

# Political Methodology III: Model Based Inference

Justin Grimmer

Associate Professor  
Department of Political Science  
Stanford University

May 23rd, 2017

## 1) Task:

- Estimate conditional expectation function (regression)  $E[Y|\mathbf{X}]$
- Estimate empirical distribution function  $\hat{f}(x)$
- Estimate conditional empirical distribution function  $\hat{f}(x|\mathbf{Z})$

## 2) Objective Function

- Mean Square Error  $\rightsquigarrow$  Average (Predictive) Risk
- Balance Bias-Variance Tradeoff

## 3) Optimization

- Local Linear Regression

## 4) Validation

- Bias/Variance Tradeoff in fitting data
- LOOCV will be a workhorse tool

# Nonparametric Regression

# Motivating Example: Birth Weight

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?



# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?
- Are differences larger (smaller) within sub-groups (smokers, obese, etc)?

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?
- Are differences larger (smaller) within sub-groups (smokers, obese, etc)?

Suppose that we observe birth weight for a set of  $N$  observations of random variable birth weight,

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?
- Are differences larger (smaller) within sub-groups (smokers, obese, etc)?

Suppose that we observe birth weight for a set of  $N$  observations of random variable birth weight,

$$\mathbf{y} = (y_1, y_2, \dots, y_N)$$

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?
- Are differences larger (smaller) within sub-groups (smokers, obese, etc)?

Suppose that we observe birth weight for a set of  $N$  observations of random variable birth weight,

$$\mathbf{y} = (y_1, y_2, \dots, y_N)$$

We'll also observe a set of  $J$  **covariates** for each observation  $i$

# Motivating Example: Birth Weight

- **Birth Weight** is associated with a wide array of developmental outcomes
- **Birth Weight** varies significantly across socioeconomic and racial groups
- Understand how Birth Weight varies with a variety of attributes
- What are differences in birth weight across racial groups?
- Are differences larger (smaller) within sub-groups (smokers, obese, etc)?

Suppose that we observe birth weight for a set of  $N$  observations of random variable birth weight,

$$\mathbf{y} = (y_1, y_2, \dots, y_N)$$

We'll also observe a set of  $J$  **covariates** for each observation  $i$

$$\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iJ})$$

# Conditional Expectation Functions and Regressions

## Definition

*For any two random variables  $Y$  and  $X$  define the regression of  $Y$  on  $X$  as*

$$f(x) \equiv E[Y|X = x] = \int_{-\infty}^{\infty} y f(y|X = x) dy$$

# Conditional Expectation Functions and Regressions

## Definition

*For any two random variables  $Y$  and  $X$  define the regression of  $Y$  on  $X$  as*

$$f(x) \equiv E[Y|X = x] = \int_{-\infty}^{\infty} yf(y|X = x)dy$$

- **Linear Regression** is a special case with:

$$f(x) = \alpha + \beta x$$

- $E[Y|X]$  is called a **conditional expectation function**

# Conditional Expectation Functions and Regressions

## Definition

*For any two random variables  $Y$  and  $X$  define the regression of  $Y$  on  $X$  as*

$$f(x) \equiv E[Y|X = x] = \int_{-\infty}^{\infty} y f(y|X = x) dy$$

- **Linear Regression** is a special case with:

$$f(x) = \alpha + \beta x$$

- $E[Y|X]$  is called a **conditional expectation function**

We're going to consider more flexible functional forms for  $f(x)$



# Conditional Expectation Functions and Regressions

## Definition

*For any two random variables  $Y$  and  $X$  define the regression of  $Y$  on  $X$  as*

$$f(x) \equiv E[Y|X = x] = \int_{-\infty}^{\infty} yf(y|X = x)dy$$

- **Linear Regression** is a special case with:

$$f(x) = \alpha + \beta x$$

- $E[Y|X]$  is called a **conditional expectation function**

We're going to consider more flexible functional forms for  $f(x)$

Our estimator  $\hat{f}(x)$  will estimate  $f(x)$

# Estimating $E[Y|X = x]$

Easy case: suppose  $X$  is discrete

# Estimating $E[Y|X = x]$

**Easy case:** suppose  $X$  is discrete

Further suppose that  $N \rightarrow \infty$ . Infinite number of realizations of  $Y$ ,  $(y_1, \dots, y_N)$

# Estimating $E[Y|X = x]$

**Easy case:** suppose  $X$  is discrete

Further suppose that  $N \rightarrow \infty$ . Infinite number of realizations of  $Y$ ,  
 $(y_1, \dots, y_N)$

And accompanying covariates  $(x_1, \dots, x_N)$ , then

# Estimating $E[Y|X = x]$

**Easy case:** suppose  $X$  is discrete

Further suppose that  $N \rightarrow \infty$ . Infinite number of realizations of  $Y$ ,  $(y_1, \dots, y_N)$

And accompanying covariates  $(x_1, \dots, x_N)$ , then

$$\begin{aligned}\hat{f}(x) &= \text{average in each bin} \\ &= \frac{\sum_{i=1}^N y_i I(x_i = x)}{\sum_{i=1}^N I(x_i = x)}\end{aligned}$$

# Bias-Variance Tradeoff

# Bias-Variance Tradeoff

**Hard Case:** suppose  $X$  is continuous

# Bias-Variance Tradeoff

**Hard Case:** suppose  $X$  is continuous

Then all values of  $x_1, x_2, \dots, x_N$  are distinct.



