

# Workshop: Intro to Bayesian Learning

## Lecture 3 - Bayesian Regression and Regularization

Mattias Villani

**Department of Statistics  
Stockholm University**



# Overview

- Bayesian linear regression
- Regularization priors

# Linear regression

## ■ Linear Gaussian regression

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i, \quad \varepsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \sigma^2)$$

## ■ The linear regression model in **matrix form**

$$\underset{(n \times 1)}{\mathbf{y}} = \underset{(n \times k)}{\mathbf{X}} \underset{(k \times 1)}{\boldsymbol{\beta}} + \underset{(n \times 1)}{\boldsymbol{\varepsilon}}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_n^\top \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

## ■ Usually $x_{i1} = 1$ , for all $i$ . $\beta_1$ is the intercept.

## ■ **Likelihood**

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

# Posterior in linear regression - uniform prior

## Gaussian linear regression with non-informative prior

**Model:**  $y = \mathbf{X}\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$

**Prior:**  $p(\beta, \sigma^2) \propto 1/\sigma^2$

**Posterior:**

$$\beta | \sigma^2, \mathbf{y}, \mathbf{X} \sim N(\hat{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$$

$$\sigma^2 | \mathbf{y}, \mathbf{X} \sim \text{Inv-}\chi^2(n - p, s^2)$$

$$\hat{\beta} \equiv (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$s^2 \equiv (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) / (n - p)$$

**Marginal posterior:**

$$\beta | \mathbf{y} \sim t_{n-k} \left( \hat{\beta}, s^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right)$$

# Linear regression - conjugate prior

- **Joint prior** for  $\beta$  and  $\sigma^2$

$$\begin{aligned}\beta|\sigma^2 &\sim N(\mu_0, \sigma^2 \Omega_0^{-1}) \\ \sigma^2 &\sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)\end{aligned}$$

- Common choices:

- ▶  $\Omega_0 = \kappa I_p$  (**Ridge**)
- ▶  $\Omega_0 = \frac{\kappa}{n} \mathbf{X}^\top \mathbf{X}$  (**Zellner's prior**).
- ▶  $\Omega_0 = \frac{1}{n} \mathbf{X}^\top \mathbf{X}$  (Noninformative **Unit information prior**)

# Posterior in linear regression - conjugate prior

## Gaussian linear regression with conjugate prior

**Model:**  $y = \mathbf{X}\beta + \epsilon, \quad \epsilon \stackrel{iid}{\sim} N(0, \sigma^2 I_n)$

**Prior:**

$$\begin{aligned}\beta | \sigma^2 &\sim N(\mu_0, \sigma^2 \Omega_0^{-1}) \\ \sigma^2 &\sim \text{Inv-}\chi^2(\nu_0, \sigma_0^2)\end{aligned}$$

**Posterior:**

$$\begin{aligned}\beta | \sigma^2, \mathbf{y}, \mathbf{X} &\sim N(\mu_n, \sigma^2 \Omega_n^{-1}) \\ \sigma^2 | \mathbf{y}, \mathbf{X} &\sim \text{Inv-}\chi^2(\nu_n, \sigma_n^2)\end{aligned}$$

$$\Omega_n = \mathbf{X}^\top \mathbf{X} + \Omega_0,$$

$$\mu_n = (\Omega_n^{-1} \mathbf{X}^\top \mathbf{X}) \hat{\beta} + \Omega_n^{-1} \Omega_0 \mu_0$$

