

STA 602 - Intro to Bayesian Statistics

Lecture 9

Li Ma

Duke University

Gaussian model with unknown mean and unknown variance

- ▶ Sampling model for n readings given the mean θ and variance σ^2 is

$$X_1, X_2, \dots, X_n \mid \theta, \sigma^2 \stackrel{\text{iid}}{\sim} \text{N}(\theta, \sigma^2)$$

- ▶ Prior distribution for the mean μ

$$\theta \mid \sigma^2 \sim \text{N}(\mu_0, \tau_0^2)$$

- ▶ Now let's assume that the sampling variance σ^2 is also unknown.
- ▶ So now we have a two parameter model (θ, σ^2) .
- ▶ To complete Bayesian inference, we need to specify a joint prior for (θ, σ^2) .

A conjugate prior

- If we let $\tau_0^2 = \sigma^2 / \kappa_0$,

$$\theta \mid \sigma^2 \sim \mathcal{N}(\mu_0, \sigma^2 / \kappa_0),$$

then we have found the conditional posterior

$$\theta \mid \mathbf{x}, \sigma^2 \sim \mathcal{N}(\mu_n, \tau_n^2)$$

where

$$\mu_n = \frac{\kappa_0}{\kappa_n} \mu_0 + \frac{n}{\kappa_n} \bar{x}$$

where $\kappa_n = \kappa_0 + n$ and

$$\tau_n^2 = \sigma^2 / \kappa_n.$$

A conjugate prior

- ▶ There exists a marginal prior on σ^2 as well to make the joint prior on (θ, σ^2) fully conjugate.

$$\frac{1}{\sigma^2} \sim \text{Gamma}(v_0/2, v_0\sigma_0^2/2).$$

- ▶ Equivalently sometimes we say that σ^2 follows an *inverse-Gamma* (IG) prior

$$\sigma^2 \sim \text{IG}(v_0/2, v_0\sigma_0^2/2).$$

An equivalent parametrization

- ▶ Recall that χ_v^2 distribution (i.e., with v degrees of freedom) is $\text{Gamma}(v/2, 1/2)$.
- ▶ So $\text{Gamma}(v_0/2, v_0\sigma_0^2/2)$ is a scaled $\chi_{v_0}^2$ distribution

$$\text{Gamma}(v_0/2, v_0\sigma_0^2/2) =_d \frac{\chi_{v_0}^2}{v_0\sigma_0^2} = \frac{\chi_{v_0}^2}{v_0} \cdot \frac{1}{\sigma_0^2}.$$

- ▶ It has mean $1/\sigma_0^2$, and its variance is smaller as the degrees of freedom v_0 increases. (What happens when $v_0 \uparrow \infty$?)
- ▶ v_0 is a *prior degrees of freedom (d.f.)*, which quantifies the strength of our prior knowledge for the variance.

The marginal posterior of σ^2

- ▶ One can find the marginal posterior of σ^2 to be

$$\frac{1}{\sigma^2} | \mathbf{x} \sim \text{Gamma}(v_n/2, v_n \sigma_n^2/2)$$

or

$$\sigma^2 | \mathbf{x} \sim \text{IG}(v_n/2, v_n \sigma_n^2/2)$$

where $v_n = v_0 + n$ (the posterior d.f.)

$$\begin{aligned}\sigma_n^2 &= \frac{1}{v_n} \left[v_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{x} - \mu_0)^2 \right] \\ &= \frac{1}{v_n} \left[v_0 \sigma_0^2 + \sum_{i=1}^n (x_i - \bar{x})^2 + \frac{\kappa_0}{\kappa_n} \cdot n (\bar{x} - \mu_0)^2 \right]\end{aligned}$$

- ▶ The posterior mean of σ^2 has contribution from three pieces
 - ▶ prior mean σ_0^2 (with weight proportional to prior d.f.)
 - ▶ sample variance $s^2 = \sum_i (x_i - \bar{x})^2 / (n-1)$ (with weight proportional to $n-1$)
 - ▶ the deviation of sample mean \bar{x} from its prior mean μ_0 . (Because $\tau_0^2 = \sigma^2 / \kappa_0$, σ^2 is tied to the variability of θ .)

