

# STAT 224 Lecture/Activities on Friday, 10/01

Your Name Goes Here

## 1. The NC Births Data

The NC Births data came from a random sample of 1000 birth records released by the State of North Carolina in 2004. The variables include:

- **weight**: weight of the baby at birth in pounds.
- **gender**: gender of the baby, female or male.
- **habit** : status of the mother as a nonsmoker or a smoker.
- **marital**: whether mother is married or not married at birth.
- **whitemom**: whether mom is white or not white.
- **fage**: father's age in years.
- **mage**: mother's age in years.
- **gained**: weight gained by mother during pregnancy in pounds.

See <http://www.stat.uchicago.edu/~yibi/s220/labs/lab08.html> for the complete description of variables

## 2. Loading Data

You can download the data at

<https://www.openintro.org/stat/data/csv/nbirths.csv>

Please save it at **the same folder** as this RMD file.

Then you can easily **change the working directory** by

[Session] – [Set Working Directory] – [To Source File Location]

Run the command below to load the data

```
nc = read.csv("ncbirths.csv")
```

## 3. Data Summary by Group

```
library(mosaic)
favstats(weight ~ gender + habit, data=nc)
```

```
##      gender.habit  min   Q1 median   Q3   max     mean      sd   n missing
## 1 female.nonsmoker 1.00 6.31   7.13 7.81 11.63 6.939393 1.511282 445      0
## 2  male.nonsmoker 1.38 6.69   7.50 8.38 11.75 7.357290 1.498625 428      0
## 3  female.smoker  2.19 6.00   6.88 7.44  8.38 6.675263 1.078271  57      0
## 4   male.smoker  1.69 6.25   7.31 8.13  9.19 6.955507 1.593304  69      0
```

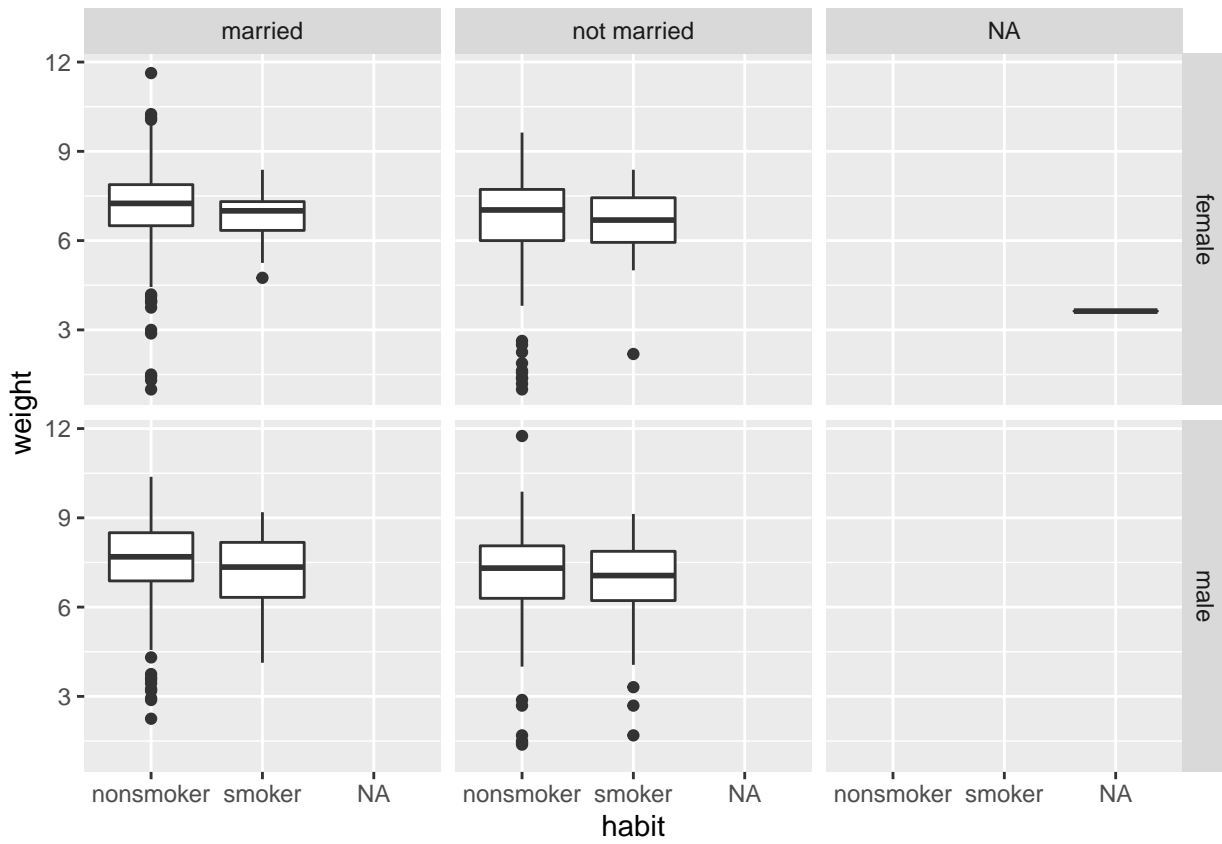
- Describe the effect of **gender** on **weight** controlling for mom's smoking status (**habit**)
- Describe the effect of the mom's smoking **habit** on **weight** controlling for baby's gender

