta_i}{\lambda}\right) \quad \frac{\beta}{\lambda} \sim N(0, 1)$$

$$② y_i = x_i^T \beta + u_i$$

$$y_i | x_i^T \beta = k + u_i = N(k, 1)$$

$$= N(k, 1)$$

$\{y_1, \dots, y_m\}$  IID

$$\begin{aligned} \therefore f_{Y|Bx} &= \prod_{i=1}^m f_{Y_i|Bx} = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{\left\{-\frac{1}{2}(y_i - k)^2\right\}} \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{\left\{-\frac{1}{2}(y_i - x_i^T \beta)^2\right\}} \quad \text{PDF STANDARD NORMAL} \\ &= \prod_{i=1}^m \phi(y_i - x_i^T \beta) \end{aligned}$$

$$③ \hat{\beta} = \arg \max \prod_{i=1}^m \phi(y_i - x_i^T \beta)$$

$$LIK = \prod_{i=1}^m \phi(y_i - x_i^T \beta)$$

$$LIK = (2\pi)^{-\frac{m}{2}} \exp \left\{ -\frac{1}{2} \sum (y_i - x_i^T \beta)^2 \right\}$$

NOW I TAKE LOG

$$\log(LIK) = \log(2\pi)^{-\frac{m}{2}} - \frac{1}{2} \sum (y_i - x_i^T \beta)^2$$

max  $\log(LIK)$  IS EQUIVALENT TO

$$\min \frac{1}{2} \sum (y_i - x_i^T \hat{\beta})^2$$

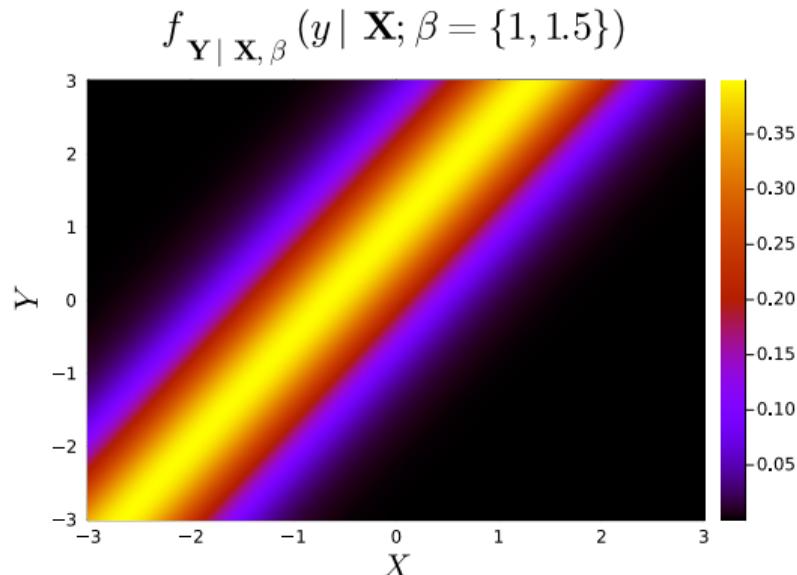
$$FOC \quad - \sum x_i (y_i - x_i^T \hat{\beta}) = 0$$

$$\hat{\beta} = (\frac{1}{m} \sum x_i x_i^T)^{-1} (\frac{1}{m} \sum x_i y_i) = \hat{\beta}_{OLS}$$

# Ridge Regression - Likelihood Function

---

The **likelihood** describes the probability of the data that has already been observed given certain parameter values  $\beta$ . Given different values of  $x_i$  and  $y_i$ , the points with highest probability lies on  $y_i = 1 + 1.5x_i$ .



## ④ POSTERIOR DISTRIBUTION

INDEPENDENT

$$a \cdot f_{Y|\beta X}(y, \beta, x) = f_{Y|BX}(y|\beta x) f_{\beta X}(\beta, x)$$

$$= f_{Y|BX}(y|\beta x) f_{\beta}(b) \cdot f_x(x)$$

$$b \cdot f_{Y|\beta X}(y, \beta, x) = f_{\beta|YX}(b|Yx) \cdot f_{YX}(y, x)$$

$$f_{\beta|YX} = \frac{f_{Y|\beta X}(y, \beta, x)}{f_{YX}(y, x)} =$$

$$f_{\beta|YX} = \frac{f_{Y|BX}(y|\beta x) f_{\beta}(b) \cdot f_x(x)}{f_{Y|X}(y|x) \cdot f_x(x)} \propto f_{Y|BX} \cdot f_{\beta}$$