$$\nu_n = \nu_0 + n$$

**Marginal posterior:**  $\beta | \mathbf{y} \sim t_{\nu_n}(\mu_n, \sigma_n^2 \Omega_n^{-1})$

# Julia code for linear regression - conjugate prior

```
function BayesLinReg(y::Vector, X,  $\mu_o$ ,  $\Omega_o$ ,  $v_o$ ,  $\sigma^2_o$ , nSim)

    # Define ScaledInverseChiSquare distribution
    ScaledInverseChiSq(v, $\tau^2$ ) = InverseGamma(v/2,v* $\tau^2$ /2)

    # Compute posterior hyperparameters
    n = length(y)
    p = size(X,2)
    XX = X'*X
     $\hat{\beta}$  = X \ y
     $\Omega_n$  = Symmetric(XX +  $\Omega_o$ )
     $\mu_n$  =  $\Omega_n \backslash (XX * \hat{\beta} + \Omega_o * \mu_o)$ 
     $v_n$  =  $v_o$  + n
     $\sigma^2_n$  = (  $v_o * \sigma^2_o$  + (y-X* $\hat{\beta}$ )'*(y-X* $\hat{\beta}$ ) + ( $\mu_n - \hat{\beta}$ )'*XX*( $\mu_n - \hat{\beta}$ ) +
    | ( $\mu_n - \mu_o$ )'* $\Omega_o$ *( $\mu_n - \mu_o$ ) )/v_n
    inv $\Omega_n$  = inv( $\Omega_n$ )

    # Sampling from posterior
     $\sigma^2_{sim}$  = zeros(nSim)
     $\beta_{sim}$  = zeros(nSim,p)
    for i ∈ 1:nSim
        # Simulate from p( $\sigma^2$ |y,X)
         $\sigma^2$  = rand(ScaledInverseChiSq( $v_n$ , $\sigma^2_n$ ))
         $\sigma^2_{sim}[i]$  =  $\sigma^2$ 

        # Simulate from p( $\beta$ | $\sigma^2$ ,y,X)
         $\beta$  = rand(MvNormal( $\mu_n$ , $\sigma^2$ *inv $\Omega_n$ ))
         $\beta_{sim}[i,:]$  =  $\beta'$ 
    end

    return  $\mu_n$ ,  $\Omega_n$ ,  $v_n$ ,  $\sigma^2_n$ ,  $\beta_{sim}$ ,  $\sigma^2_{sim}$ 
end
```

end

# R for linear regression - conjugate prior

```
# Function to simulate from the scaled inverse Chi-square distribution
rScaledInvChi2 <- function(n, df, scale){
  return((df*scale)/rchisq(n,df=df))
}

BayesLinReg <- function(y, X, mu_0, Omega_0, v_0, sigma2_0, nIter){

  # Compute posterior hyperparameters
  n = length(y) # Number of observations
  nCovs = dim(X)[2] # Number of covariates
  XX = t(X)%*%X
  betaHat <- solve(XX,t(X)%*%y)
  Omega_n = XX + Omega_0
  mu_n = solve(Omega_n,XX%*%betaHat+Omega_0%*%mu_0)
  v_n = v_0 + n
  sigma2_n = as.numeric((v_0*sigma2_0 + ( t(y)%*%y + t(mu_0)%*%Omega_0%*%mu_0 -
                                          t(mu_n)%*%Omega_n%*%mu_n))/v_n)

  invOmega_n = solve(Omega_n)

  # The actual sampling
  sigma2Sample = rep(NA, nIter)
  betaSample = matrix(NA, nIter, nCovs)
  for (i in 1:nIter){

    # Simulate from p(sigma2 | y, X)
    sigma2 = rScaledInvChi2(n=1, df = v_n, scale = sigma2_n)
    sigma2Sample[i] = sigma2

    # Simulate from p(beta | sigma2, y, X)
    beta_ = rmvnorm(n=1, mean = mu_n, sigma = sigma2*invOmega_n)
    betaSample[i,] = beta_

  }
  return(results = list(sigma2Sample = sigma2Sample, betaSample=betaSample))
}
```



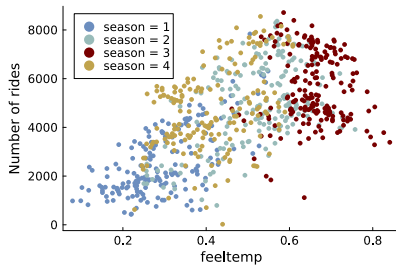
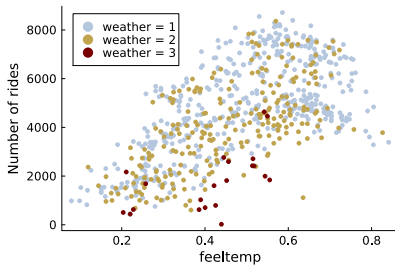
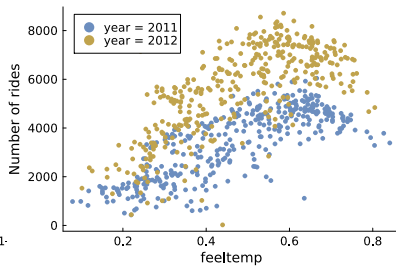
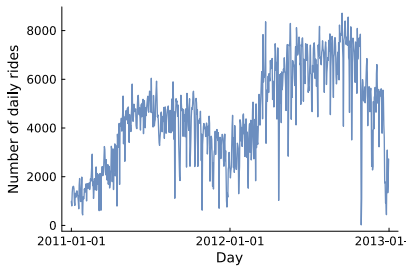
# Bike share data

- **Bike share data.** Predict the number of bike rides.
- Response variable: number of rides on 731 days.