### On Your Own

**Q1:** Find appropriate data summary and use it to describe the effect of mother's marital status (**marital**) on **weight** controlling for the **gender** of baby

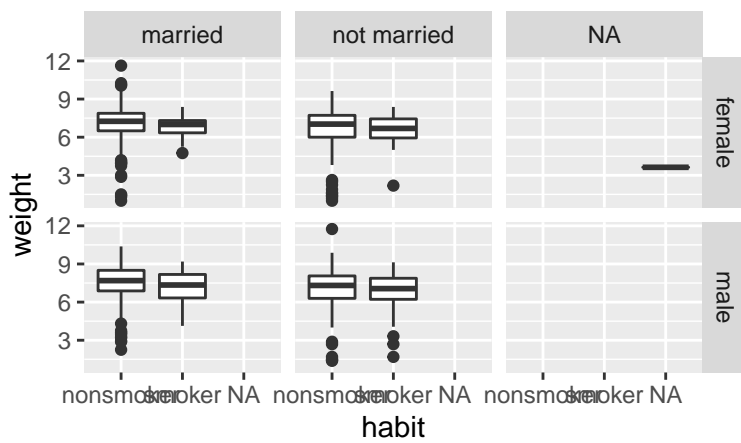
## 4. Side-By-Side Boxplots w/ Facet

```
ggplot(nc, aes(x=habit, y=weight)) +
  geom_boxplot() +
  facet_grid(gender ~ marital)
```



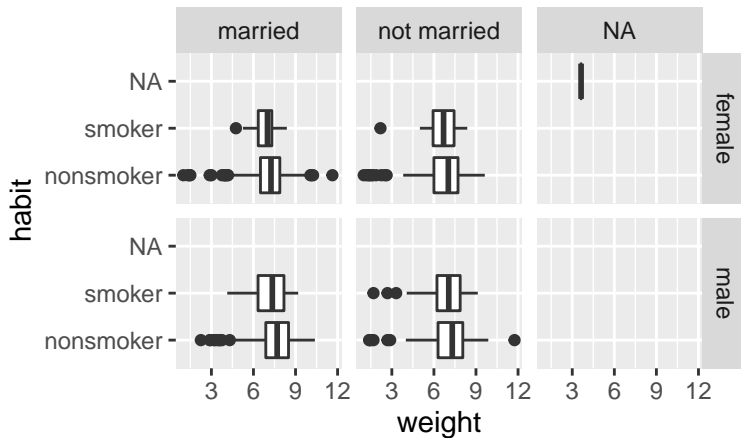
Same Plot Resized:

```
ggplot(nc, aes(x=habit, y=weight)) +
  geom_boxplot() +
  facet_grid(gender ~ marital)
```