Deriving the marginal posterior of σ^2 (or $1/\sigma^2$)

- ▶ There are several ways to derive the marginal posterior.
- ▶ The most basic way is by marginalizing out θ in the joint posterior of $(\theta, 1/\sigma^2)$.
- ▶ For notational simplicity, let $\gamma = 1/\sigma^2$.
- ▶ By Bayes theorem,

$$\begin{aligned} p(\theta, \gamma | \mathbf{x}) &\propto p(\theta | \gamma) p(\gamma) p(\mathbf{x} | \theta, \gamma) \\ &\propto \gamma^{\frac{1}{2}} e^{-\frac{1}{2} \kappa_0 \gamma (\theta - \mu_0)^2} \cdot \gamma^{\frac{\nu_0}{2} - 1} e^{-\frac{\nu_0 \sigma_0^2}{2} \gamma} \cdot \gamma^{\frac{n}{2}} e^{-\frac{1}{2} \gamma \sum_i (x_i - \theta)^2} \end{aligned}$$

Note that the constant must not involve either θ or γ .

- ▶ Use the fact that

$$\sum_i (x_i - \theta)^2 = \sum_i (x_i - \bar{x})^2 + n(\bar{x} - \theta)^2 = (n-1)s^2 + n(\bar{x} - \theta)^2$$

- We have

$$\begin{aligned}
 p(\theta, \gamma | \mathbf{x}) &\propto \gamma^{\frac{1}{2}} e^{-\frac{1}{2} \kappa_0 \gamma (\theta - \mu_0)^2} \cdot \gamma^{\frac{\nu_0}{2} - 1} e^{-\frac{\nu_0 \sigma_0^2}{2} \gamma} \cdot \gamma^{\frac{n}{2}} e^{-\frac{1}{2} \gamma \cdot (n-1) s^2} \cdot e^{-\frac{1}{2} \gamma \cdot n (\bar{x} - \theta)^2} \\
 &\propto \underbrace{\gamma^{\frac{1}{2}} e^{-\frac{1}{2} [\kappa_0 \gamma (\theta - \mu_0)^2 + n \gamma (\bar{x} - \theta)^2]}}_{\text{Involves } \theta \text{ and } \gamma} \cdot \underbrace{\gamma^{\frac{\nu_0 + n}{2} - 1} e^{-\frac{\gamma}{2} [\nu_0 \sigma_0^2 + (n-1) s^2]}}_{\text{Involves only } \gamma}
 \end{aligned}$$

- Now let's focus on the first term. The idea is to turn it into the form of

$$\gamma^{1/2} e^{-\frac{(\theta - A)^2}{2} \cdot B \gamma}$$

through completion of squares, because

$$\int \gamma^{1/2} e^{-\frac{(\theta - A)^2}{2} \cdot B \gamma} d\theta = \sqrt{2\pi/B} \quad \text{a constant not involving } \theta \text{ and } \gamma$$

from the fact that

$$\int \left(\frac{B\gamma}{2\pi} \right)^{1/2} e^{-\frac{1}{2} B \gamma (\theta - A)^2} = 1$$

Finding the values of A and B (which we have already done before!)

- ▶ Let's complete the square for the exponent in the first term

$$\begin{aligned} & \kappa_0 \gamma (\theta - \mu_0)^2 + n \gamma (\bar{x} - \theta)^2 \\ &= (\kappa_0 + n) \gamma \left[\left(\theta - \frac{\kappa_0 \mu_0 + n \bar{x}}{\kappa_0 + n} \right)^2 - \left(\frac{\kappa_0 \mu_0 + n \bar{x}}{\kappa_0 + n} \right)^2 + \frac{\kappa_0 \mu_0^2 + n \bar{x}^2}{\kappa_0 + n} \right] \\ &= (\kappa_0 + n) \gamma \left[\left(\theta - \frac{\kappa_0 \mu_0 + n \bar{x}}{\kappa_0 + n} \right)^2 + \frac{\kappa_0 n}{(\kappa_0 + n)^2} (\bar{x} - \mu_0)^2 \right] \\ &= \kappa_n \gamma \left[(\theta - \mu_n)^2 + \frac{\kappa_0 n}{\kappa_n^2} (\bar{x} - \mu_0)^2 \right]. \end{aligned}$$

- ▶ Remark: Previously when σ^2 is known, we have treated the second term, which does not involve θ as a “constant”. Note that here we cannot do that because σ^2 is also an unknown parameter!