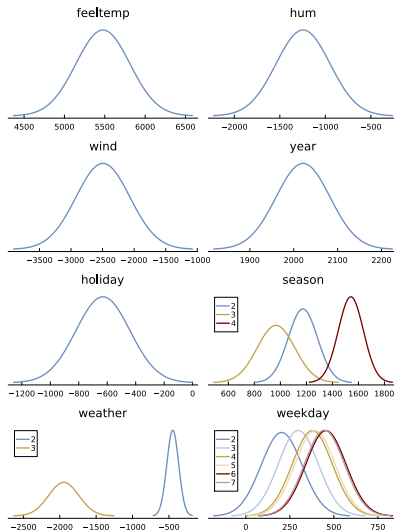
variable	description	type	values	comment
nrides	# of rides	counts	$\{0, 1, \dots\}$	min= 22, max= 8714
feeltemp	perceived temp	cont.	$[0, 1]$	min= 0.07, max= 0.85
hum	humidity	cont.	$[0, 1]$	min= 0.00, max= 0.98
wind	wind speed	cont.	$[0, 1]$	min= 0.02, max= 0.51
year	year	binary	$\{0, 1\}$	year 2011 = 0
season	season	cat.	$\{1, 2, 3, 4\}$	winter $\rightarrow$ fall
weather	weather	ordinal	$\{1, 2, 3\}$	clear $\rightarrow$ rain/snow
weekday	day of week	cat.	$\{0, \dots, 6\}$	sunday $\rightarrow$ saturday
holiday	holiday	binary	$\{0, 1\}$	holiday = 1

- Prior:
  - ▶  $\mu_0 = (1000, 0, \dots, 0)^\top$
  - ▶  $\Omega_0 = \frac{\kappa_0}{n} \mathbf{X}^\top \mathbf{X}$  with  $\kappa_0 = 1$  (unit information prior)
  - ▶  $\sigma_0^2 = 1000^2$  and  $\nu_0 = 5$ .

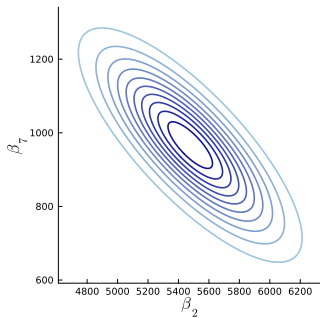
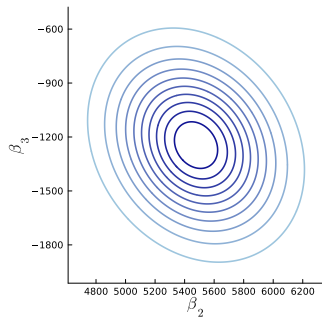
# Bike share data



# Bike share data - marginal posteriors of $\beta$



# Bike share data - joint posteriors of $\beta$



# Interactive - Bayesian regression

prior type: ☒ Zellner g-prior ☐ Identity

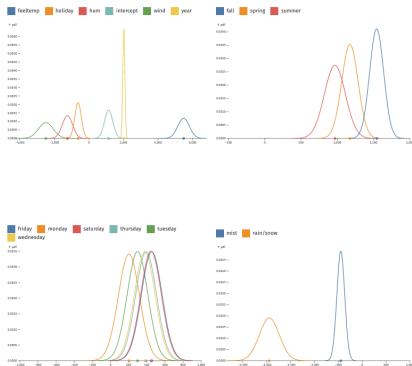
prior sample size,  $\kappa_0$

prior intercept,  $\mu_0$

$\lambda_0$

$\sigma_0$

Marginal posteriors for selected variables. Least squared estimates are indicated by points.



## Ridge regression = iid normal prior

- **Shrinkage/regularization prior** [ $\Omega_0 = \lambda I_p$ ]

$$\beta_i | \lambda, \sigma^2 \stackrel{\text{iid}}{\sim} N\left(0, \frac{\sigma^2}{\lambda}\right)$$

- Posterior mean is the **ridge regression** estimator

$$\mu_n = \left(\mathbf{X}^\top \mathbf{X} + \lambda I\right)^{-1} \mathbf{X}^\top \mathbf{y}$$

- **Shrinkage** toward zero

$$\text{As } \lambda \rightarrow \infty, \mu_n \rightarrow \mathbf{0}$$

- When  $\mathbf{X}^\top \mathbf{X} = I_p$

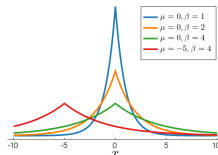
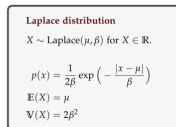
$$\mu_n = (1 - \phi) \hat{\beta}, \quad \text{for } \phi = \frac{\lambda}{1 + \lambda}$$

- **Shrinkage factor**  $\phi \in [0, 1]$ .
- Lecture 5: Bayesian learning of  $\lambda$ .