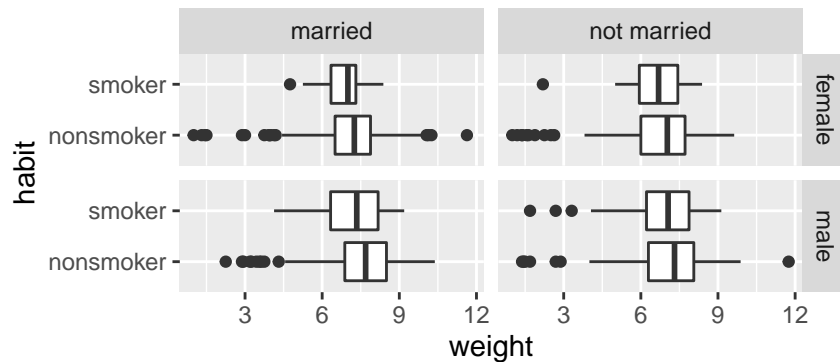
## Same Boxplots Turned Horizontal

```
ggplot(nc, aes(y=habit, x=weight)) +
  geom_boxplot() +
  facet_grid(gender ~ marital)
```



## Same Plot w/ Missing Values Removed

```
ggplot(subset(nc, !is.na(habit)), aes(y=habit, x=weight)) +
  geom_boxplot() +
  facet_grid(gender ~ marital)
```



Based on the plot, describe the effect of smoking on baby's birth weight, controlling for **gender** and mom's marital status.

Answer #1: For babies of the same gender and same marital status of their moms, babies of smoking moms have a lower median birth weight than those of nonsmoking moms.

Answer #2: Babies of smoking moms have a lower median birth weight than those of nonsmoking moms, controlling/adjusted for the gender of the baby and the mom's marital status

## On Your Own

Q2: Find appropriate boxplots of the data and use them to describe the effect of mother's marital status (`marital`) on weight controlling for the gender of baby and mother's smoking habit.

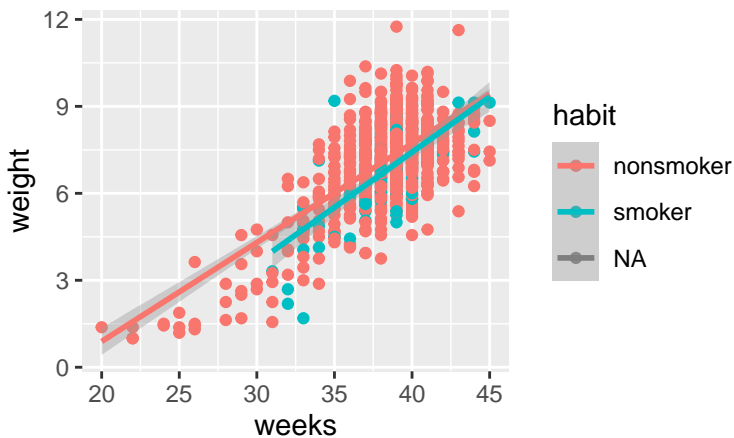
## 5. Controlling for a Numerical Predictor

```
ggplot(nc, aes(x=weeks, y=weight, color=habit)) +  
  geom_point() +  
  geom_smooth(method='lm')
```

```
## 'geom_smooth()' using formula 'y ~ x'
```

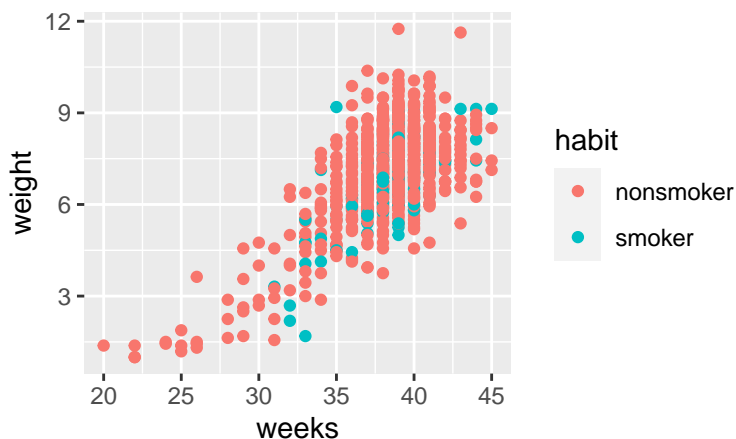
```
## Warning: Removed 2 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 2 rows containing missing values (geom_point).
```



Adding `warning=FALSE` inside `{r}` to get rid of the warning message in the knitted output. The missing value is also removed.