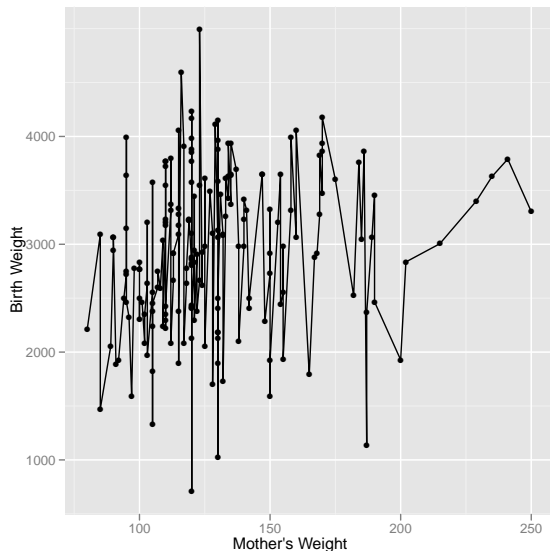
# Bias-Variance Tradeoff

**Hard Case:** suppose  $X$  is continuous

Then all values of  $x_1, x_2, \dots, x_N$  are distinct.

If we calculate conditional mean naively, we get the observed values of  $y_1, \dots, y_N$ :

# Bias-Variance Tradeoff



```
< R Code >  
library(ggplot2)  
qplot(lwt, bwt,  
  geom=c('point'),  
  xlab='Mothers Weight',  
  ylab='Birth Weight')  
+geom_line()
```

# Bias-Variance Tradeoff

Minimize variance:

# Bias-Variance Tradeoff

Minimize variance: assume  $\hat{f}(x)$  is constant

# Bias-Variance Tradeoff

Minimize variance: assume  $\hat{f}(x)$  is constant

$$\hat{f}(x) = \frac{\sum_{i=1}^N Y_i}{N}$$

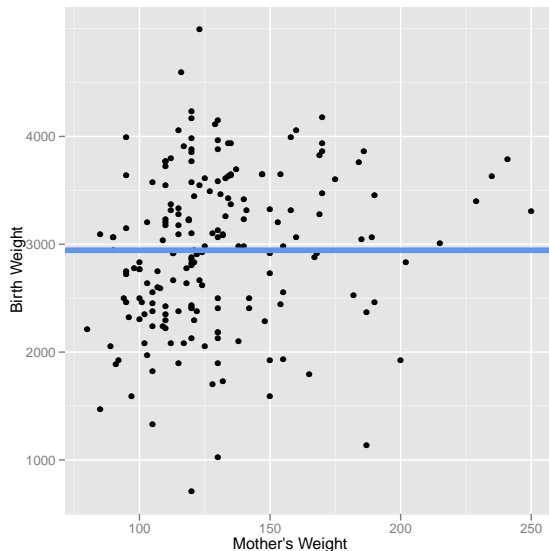
# Bias-Variance Tradeoff

Minimize variance: assume  $\hat{f}(x)$  is constant

$$\hat{f}(x) = \frac{\sum_{i=1}^N Y_i}{N}$$

This yields:

# Bias-Variance Tradeoff



```
library(ggplot2)
qplot(lwt, bwt,
      geom=c('point'),
      xlab='Mother's Weight',
      ylab='Birth Weight')
+geom_abline
(intercept=mean(bwt),
 slope=0, lwd=2,
 col='cornflowerblue')
```

# Compromise



# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

- Suppose  $N$  draws  $\mathbf{X} = (X_1, X_2, \dots, X_{10})$  where

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

- Suppose  $N$  draws  $\mathbf{X} = (X_1, X_2, \dots, X_{10})$  where

$$X_i \sim_{iid} \text{Uniform}(-1, 1)$$

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

- Suppose  $N$  draws  $\mathbf{X} = (X_1, X_2, \dots, X_{10})$  where

$$X_i \sim_{iid} \text{Uniform}(-1, 1)$$

- How many points between  $-0.1$  and  $0.1$  on all dimensions?



# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

- Suppose  $N$  draws  $\mathbf{X} = (X_1, X_2, \dots, X_{10})$  where

$$X_i \sim_{iid} \text{Uniform}(-1, 1)$$

- How many points between  $-0.1$  and  $0.1$  on all dimensions?

$$N \times \Pr(\mathbf{X} \in [-0.1, 0.1] \times [-0.1, 0.1] \times \dots \times [-0.1, 0.1]) =$$

# Curse of Dimensionality

- We're going to learn around a **neighborhood** of a point
- Problem: high-dimensional space is sparse
- Suppose We have  $N$  draws of  $X \sim \text{Uniform}(-1, 1)$
- How many points are between  $[-0.1, 0.1]$ ?

$$N \times \Pr(X \in [-0.1, 0.1]) = N \times \frac{0.2}{2} = \frac{N}{10}$$

- Suppose  $N$  draws  $\mathbf{X} = (X_1, X_2, \dots, X_{10})$  where

$$X_i \sim_{iid} \text{Uniform}(-1, 1)$$

- How many points between  $-0.1$  and  $0.1$  on all dimensions?

$$N \times \Pr(\mathbf{X} \in [-0.1, 0.1] \times [-0.1, 0.1] \times \dots \times [-0.1, 0.1]) = \\ N \times \left(\frac{1}{10}\right)^{10} = \frac{N}{10,000,000,000}$$

# Bias-Variance Tradeoff

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\begin{aligned}\hat{\beta} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y} \\ \hat{f}(\mathbf{X}) \equiv \hat{\mathbf{Y}} &= \underbrace{\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'}_H \mathbf{Y}\end{aligned}$$

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

$$\hat{f}(\mathbf{X}) \equiv \hat{\mathbf{Y}} = \underbrace{\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'}_H \mathbf{Y}$$

$$J+1 = \text{tr}(\mathbf{H})$$

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

$$\hat{f}(\mathbf{X}) \equiv \hat{\mathbf{Y}} = \underbrace{\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'}_H \mathbf{Y}$$

$$J + 1 = \text{tr}(\mathbf{H})$$

$$\hat{f}(\mathbf{x}) = \mathbf{x}' \hat{\boldsymbol{\beta}}$$



# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

$$\hat{f}(\mathbf{X}) \equiv \hat{\mathbf{Y}} = \underbrace{\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'}_H \mathbf{Y}$$

$$J+1 = \text{tr}(\mathbf{H})$$

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \mathbf{x}' \hat{\boldsymbol{\beta}} \\ &= \underbrace{\mathbf{x}'}_{1 \times J} \underbrace{(\mathbf{X}'\mathbf{X})^{-1}}_{J \times J} \underbrace{\mathbf{X}'}_{J \times N} \underbrace{\mathbf{Y}}_{N \times 1} \end{aligned}$$

# Modeling the Relationship

$$Y = f(x_i) + \epsilon_i$$

Suppose  $f(x_i) \equiv \beta_0 + \beta_{j=1}^J x_{ij}$  and  $\mathbf{x} = (1, x_1, \dots, x_J)$ .

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}$$

$$\hat{f}(\mathbf{X}) \equiv \hat{\mathbf{Y}} = \underbrace{\mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'}_H \mathbf{Y}$$

$$J+1 = \text{tr}(\mathbf{H})$$

$$\begin{aligned} \hat{f}(\mathbf{x}) &= \mathbf{x}' \hat{\boldsymbol{\beta}} \\ &= \underbrace{\mathbf{x}'}_{1 \times J} \underbrace{(\mathbf{X}'\mathbf{X})^{-1}}_{J \times J} \underbrace{\mathbf{X}'}_{J \times N} \underbrace{\mathbf{Y}}_{N \times 1} \\ &= \sum_{i=1}^N h_i(\mathbf{x}) Y_i \end{aligned}$$

# Modeling the Relationship

## Definition

We will say  $\hat{f}(\mathbf{x})$ , an estimator of  $f(\mathbf{x})$ , is a **linear smoother** if for each  $\mathbf{x}$  there exists a vector  $\mathbf{h}(\mathbf{x}) = (h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_N(\mathbf{x}))$  such that

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^N h_i(\mathbf{x}) Y_i$$

This implies that fitted values  $\hat{\mathbf{Y}} = (\hat{f}(\mathbf{x}_1), \hat{f}(\mathbf{x}_2), \dots, \hat{f}(\mathbf{x}_N))$  and

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$$

where  $i^{\text{th}}$  row is  $\mathbf{h}(\mathbf{x}_i) = (h_1(\mathbf{x}_1), h_2(\mathbf{x}_1), \dots, h_N(\mathbf{x}_1))$

# Modeling the Relationship

## Definition

*The matrix  $\mathbf{H}$  is the smoothing matrix. The row  $h(x_i)$  is called the effective kernel for observation  $i$ , and the effective number of parameters  $\nu = \text{tr}(\mathbf{H})$ .*

# “Regressogram”

- Suppose we have one variable  $x_i$
- Suppose we divide our data into  $K$  bins,  $B_k = \{x : a < x < b\}$ ,  $k = 1, 2, \dots, K$  so that there are 3 observations in each bin  $|B_k| = 3$  for all  $k$
- We will say that the fitted value for  $x_i$  is then

$$\hat{f}(x_i) = \sum_{x_i \in B_k} \frac{Y_i}{3}$$

# “Regressogram”

- In smoothing terms (sorting observations according to  $x_i$  values)

$$\mathbf{h}(x_i) = (0, 0, 0, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0, 0, 0)$$
$$\mathbf{H} = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 \\ \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \\ 0 & 0 & 0 & 0 & 0 & 0 & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} \end{pmatrix}$$