- Therefore

$$p(\theta, \gamma | \mathbf{x}) \propto \underbrace{\gamma^{\frac{1}{2}} e^{-\frac{1}{2} B \gamma (\theta - A)^2}}_{\text{Involves } \theta \text{ and } \gamma} \cdot \underbrace{\gamma^{\frac{\nu_0 + n}{2} - 1} e^{-\frac{\gamma}{2} [\nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{x} - \mu_0)^2]}}_{\text{Involves only } \gamma}$$

where

$$A = \mu_n = \frac{\kappa_0 \mu_0 + n \bar{x}}{\kappa_0 + n} \quad \text{and} \quad B = \kappa_n = \kappa_0 + n.$$

- *Caveat*: Just because the first term looks like a Gaussian density in θ , we still need to keep track of the γ in its normalizing constant!

- Moreover,

$$\begin{aligned}
 p(\gamma|\mathbf{x}) &= \int p(\theta, \gamma|\mathbf{x}) d\theta \\
 &\propto \sqrt{2\pi/B} \cdot \gamma^{\frac{v_0+n}{2}-1} e^{-\frac{\gamma}{2}[v_0\sigma_0^2+(n-1)s^2+\frac{\kappa_0 n}{\kappa_n}(\bar{x}-\mu_0)^2]} \\
 &\propto \gamma^{\frac{v_0+n}{2}-1} e^{-\frac{\gamma}{2}[v_0\sigma_0^2+(n-1)s^2+\frac{\kappa_0 n}{\kappa_n}(\bar{x}-\mu_0)^2]} \\
 &= \gamma^{\frac{v_n}{2}-1} e^{-\frac{\gamma v_n \sigma_n^2}{2}},
 \end{aligned}$$

which is a Gamma($v_n/2, v_n \sigma_n^2/2$) where $v_n = v_0 + n$.

- Note that we never needed to know the value of A and B !

The marginal posterior of θ

- ▶ One can also find the marginal posterior of θ , $p(\theta | \mathbf{x})$, by integrating out σ^2 (or γ) in the joint posterior. It is a (scaled and shifted) t -distribution with v_n degrees of freedom.
- ▶ Specifically,

$$\frac{\theta - \mu_n}{\sigma_n / \sqrt{\kappa_n}} \sim t_{v_n}.$$

- ▶ This can be done directly through integrating out γ .
- ▶ This time, we will need the values of A and B !

- To see this, note that by plugging in the values of A and B ,

$$p(\theta, \gamma | \mathbf{x}) = p(\theta | \gamma, \mathbf{x})p(\gamma | \mathbf{x}) \propto \gamma^{\frac{v_n-1}{2}} e^{-\frac{1}{2}\gamma[\kappa_n(\theta-\mu_n)^2 + v_n\sigma_n^2]}.$$

- Hence

$$\begin{aligned} p(\theta | \mathbf{x}) &= \int p(\theta, \gamma | \mathbf{x}) d\gamma \propto \int \gamma^{\frac{v_n-1}{2}} e^{-\frac{1}{2}\gamma[\kappa_n(\theta-\mu_n)^2 + v_n\sigma_n^2]} d\gamma \\ &\propto \frac{\Gamma\left(\frac{v_n+1}{2}\right)}{\left\{\frac{1}{2}[\kappa_n(\theta-\mu_n)^2 + v_n\sigma_n^2]\right\}^{\frac{v_n+1}{2}}} \\ &\propto \left[1 + \left(\frac{\theta - \mu_n}{\sigma_n/\sqrt{\kappa_n}}\right)^2 \cdot \frac{1}{v_n}\right]^{-\frac{v_n+1}{2}}. \end{aligned}$$

- Recall that the density of the t_v distribution is

$$p(t) \propto (1 + t^2/v)^{-\frac{v+1}{2}}$$

- So by change of variable, the posterior density of $t = \frac{\theta - \mu_n}{\sigma_n/\sqrt{\kappa_n}}$ is exactly proportional to the density of a t_n -distribution. Thus

$$\frac{\theta - \mu_n}{\sigma_n/\sqrt{\kappa_n}} \sim t_{v_n}.$$

t -distribution as a *scale-mixture* of normals

- ▶ More generally, a scale-mixture of Gaussian with IG prior on the variance leads to a t -distribution. That is, if

$$\begin{aligned} Y | b &\sim N(a, b^2) \\ b^2 &\sim \text{IG}(\nu/2, c^2 \nu/2) \end{aligned}$$