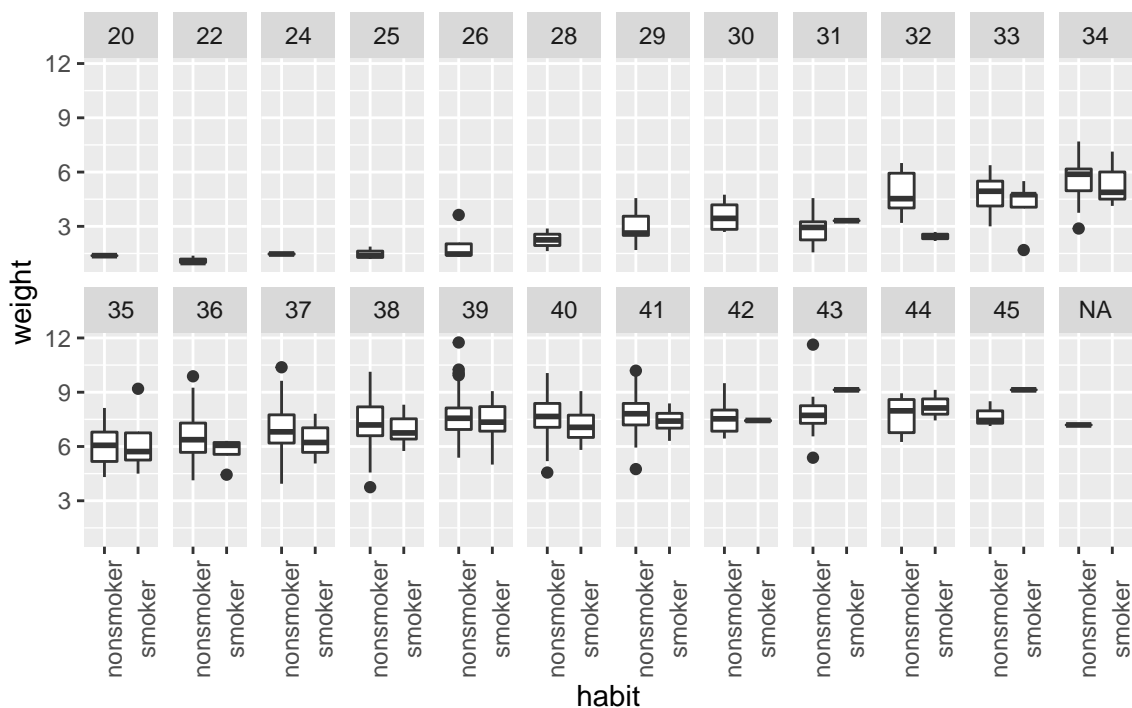
```
nc = subset(nc, !is.na(habit))  
ggplot(nc, aes(x=weeks, y=weight, color=habit)) + geom_point()
```



Hard to see whether mom's smoking habit has any effect on birth weights after accounting for the length of pregnancy (`weeks`).

## 6. Another Way to Control for weeks

```
ggplot(nc, aes(x=habit, y=weight)) +
  geom_boxplot() +
  facet_wrap(~weeks, nrow=2) +
  theme(axis.text.x = element_text(angle = 90))
```