# Local Averages

- Define  $B_x = \{i : |x_i - x| < b\}$  and  $|B_x| = n_x$ .
- If  $n_x > 0$ ,

$$\hat{f}(x) = \sum_{i \in B_x} \frac{Y_i}{n_x}$$

In smoothing terms

$$\begin{aligned} h_i(x) &= 0 \text{ if } |x_i - x| > b \\ h_i(x) &= \frac{1}{n_x} \text{ if } |x_i - x| < b \end{aligned}$$

# More Sophisticated Weights $\rightsquigarrow$ Kernels

- Previous examples are binary (in/out)



# More Sophisticated Weights $\rightsquigarrow$ Kernels

- Previous examples are binary (in/out)
- Continuous weights via **Kernel**

# More Sophisticated Weights $\rightsquigarrow$ Kernels

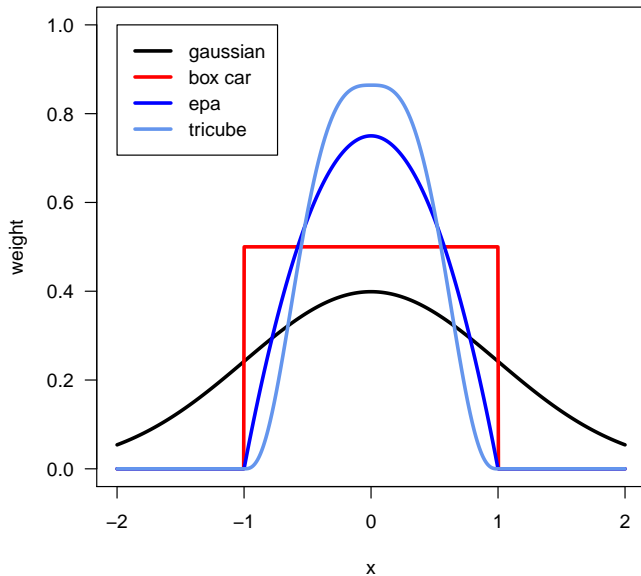
- Previous examples are binary (in/out)
- Continuous weights via **Kernel**
- Some famous examples (Define  $I(x) \equiv 1$  if  $|x| < 1$ , otherwise 0)

$$K(x) = \frac{1}{2}I(x) \text{ (Box Car)}$$

$$K(x) = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-x^2}{2}\right) \text{ (Gaussian)}$$

$$K(x) = \frac{3}{4}(1 - x^2)I(x) \text{ (Epanechnikov)}$$

$$K(x) = \frac{70}{81}(1 - |x|^3)^3 I(x) \text{ (Tri-cube)}$$



# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

Consider the box-car kernel

# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

Consider the box-car kernel

$$\frac{1}{2}I\left(\frac{x - x_i}{1}\right) = \frac{1}{2} \text{ if } |x - x_i| < 1$$

# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

Consider the box-car kernel

$$\begin{aligned}\frac{1}{2}I\left(\frac{x - x_i}{1}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 1 \\ \frac{1}{2}I\left(\frac{x - x_i}{100}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 100\end{aligned}$$



# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

Consider the box-car kernel

$$\begin{aligned}\frac{1}{2}I\left(\frac{x - x_i}{1}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 1 \\ \frac{1}{2}I\left(\frac{x - x_i}{100}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 100 \\ \frac{1}{2}I\left(\frac{x - x_i}{0.01}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 0.01\end{aligned}$$

# Kernels and Bandwidth

Define  $b > 0$  to be the **bandwidth**.

$$K\left(\frac{x - x_i}{b}\right) = \text{Scales difference between } x \text{ and } x_i \text{ by } b$$

Consider the box-car kernel

$$\begin{aligned}\frac{1}{2}I\left(\frac{x - x_i}{1}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 1 \\ \frac{1}{2}I\left(\frac{x - x_i}{100}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 100 \\ \frac{1}{2}I\left(\frac{x - x_i}{0.01}\right) &= \frac{1}{2} \text{ if } |x - x_i| < 0.01\end{aligned}$$