Then

$$\frac{Y - a}{c} \sim t_\nu.$$

- ▶ We just proved it on the previous slide! (Simply repeat by replacing μ_n with a , $1/\gamma$ with b^2 , ν_n with ν , and c^2 with σ_n^2 .)
- ▶ In particular, a t -distribution with ν d.f. is exactly the marginal distribution of a $\text{Normal}(0, \sigma^2)$ and a $\text{Gamma}(\nu/2, \nu/2)$ on $1/\sigma^2$.

$$\begin{aligned} t_\nu(\cdot) &= \int N(\cdot | 0, \sigma^2) \times \text{IG}(\sigma^2 | \nu/2, \nu/2) d\sigma^2 \\ &= \int N(\cdot | 0, 1/\gamma) \times \text{Gamma}(\gamma | \nu/2, \nu/2) d\gamma. \end{aligned}$$

Going the other way around knowing this property of t distributions

- Applying this to our posterior

$$\begin{aligned} p(\theta|\mathbf{x}) &= \int \underbrace{p(\theta|\gamma, \mathbf{x})}_{\text{N}(\mu_n, \tau_n^2 = \frac{1}{\kappa_n \gamma})} \underbrace{p(\gamma|\mathbf{x})}_{\text{Gamma}(v_n/2, v_n \sigma_n^2/2)} d\gamma \\ &= \int \underbrace{p(\theta|\tau_n^2, \mathbf{x})}_{\text{N}(\mu_n, \tau_n^2)} \underbrace{p(\tau_n^2|\mathbf{x})}_{\text{IG}\left(\frac{v_n}{2}, \frac{v_n \sigma_n^2 + \kappa_n}{2}\right)} d(\tau_n^2). \end{aligned}$$

- By setting $a = \mu_n$, $b = \tau_n$, and $c = \sigma_n / \sqrt{\kappa_n}$, we have

$$\frac{\theta - a}{c} = \frac{\theta - \mu_n}{\sigma_n / \sqrt{\kappa_n}} \sim t_{v_n}.$$

Monte Carlo for t -distributions as a scale-mixture of Gaussian

- ▶ We can use this fact to draw Monte Carlo samples from a t -distribution with a shift a and a scaled parameter c ?
- ▶ For $i = 1, 2, \dots, S$,
 - ▶ First draw $\sigma^{2(i)} \sim \text{IG}(\nu/2, \nu/2)$.
 - ▶ Then draw $t^{(i)} \mid \sigma^{(i)} \sim \text{N}(0, \sigma^{2(i)})$ and set $\theta^{(i)} = a + ct^{(i)}$.
- ▶ The last step is equivalent to drawing $\theta^{(i)} \mid \sigma^{(i)} \sim \text{N}(a, c^2 \sigma^{2(i)})$.
- ▶ Online animation for the special case $a = 0$ and $c = 1$:
<http://www.sumsar.net/blog/2013/12/t-as-a-mixture-of-normals/>

Predictive distribution

- ▶ One can also show that the predictive distribution for a new observation x_{new} is also a (scaled and shifted) t distribution with v_n degrees of freedom.
- ▶ Hint:

$$\begin{aligned} p(x_{new}|\mathbf{x}) &= \int \int p(x_{new}|\boldsymbol{\theta}, \gamma) p(\boldsymbol{\theta}, \gamma|\mathbf{x}) d\boldsymbol{\mu} d\gamma \\ &= \int \underbrace{\int p(x_{new}|\boldsymbol{\theta}, \gamma) p(\boldsymbol{\theta}|\gamma, \mathbf{x}) d\boldsymbol{\theta}}_{N(\boldsymbol{\mu}_n, \boldsymbol{\tau}_n^2 + \boldsymbol{\sigma}^2 = (1/\boldsymbol{\kappa}_n + 1)/\gamma)} \underbrace{p(\gamma|\mathbf{x})}_{\text{Gamma}(v_n/2, v_n \boldsymbol{\sigma}_n^2/2)} d\gamma \end{aligned}$$

and therefore

$$\frac{x_{new} - \boldsymbol{\mu}_n}{\boldsymbol{\sigma}_n \sqrt{1 + 1/\boldsymbol{\kappa}_n}} \sim t_{v_n}.$$