- can ignore week 20-30 where there were no smokers

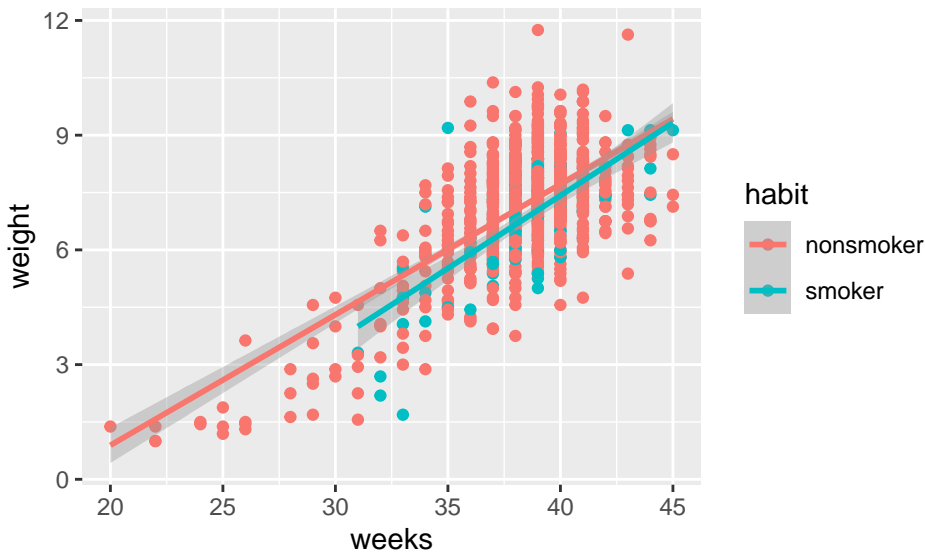
- For 32 to 42 weeks of pregnancies, the median birth weight were slightly lower for those born to smoking moms, comparing babies w/ the same weeks of pregnancy.
- For 31, 43, 44, 45 weeks, smoking group has higher median birth weights, but those groups had only a few observations

```
xtabs(~habit + weeks, data=nc)
```

```
##           weeks
## habit      20  22  24  25  26  28  29  30  31  32  33  34
## nonsmoker   1   3   2   3   4   3   5   4   5   6  11  16
## smoker      0   0   0   0   0   0   0   0   1   2   5   3
##           weeks
## habit      35  36  37  38  39  40  41  42  43  44  45
## nonsmoker  28  42  93 155 201 150  89  22  16  10   3
## smoker     4   4  12  23  36  18  10   3   1   3   1
```

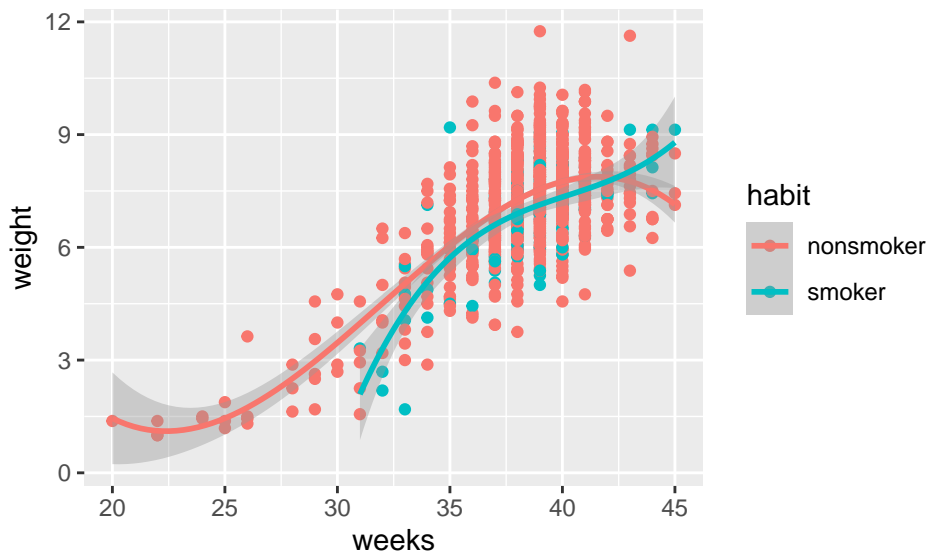
## 7. Controlling for the Effect of weeks Using a Model

```
ggplot(subset(nc, !is.na(habit)), aes(x=weeks, y=weight, color=habit)) +
  geom_point() +
  geom_smooth(method='lm', formula='y~x')
```



Birth weights seem to increase with the length of pregnancy (**weeks**) in a **nonlinear** manner.

```
ggplot(nc, aes(x=weeks, y=weight, color=habit)) +
  geom_point() +
  geom_smooth(method='lm', formula='y~x+I(x^2)+I(x^3)')
```