$b \rightsquigarrow$  controls the sensitivity of Kernel

# Local Regression

## Definition

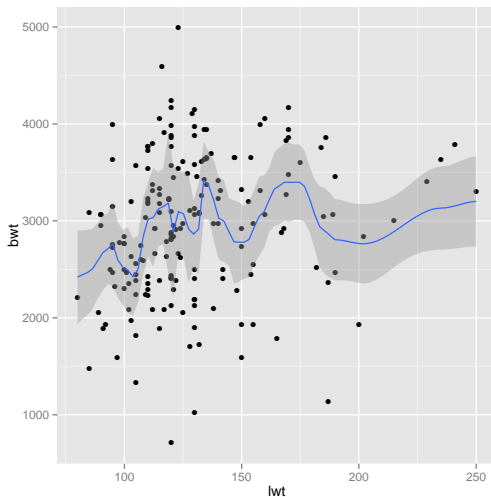
*Suppose  $b > 0$ . Define the local regression estimator  $\hat{f}(x)$  of  $f(x)$  as*

$$\hat{f}(x) = \sum_{i=1}^N h_i(x) Y_i$$

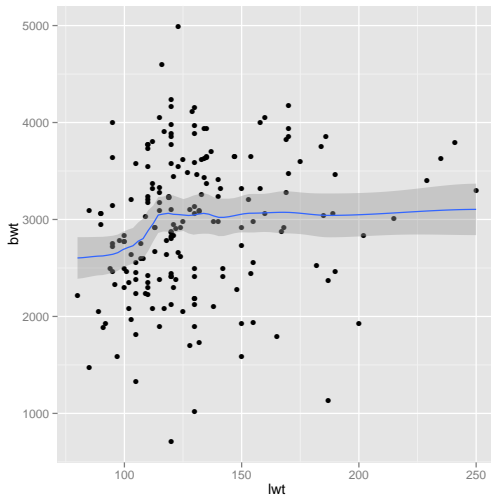
*where*

$$h_i(x) = \frac{K\left(\frac{x-x_i}{b}\right)}{\sum_{j=1}^N K\left(\frac{x-x_j}{b}\right)}$$

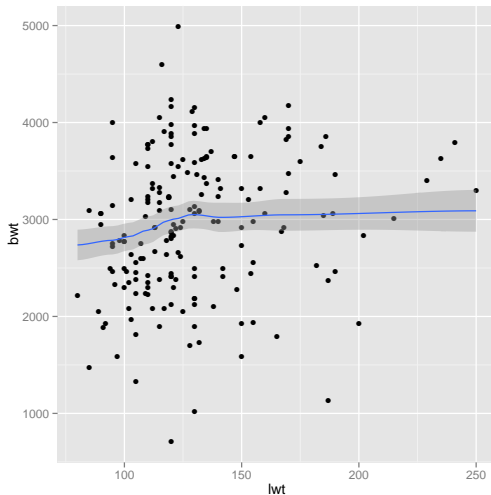
*$\mathbf{H} = N \times N$  matrix with  $H_{ij} = h_j(x_i)$*



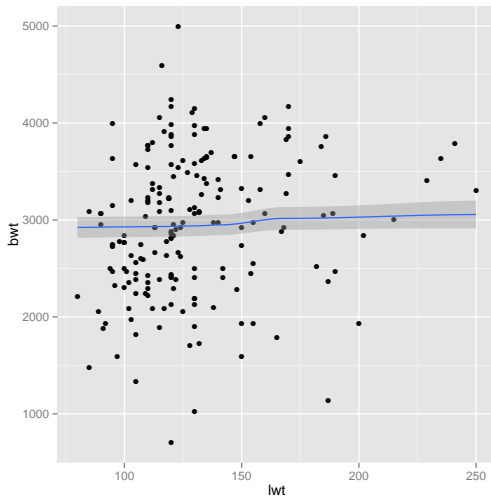
```
qplot(lwt, bwt,  
  geom='point') +  
  stat_smooth(degree = 0,  
    span = 0.1)
```



```
qplot(lwt, bwt,  
geom='point') +  
stat_smooth(degree = 0,  
span = 0.4)
```



```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 0,  
             span = 0.7)
```



```
qplot(lwt, bwt,  
geom='point') +  
stat_smooth(degree = 0,  
span = 1)
```

# Local Regression

$$w_i(x) = K\left(\frac{x_i - x}{b}\right)$$

$$\hat{a} = \arg \min_a \sum_{i=1}^N w_i(x)(Y_i - a)^2$$

$$\hat{a} = \frac{\sum_{i=1}^N w_i(x)Y_i}{\sum_{i=1}^N w_i(x)}$$

$$\hat{a} = \hat{f}(x)$$



# Local Polynomial regression

$$(\hat{a}, \hat{b}, \hat{c}) = \arg \min_{a,b,c} \sum_{i=1}^N w_i(x) (Y_i - a - bx_i - cx_i^2)^2$$

Define  $\mathbf{X}_x$  to be an  $N \times 3$  matrix:

$$\mathbf{X}_x = \begin{pmatrix} 1 & x_1 - x & (x_1 - x)^2 \\ 1 & x_2 - x & (x_2 - x)^2 \\ \vdots & \vdots & \vdots \\ 1 & x_N - x & (x_N - x)^2 \end{pmatrix}$$

$\mathbf{W}_x = N \times N$  matrix where  $W_{ii} = w_i(x)$ .