# Lasso regression = Laplace prior

- **Lasso** is equivalent to posterior mode under Laplace prior

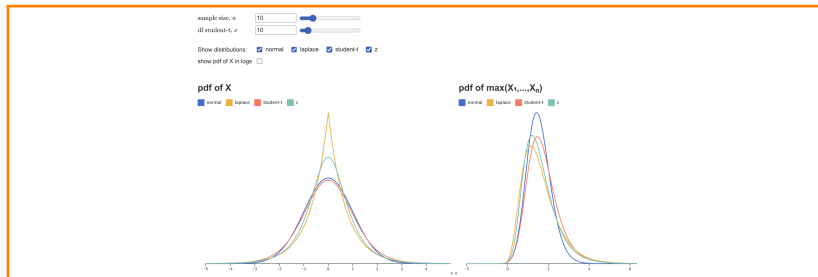
$$\beta_i | \lambda, \sigma^2 \stackrel{\text{iid}}{\sim} \text{Laplace} \left( 0, \frac{\sigma^2}{\lambda} \right)$$



- **Laplace prior:**
  - ▶ heavy tails
  - ▶ many  $\beta_i$  close to zero, but some  $\beta_i$  can be very large.
- **Normal prior:**
  - ▶ light tails
  - ▶ all  $\beta_i$ 's are similar in magnitude and no  $\beta_i$  very large.



# Interactive - tails of distributions





# Horseshoe prior

- Normal and Laplace - one **global shrinkage** parameter  $\lambda$ .
- **Global-Local shrinkage**: global + local shrinkage for each  $\beta_j$ .
- **Horseshoe prior**:

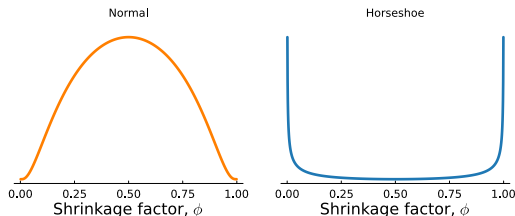
$$\beta_j | \lambda_j^2, \tau^2 \stackrel{\text{ind}}{\sim} N(0, \sigma^2 \tau^2 \lambda_j^2)$$

$$\lambda_j \sim C^+(0, 1)$$

$$\tau \sim C^+(0, 1)$$

- The posterior mean for  $\beta$  satisfies approximately

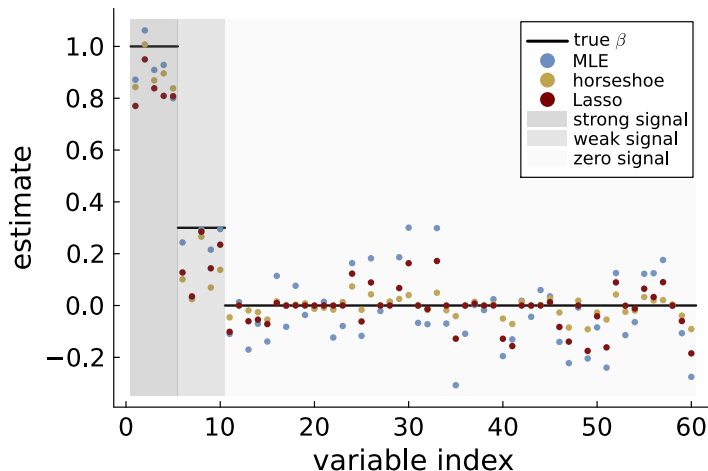
$$\mu_{n,j} \approx (1 - \phi_j) \hat{\beta}_j$$



# Simulated example: Lasso vs Horseshoe

■ Linear regression with  $p = 60$  uncorrelated covariates:

- ▶ 5 strong signals ( $t$ -ratio: 10)
- ▶ 5 mild signals ( $t$ -ratio: 3)
- ▶ 50 noise signals ( $t$ -ratio: 0)



# Variable selection by spike-and-slab prior

## ■ Spike-and-slab prior

$$\beta_j | \sigma^2, \tau^2, z_j \sim \begin{cases} 0 & \text{if } z_j = 0 \\ N(0, \sigma^2 \tau^2) & \text{if } z_j = 1 \end{cases}$$

## ■ Prior for the **variable selection indicators**

$$z_j \stackrel{iid}{\sim} \text{Bernoulli}(\omega)$$

## ■ This is a **mixture prior** for the $\beta_j$

$$p(\beta_j) = (1 - \omega)\delta_0(\beta_j) + \omega N(\beta_j | 0, \sigma^2 \tau^2)$$