## BAYES ESTIMATOR

$$\hat{\beta}^{\text{BAYES}} = \underset{\beta}{\operatorname{argmax}} f_{Y|BX} \cdot f_{\beta}$$

$$\left( \prod_{i=1}^m \phi(y_i - x_i^T \beta) \right) \left( \frac{1}{\lambda^d} \prod_{j=1}^d \phi\left(\frac{\beta_j}{\lambda}\right) \right) \downarrow$$

$$\left( \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(y_i - x_i^T \beta)^2} \right) \left( \frac{1}{\lambda^d} \prod_{j=1}^d \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(\beta_j/\lambda)^2} \right) \downarrow$$

$$\left( \frac{1}{\sqrt{2\pi}} \right)^m \exp \left\{ -\frac{1}{2} \sum (y_i - x_i^T \beta)^2 \right\} \left( \frac{1}{\sqrt{2\pi} \lambda} \right)^d \exp \left\{ -\frac{1}{2} \sum_{j=1}^d (\beta_j/\lambda)^2 \right\}$$

$$\log\left(\frac{1}{\sqrt{2\pi}}\right)^m - \frac{1}{2} \sum (y_i - x_i^T \beta)^2 + \log\left(\frac{1}{\sqrt{2\pi}\lambda}\right)^d - \frac{1}{2} \sum_{j=1}^d (\beta_j/\lambda)^2$$

max LIK = max (log(LIK)) WHY? LIKELIHOOD POSITIVE

$$= \max \left\{ -\frac{1}{2} \sum (y_i - x_i^T \beta)^2 - \frac{1}{2} \sum_{j=1}^d (\beta_j/\lambda)^2 \right\}$$

$$= \min \left\{ \frac{1}{2} \sum (y_i - x_i^T \beta)^2 + \frac{1}{2} \sum_{j=1}^d (\beta_j/\lambda)^2 \right\}$$

$$= \min \left\{ \frac{1}{2} \sum (y_i - x_i^T \beta)^2 + \frac{1}{2\lambda^2} \beta^T \beta \right\}$$

PENALIZATION

$$\text{FOC} \quad \sum x_i (y_i - x_i^T \beta) + \frac{1}{2\lambda^2} 2 I_d \beta = 0$$

$$\left( \sum x_i x_i^T + \frac{1}{\lambda^2} I_d \right) \beta = \sum x_i y_i$$

$$\hat{\beta}^{\text{BAYES}} = \left( \sum x_i x_i^T + \frac{1}{\lambda^2} I_d \right)^{-1} \left( \sum x_i y_i \right)$$

$$\beta^T \beta = \|\beta\|_{L^2}^2 \quad \text{PENALIZING ITS SQUARE EUCLIDEAN NORM.}$$

$$\frac{x^T A x}{\delta x} = 2 A x$$

$\lambda \downarrow$ , PENALIZATION TERM BIGGER.

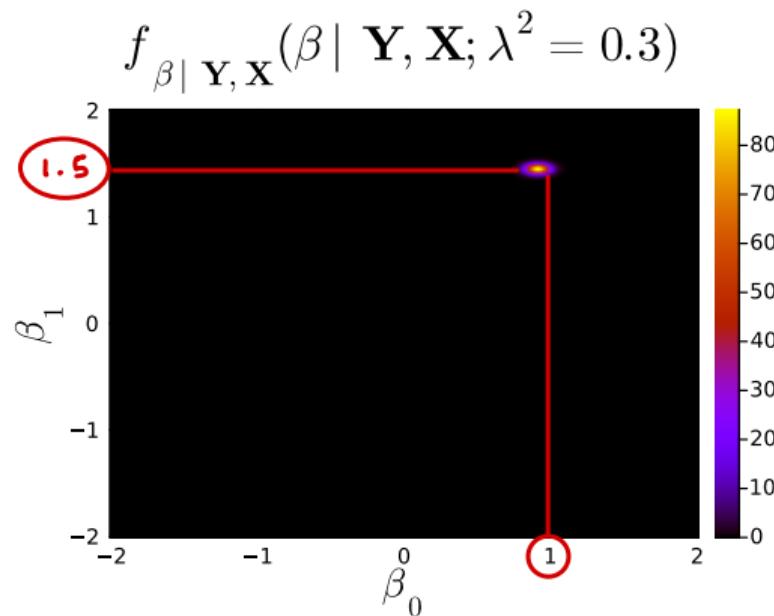
$\lambda \downarrow$ , MORE CONFIDENT ABOUT MY PRIOR WHICH SHRINKS PARAMETERS TO ZERO

↓  
 $\therefore$  I WILL PENALIZE MORE DEVIATIONS FROM ZEROS.