$$\hat{f}(x) = \left( \mathbf{X}_x' \mathbf{W}_x \mathbf{X}_x \right)^{-1} \mathbf{X}_x' \mathbf{W}_x \mathbf{Y}$$

# Local Polynomial Regression

## Definition

*Suppose  $b > 0$ . Define the local polynomial regression estimator  $\hat{f}(x)$  of  $f(x)$  as*

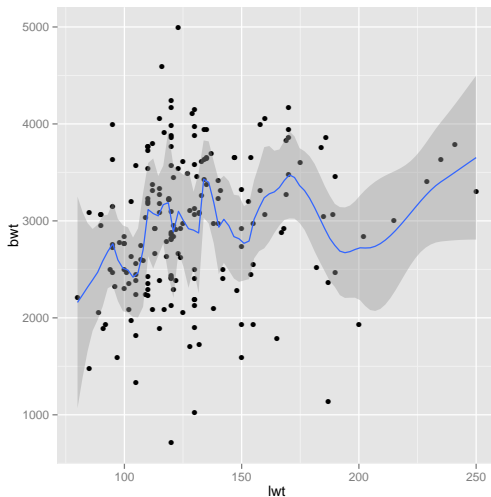
$$\hat{f}(x) = \left( \mathbf{X}'_x \mathbf{W}_x \mathbf{X}_x \right)^{-1} \mathbf{X}'_x \mathbf{W}_x \mathbf{Y}$$

*We can write this as  $\hat{f}(x) = \sum_{i=1}^N h_i(x) Y_i$ , where  $\mathbf{h}(x) = (h_1(x), h_2(x), \dots, h_N(x))$*

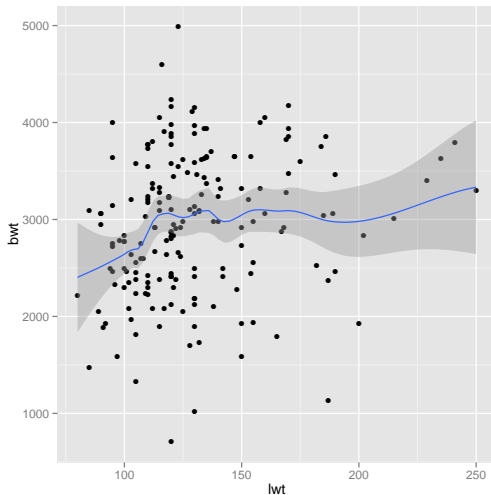
$$\mathbf{h}(x) = \mathbf{e}_1 \left( \mathbf{X}'_x \mathbf{W}_x \mathbf{X}_x \right)^{-1} \mathbf{X}'_x \mathbf{W}_x$$

*where  $\mathbf{e}_1 = (1, 0, \dots, 0)$ .*

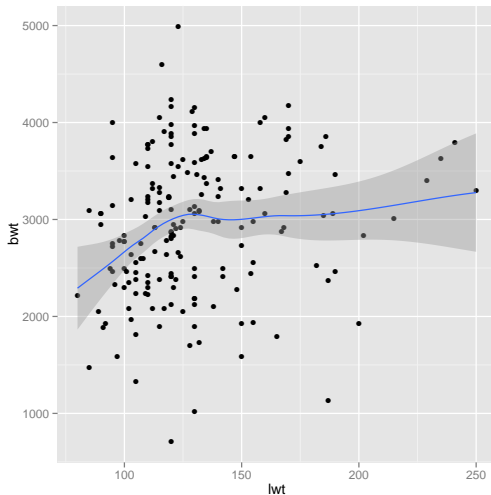
*$\mathbf{H}$  is an  $N \times N$  matrix with typical entry  $H_{ij} = h_j(x_i)$ .*



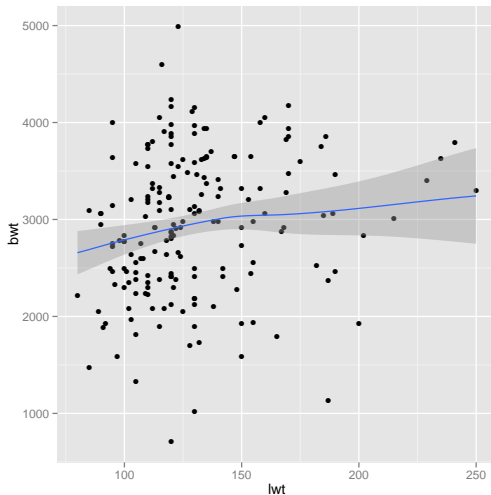
```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 1,  
             span = 0.1)
```