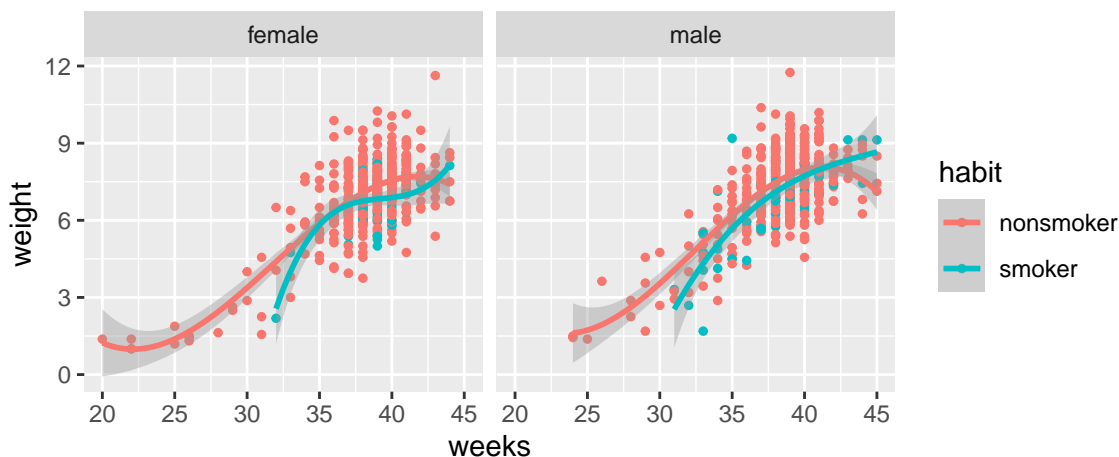
Describe the effect of mother's smoking `habit` on birth weights of babies, after adjusting for the length of pregnancy (`weeks`)

Answer #1: Comparing babies with the same of pregnancy, babies of smoking moms have lower birth weights on average than those of nonsmoking moms.

Answer #2: Babies of smoking moms have a lower mean birth weight than those of nonsmoking moms, controlling/adjusted for the length of pregnancy

## 8. Controlling for Gender As Well ...

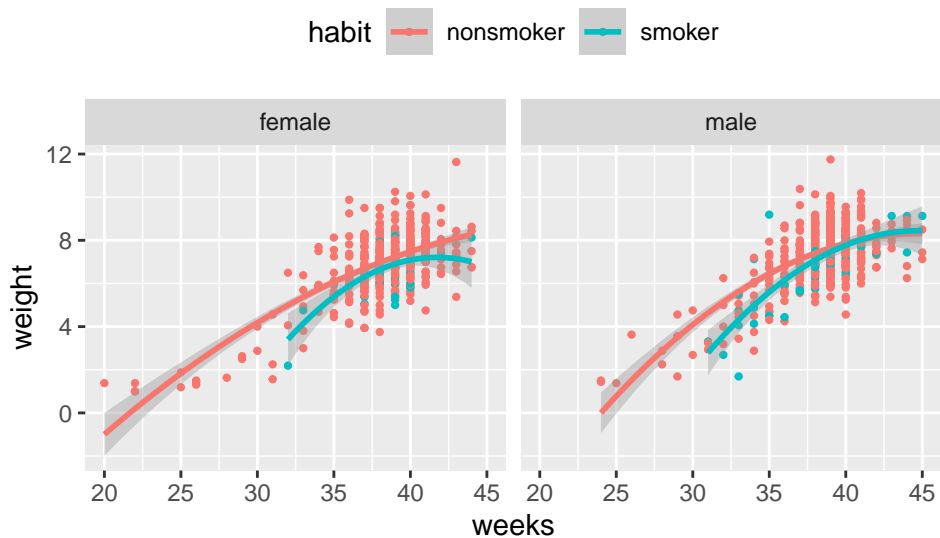
```
ggplot(nc, aes(x=weeks, y=weight, color=habit)) +
  geom_point(size = 1) +
  geom_smooth(method='lm', formula='y~x+I(x^2)+I(x^3)')+
  facet_wrap(~gender)
```