Monte Carlo sampling from the posterior

- ▶ We want to draw S samples from the joint posterior (θ, σ^2) given the data \mathbf{x} .

$$(\theta^{(1)}, \sigma^{2(1)}), (\theta^{(2)}, \sigma^{2(2)}), \dots, (\theta^{(S)}, \sigma^{2(S)}).$$

- ▶ For $i = 1, 2, \dots, S$, we proceed in two steps
 - ▶ Draw a sample of $\sigma^{2(i)}$ or equivalently $\gamma^{(i)} = 1/\sigma^{2(i)}$ from the marginal posterior $p(\sigma^2|\mathbf{x})$ or $p(\gamma|\mathbf{x})$.
 - ▶ Then draw $\theta^{(i)}$ from the conditional distribution of $p(\theta|\sigma^{2(i)}, \mathbf{x})$.
- ▶ That is, for $i = 1, 2, \dots, S$,

$$\begin{aligned}\sigma^{2(i)} &\sim \text{IG}(v_n, v_n \sigma_n^2) \\ \theta^{(i)} | \sigma^{2(i)} &\sim \text{N}(\mu_n, \sigma^{2(i)} / \kappa_n).\end{aligned}$$

Example: Air pollutant

- ▶ Suppose you decided to adopt the above conjugate prior.
- ▶ Your prior belief on the mean θ is equivalent to about $\kappa_0 = 5$ measurements. That is, $\tau_0 = \sigma^2/5$.
- ▶ Regarding the variance σ^2 of the device, however, suppose the device, you know that the precision (i.e., $1/\sigma^2$) is about $1/4$ and so $\sigma_0^2 = 4$ but you are unsure about the quality of the device, so have large uncertainty in its precision. For example, $\nu_0 = 1$.
- ▶ Suppose you recorded 10 readings

$$\mathbf{x} = (104, 105, 103, 102, 105, 107, 106, 104, 103, 106)$$

```

x <- c(104,105,103,102,105,107,106,104,103,106) # the data
n <- length(x) # sample size

mu.0 <- 100 # prior mean for mu
kappa.0 <- 5 # prior sample size
nu.0 <- 1 # prior degrees of freedom
sigma2.0 <- 4
s2 <- var(x) # sample variance of data

kappa.n <- kappa.0 + n
mu.n <- (kappa.0*mu.0 + sum(x))/kappa.n
nu.n <- nu.0 + n
sigma2.n <- (nu.0*sigma2.0 + (n-1)*s2 +
  kappa.0*n/kappa.n*(mean(x)-mu.0)^2)/nu.n

```

```
print(mean(x))
```

```
## [1] 104.5
```

```
print(mu.n)
```

```
## [1] 103
```

```
print(kappa.n)
```

```
## [1] 15
```

```
print(nu.n)
```

```
## [1] 11
```

```
print(sigma2.n)
```

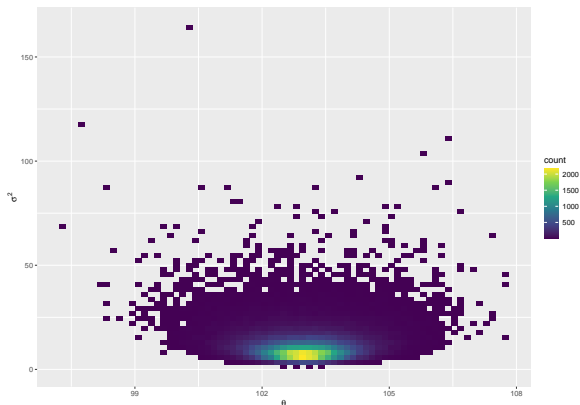
```
## [1] 8.545455
```

Draw Monte Carlo samples

```
S=100000  
sigma2.mc <- 1/rgamma(S, shape=nu.n/2, rate=nu.n*sigma2.n/2)  
theta.mc <- rnorm(S, mean=mu.n, sd=sqrt(sigma2.mc/kappa.n))
```