# Ridge Regression - Posterior Distribution

---

The posterior distribution,  $\beta | \mathbf{Y}, \mathbf{X} \sim \mathcal{N}(\dot{m}, \dot{Q})$ , belongs to the same family of probability distributions as the prior when combined with the likelihood function  $\Rightarrow$  the prior and posterior distributions are known as **conjugate distributions**.



# Formalization Bayesian Inference

---

Chain's Rule:

$$f_{\mu|\mathbf{Y}}(\mu, \mathbf{Y}) = f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) f_{\mathbf{Y}}(\mathbf{y})$$

$$f_{\mu|\mathbf{Y}}(\mu, \mathbf{Y}) = f_{\mathbf{Y}|\mu}(\mathbf{Y}|\mu) f_{\mu}(\mu)$$

Bayes' Rule:

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) = \frac{f_{\mathbf{Y}|\mu}(\mathbf{y}|\mu)}{f_{\mathbf{Y}}(\mathbf{y})} f_{\mu}(\mu) \propto f_{\mathbf{Y}|\mu}(\mathbf{y}|\mu) f_{\mu}(\mu)$$

Sample mean case:

- $\{y_1, \dots, y_n\}$  iid sample with  $y_i \sim \mathcal{N}(\mu, \sigma^2)$  and  $\sigma^2$  known.
- $\mu \sim \mathcal{N}(m, Q)$
- $\mu|\mathbf{Y} \sim ?$

# Posterior Distribution $\mu|\mathbf{Y}$

Posterior distribution:

$$f_{\mathbf{Y}|\mu} = f_{Y_1|\mu} \cdot \dots \cdot f_{Y_n|\mu} \text{ IID}$$

$$\mathcal{N}(m, Q)$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(y_i - \mu)^2\right\} \cdot \frac{1}{\sqrt{2\pi Q}} \exp\left\{-\frac{1}{2Q}(\mu - m)^2\right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - \mu)^2\right\} \cdot \frac{1}{\sqrt{2\pi Q}} \exp\left\{-\frac{1}{2Q}(\mu - m)^2\right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \left[ n(\bar{y} - \mu)^2 + \sum_{i=1}^n (y_i - \bar{y})^2 + 2(\bar{y} - \mu) \sum_{i=1}^n (y_i - \bar{y}) \right]\right\} \cdot \frac{1}{\sqrt{2\pi Q}} \exp\left\{-\frac{1}{2Q}(\mu - m)^2\right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{n}{2\sigma^2}(\bar{y} - \mu)^2\right\} \cdot (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \bar{y})^2\right\} \cdot \frac{1}{\sqrt{2\pi Q}} \exp\left\{-\frac{1}{2Q}(\mu - m)^2\right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto \exp\left\{-\frac{n}{2\sigma^2}(\bar{y} - \mu)^2 - \frac{1}{2Q}(\mu - m)^2\right\} = \exp\left\{-\frac{n}{2\sigma^2}(\bar{y}^2 + \mu^2 - 2\bar{y}\mu) - \frac{1}{2Q}(\mu^2 + m^2 - 2\mu m)\right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto \exp\left\{-\frac{1}{2} \left[ \mu^2 \left( \frac{n}{\sigma^2} + \frac{1}{Q} \right) + m^2 \left( \frac{1}{Q} \right) - 2\mu \left( \frac{n}{\sigma^2} \bar{y} + \frac{1}{Q} m \right) \right] \right\} \cdot \exp\left\{-\frac{1}{2} \left( \bar{y}^2 \frac{n}{\sigma^2} \right) \right\}$$

$$f_{\mu|\mathbf{Y}}(\mu|\mathbf{Y}) \propto \exp\left\{-\frac{1}{2Q}(\mu - m)^2\right\} \implies \mu|\mathbf{Y} \sim \mathcal{N}(m, Q)$$

KERNEL

$$\frac{1}{n} \sum y_i = \bar{y}$$

$$\frac{1}{n} \sum (y_i - \bar{y}) = 0$$

$$\therefore \sum (y_i - \bar{y}) = 0$$

$$\sum (y_i - \bar{y}) = 0$$

REMOVE TERMS WITHOUT  $\mu$

Posterior moments:

$$-\frac{1}{2Q} \mu^2 = -\frac{1}{2} \mu^2 \left( \frac{n}{\sigma^2} + \frac{1}{Q} \right) \implies \frac{1}{Q} = -\frac{1}{2} \mu^2 \left( \frac{n}{\sigma^2} + \frac{1}{Q} \right) \implies Q = [(\sigma^2/n)^{-1} + Q^{-1}]^{-1}$$

$$\frac{1}{2Q} 2\mu m = \frac{1}{2} 2\mu \left( \frac{n}{\sigma^2} \bar{y} + \frac{1}{Q} m \right) \implies \frac{m}{Q} = \frac{n}{\sigma^2} \bar{y} + \frac{1}{Q} m \implies m = Q[(\sigma^2/n)^{-1} \bar{y} + Q^{-1} m]$$