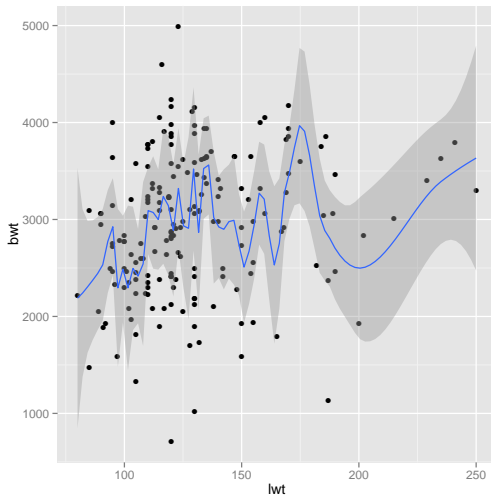
```
qplot(lwt, bwt,  
  geom='point') +  
  stat_smooth(degree = 1,  
    span = 0.4)
```



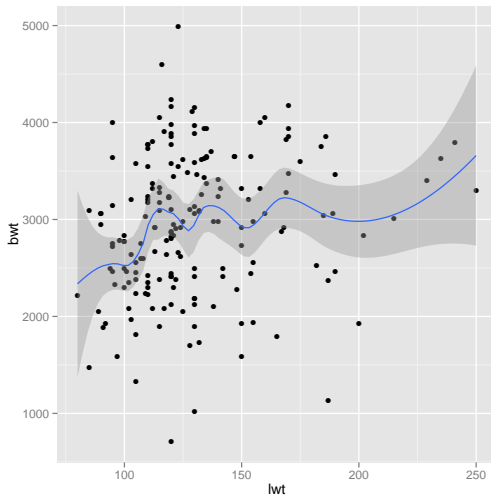
```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 1,  
             span = 0.7)
```



```
qplot(lwt, bwt,  
  geom='point') +  
  stat_smooth(degree = 1,  
    span = 1)
```

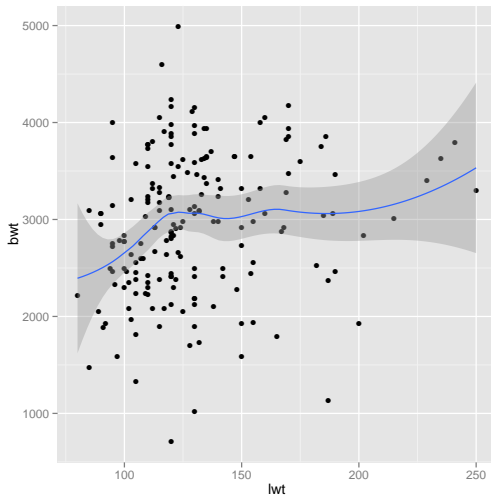


```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 2,  
             span = 0.1)
```

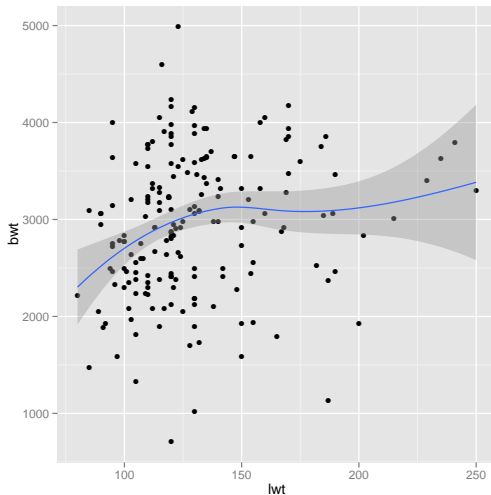


```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 2,  
             span = 0.4)
```

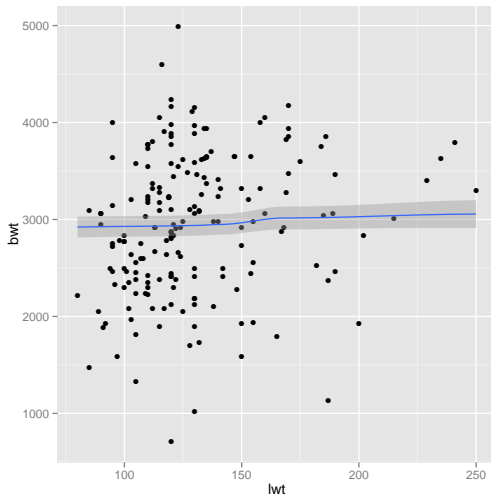




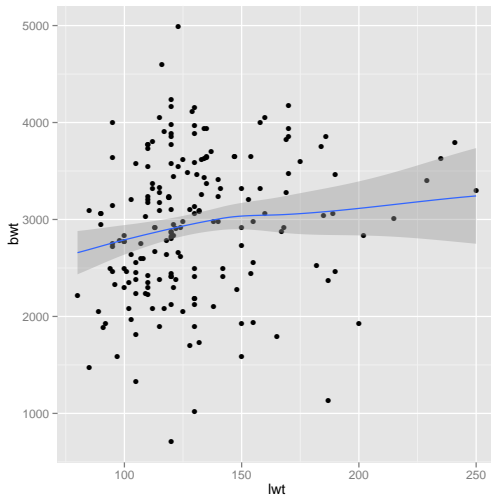
```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 2,  
             span = 0.7)
```