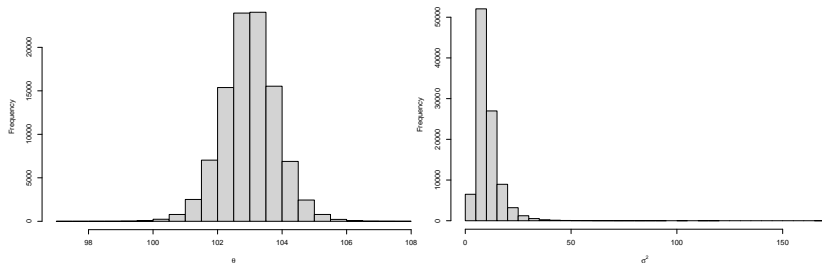
2D histogram of joint samples

```
post.samples <- data.frame(theta=theta.mc,  
                             sigma2=sigma2.mc)  
ggplot(post.samples, aes(x=theta, y=sigma2) ) +  
  labs(x=expression(theta), y=expression(sigma^2)) +  
  geom_bin2d(bins=70) +  
  scale_fill_continuous(type = "viridis")
```



Marginal histograms

```
hist(theta.mc,breaks=20,main="",xlab=expression(theta))  
hist(sigma2.mc, breaks=50, main="",xlab=expression(sigma^2))
```



With a stronger prior on σ^2 (e.g., $\nu_0 = 100$)

```
nu.0 <- 100      # prior degrees of freedom
nu.n <- nu.0 + n
sigma2.n <- (nu.0*sigma2.0 + (n-1)*s2 +
            kappa.0*n/kappa.n*(mean(x)-mu.0)^2)/nu.n
S=100000
sigma2.mc <- 1/rgamma(S,shape=nu.n/2,rate=nu.n*sigma2.n/2)
theta.mc <- rnorm(S,mean=mu.n,sd=sqrt(sigma2.mc/kappa.n))
```

Posterior parameters

```
print(mean(x))
```

```
## [1] 104.5
```

```
print(mu.n)
```

```
## [1] 103
```

```
print(kappa.n)
```

```
## [1] 15
```

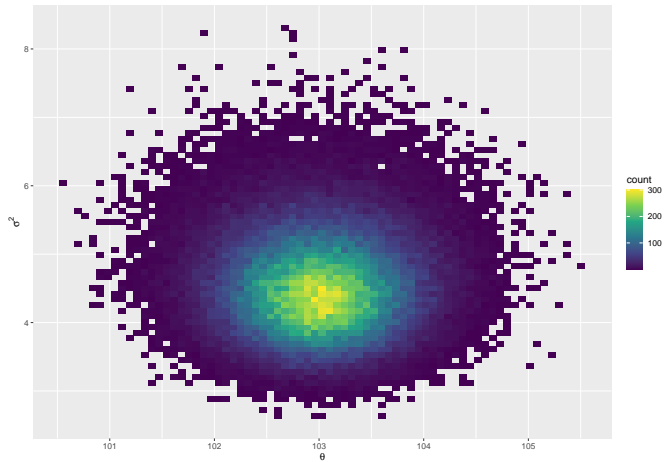
```
print(nu.n)
```

```
## [1] 110
```

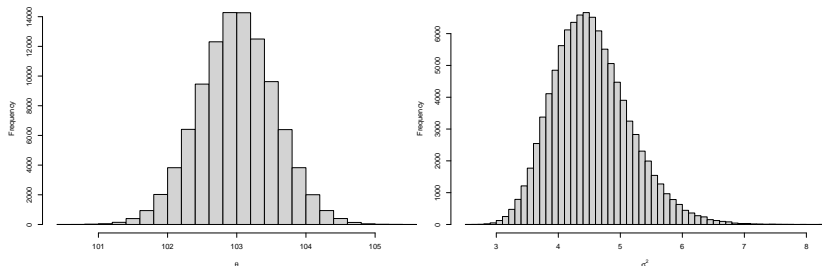
```
print(sigma2.n)
```

```
## [1] 4.454545
```

Joint histogram



Marginal histograms



- ▶ Notice the change in spread of the marginal posterior of θ .
- ▶ Notice the change in spread and shape for σ^2 .

A non-informative prior

- ▶ Suppose we have little prior knowledge about θ and σ^2 .
- ▶ Use a vague prior by letting $\kappa_0 \rightarrow 0$ and $\nu_0 \rightarrow 0$.