One can move the legend to the top.



```
ggplot(nc, aes(x=weeks, y=weight, color=habit)) +
  geom_point(size=0.8) + facet_wrap(~gender) +
  geom_smooth(method='lm', formula='y~x+I(x^2)')+
  theme(legend.position="top")
```



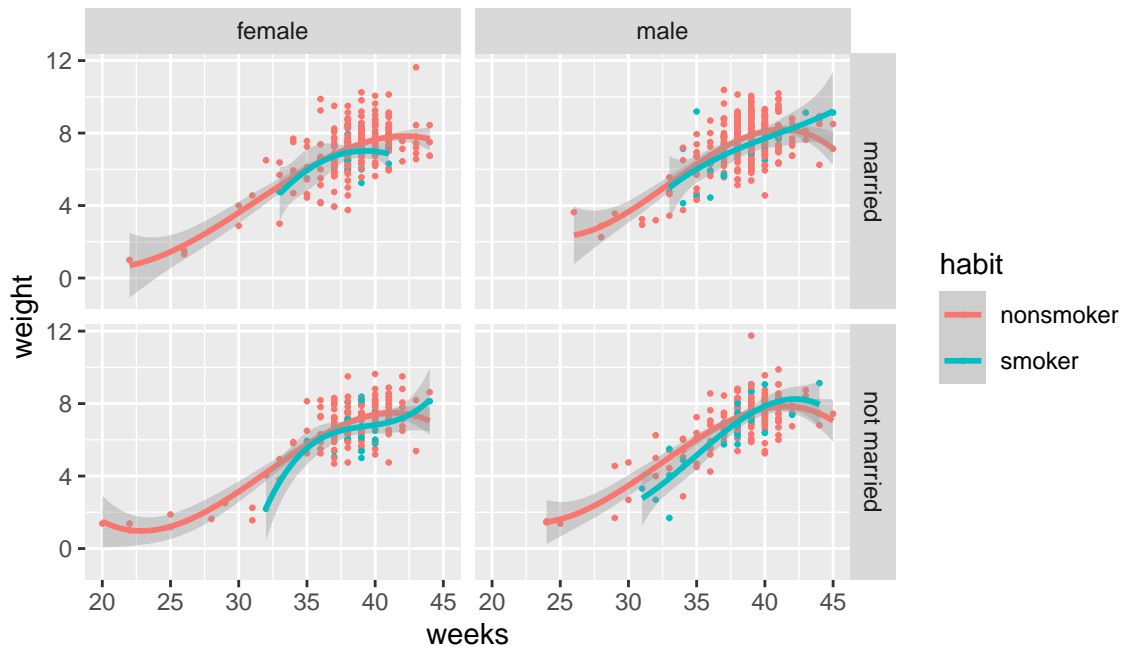
Describe the effect of mother's smoking habit on birth weights of babies, after adjusting for the length of pregnancy (`weeks`) and the gender of the baby.

### On Your Own

**Q3:** Make an appropriate plot of the data and use it to describe the effect of mother's marital status on birth weights of babies, after adjusting for the length of pregnancy (`weeks`) and the gender of the baby.

## 9. Controlling for weeks, gender, and marital

```
ggplot(nc, aes(x=weeks, y=weight, color=habit)) +
  geom_point(size=0.5) +
  geom_smooth(method='lm', formula='y~x+I(x^2)+I(x^3)')+
  facet_grid(marital ~ gender)
```



What's the effect of `habit` on `weight` after adjusting for `weeks`, `gender`, and `marital`?

### On Your Own

**Q4:** Make an appropriate plot of the data and use it to describe the effect of `gender` on `weight` after adjusting for `weeks`, `habit`, and `marital`?

**Q5:** Make another plot of the data and use it to describe the effect of `marital` on `weight` after adjusting for `weeks`, `gender`, and `habit`?

**Q6:** The variable `gained` is mother's weight gain during pregnancy in pounds. Explore to see if `gained` has any effect on baby's birth weight, after adjusting for `weeks`, `gender`, `habit`, and `marital`. If you cannot control for so many predictors at once, try to control for one or two or three or them and see.