Bayesian Inference  $\frac{O}{O + \sigma^2/m} m + \frac{\sigma^2/m}{\sigma^2/m} \bar{y} \therefore \hat{m} = \bar{y} = \hat{y}_{MLE}$

$$\hat{m} = \left( \underbrace{\frac{Q^{-1}}{Q^{-1} + (\sigma^2/n)^{-1}}}_{{(Q^{-1}/Q^{-1} + \infty)} \cdot m = O} \right) m + \left( \underbrace{\frac{(\sigma^2/n)^{-1}}{Q^{-1} + (\sigma^2/n)^{-1}}} \frac{(\sigma^2/m)^{-1}}{(\sigma^2/m)^{-1} \left\{ \frac{Q^{-1}}{(\sigma^2/m)^{-1}} + 1 \right\}} \cdot \bar{y} = \bar{y} \right) \bar{y}$$

What happens when  $n \rightarrow \infty$ ? And when  $Q \rightarrow \infty$ ?  $\hat{m} = \bar{y}$

Under a quadratic loss function, the bayesian estimate of  $\mu$  that minimizes the posterior expected loss is the **mean of the posterior distribution**  $\hat{m}$ :

$$\begin{aligned}
 E_{\mu|Y}[(\mu - \hat{\mu})^2 | Y] &= E_{\mu|Y} L(\mu - \hat{m} + \hat{m} - \hat{\mu})^2 | Y \\
 &= E_{\mu|Y} L(\mu - \hat{m})^2 | Y + E_{\mu|Y} L(\hat{m} - \hat{\mu})^2 | Y + \underline{2(\hat{m} - \hat{\mu}) E_{\mu|Y} L(\mu - \hat{m}) | Y} \\
 &= E_{\mu|Y} L(\mu - \hat{m})^2 | Y + \underbrace{E_{\mu|Y} L(\hat{m} - \hat{\mu})^2 | Y}_{\text{MINIMIZED FOR } \hat{\mu} = \hat{m} = E_{\mu|Y} L(\mu | Y)} \\
 &\quad = 2(\hat{m} - \hat{\mu})(E_{\mu|Y} L(\mu | Y) - \hat{m}) \\
 &\quad = 2(\hat{m} - \hat{\mu})(\hat{m} - \hat{m}) = 0
 \end{aligned}$$

# Link Bayesian and Frequentist Inference

---

**Bernstein-von Mises Theorem:** under some regularity conditions, given  $\tilde{\theta}$  with the posterior distribution, we have:

$$\begin{aligned}\tilde{\theta} &\xrightarrow{P} \hat{\theta}^{MLE} \\ \sqrt{N}(\tilde{\theta} - \hat{\theta}^{MLE}) &\xrightarrow{d} \mathcal{N}(0, Var(\hat{\theta}^{MLE}))\end{aligned}$$

The most important implication of the Bernstein–von Mises theorem is that the Bayesian inference is asymptotically correct from a frequentist point of view.

# Bayesian Linear Regression

---

- Previous result generalizes to linear regression case:  $y_i = \mathbf{x}_i^T \boldsymbol{\beta} + u_i$  with  $u_i \sim \mathcal{N}(0, \sigma^2)$  and  $\sigma^2$  assumed to be known.
- Assume gaussian **prior distribution**:  $f_{\boldsymbol{\beta}}(\boldsymbol{\beta}; \sigma^2) = \mathcal{N}(\boldsymbol{m}, \sigma^2 \mathbf{Q})$ .
- We get **posterior distribution**:  $f_{\boldsymbol{\beta}|\mathbf{Y}, \mathbf{X}}(\boldsymbol{\beta}|\mathbf{Y}, \mathbf{X}; \sigma^2) = \mathcal{N}(\dot{\boldsymbol{m}}, \sigma^2 \dot{\mathbf{Q}})$  where the moments of posterior distribution are:
  - (i)  $\dot{\mathbf{Q}} = (Q^{-1} + \hat{Q}_n^{-1})^{-1}$  : **INVERSE OF THE SUM OF THE PRECISION MATRICES**
  - (ii)  $\dot{\boldsymbol{m}} = \dot{\mathbf{Q}} (Q^{-1} \boldsymbol{m} + \hat{Q}_n^{-1} \hat{\boldsymbol{\beta}}^{OLS})$
  - (iii)  $\hat{Q}_n = (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T)^{-1} \approx \sigma^2 \mathbf{x}' \mathbf{x}$
- Now compare  $\dot{\boldsymbol{m}}$  with the result from the ridge regression exercise.