- ▶ $\kappa_0 \rightarrow 0$ leads to

$$p(\theta|\sigma^2) \propto 1.$$

- ▶ $\nu_0 \rightarrow 0$ leads to

$$p(\sigma^2) \propto 1/\sigma^2.$$

- ▶ Together, we have

$$p(\theta, \sigma^2) \propto 1/\sigma^2.$$

This is an *improper* prior as it integrates to infinity over $-\infty < \theta < \infty$ and $\sigma^2 > 0$.

- ▶ Nevertheless, this prior leads to a proper posterior.

The corresponding posterior

- The posterior parameters are

$$\kappa_n = \kappa_0 + n = n$$

$$\mu_n = \frac{\kappa_0}{\kappa_n} \mu_0 + \frac{n}{\kappa_n} \bar{x} = \bar{x}$$

$$\tau_n^2 = \sigma^2 / \kappa_n = \sigma^2 / n$$

$$\nu_n = \nu_0 + n = n$$

$$\begin{aligned} \sigma_n^2 &= \frac{1}{\nu_n} \left[\nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{x} - \mu_0)^2 \right] \\ &= \frac{(n-1)s^2}{n} = \frac{1}{n} \sum_i (x_i - \bar{x})^2. \end{aligned}$$

- Hence we have

$$\sigma^2 \mid \mathbf{x} \sim \text{IG}(n/2, (n-1)s^2/2)$$

$$\theta \mid \sigma^2, \mathbf{x} \sim \text{N}(\bar{x}, \sigma^2/n).$$

Example: Air pollutant

```
print(mean(x))
```

```
## [1] 104.5
```

```
print(mu.n)
```

```
## [1] 104.5
```

```
print(kappa.n)
```

```
## [1] 10
```

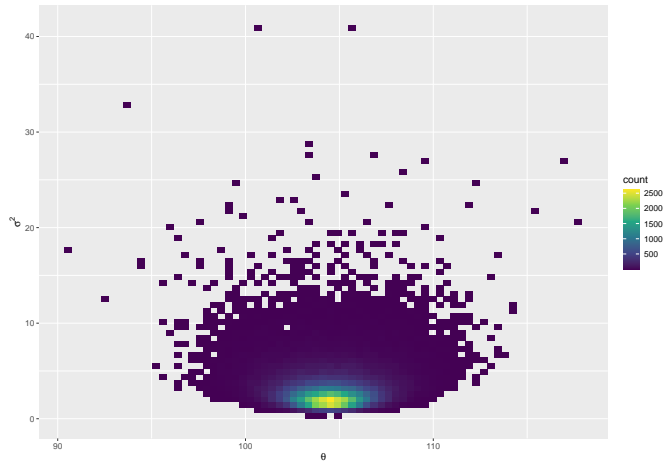
```
print(nu.n)
```

```
## [1] 10
```

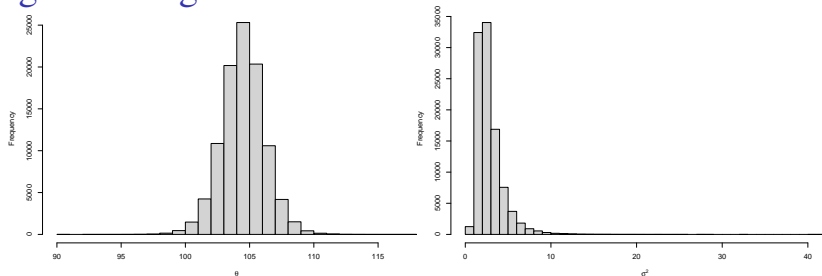
```
print(sigma2.n)
```

```
## [1] 2.25
```

Joint histogram



Marginal histograms



- ▶ Notice the change in spread of the marginal posterior of θ .
- ▶ Notice the change in spread and shape for σ^2 .
- ▶ Question: Why is the marginal posterior of σ^2 less spread out than the first case with $v_0 = 1$ and $\kappa_0 = 5$?
- ▶ Hint: Note the third term that contributes to σ_n^2 .

$$\sigma_n^2 = \frac{1}{v_n} \left[v_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_n} (\bar{x} - \mu_0)^2 \right].$$