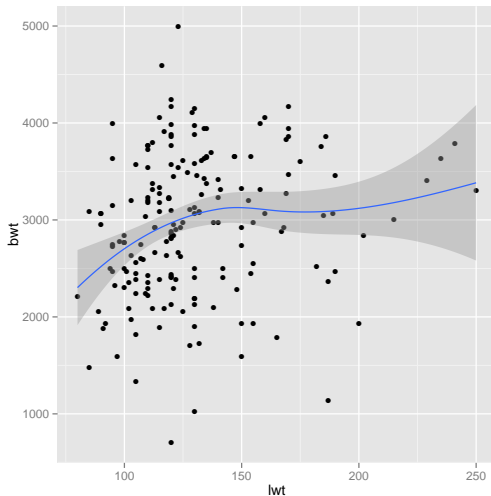
```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 2,  
             span = 1)
```



```
qplot(lwt, bwt,  
  geom='point') +  
  stat_smooth(degree = 0,  
    span = 1)
```



```
qplot(lwt, bwt,  
      geom='point') +  
  stat_smooth(degree = 1,  
             span = 1)
```



```
qplot(lwt, bwt,  
      geom='point') +  
      stat_smooth(degree = 1,  
                  span = 1)
```

# Bandwidth Selection

# Leave One Out Cross Validation LOOCV

- Define  $\hat{f}_{-i,b}(x)$  to be the linear smooth estimator without observation  $i$  and bandwidth  $b$ .
- The LOOCV statistic will be

$$\text{CV}(b) = \frac{1}{N} \sum_{i=1}^N \left( Y_i - \hat{f}_{-i,b}(x) \right)^2$$

## Theorem

*Suppose  $\hat{f}(x)$  is a linear smoother. Then  $CV(b)$  can be written as*

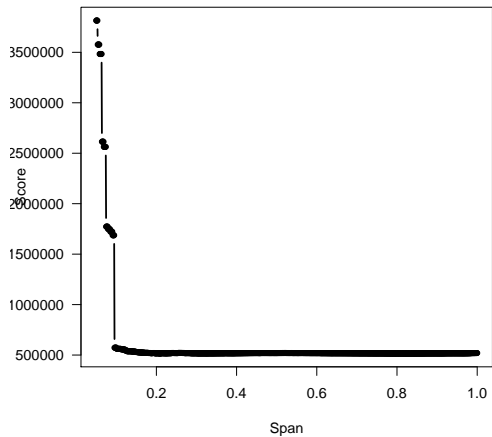
$$CV(b) = \frac{1}{N} \sum_{i=1}^N \left( \frac{Y_i - \hat{f}_b(x)}{1 - H_{ii}} \right)^2$$

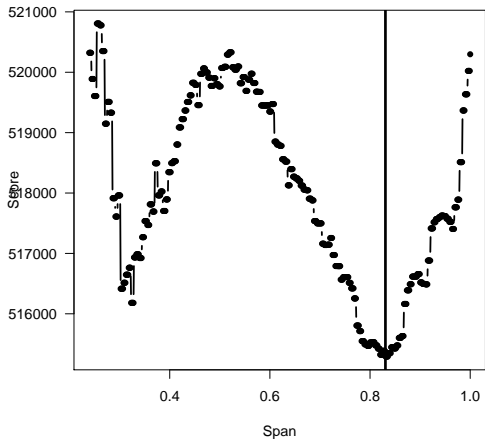
*And is often approximated with the Generalized Cross Validation (GCV):*

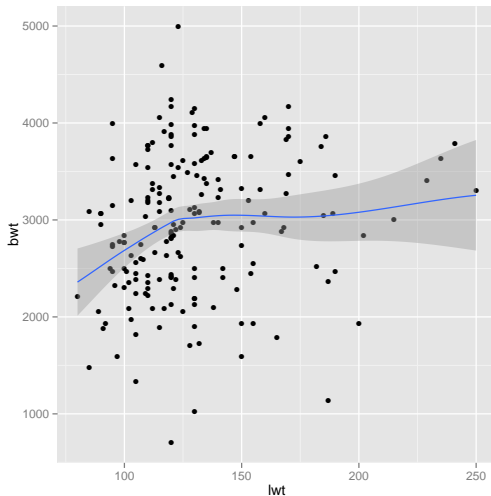
$$GCV(b) = \frac{1}{N} \sum_{i=1}^N \left( \frac{Y_i - \hat{f}_b(x)}{1 - \nu/N} \right)^2$$

*where  $\nu = \text{tr}(\mathbf{H})$*









```
qplot(lwt, bwt,  
  geom='point') +  
  stat_smooth(degree = 1,  
    span = 0.83)
```

# Loess and Conditioning

Many models:

# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship

# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship  
Possible to visualize relationships in **lattice**:

# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship

Possible to visualize relationships in **lattice**:

Suppose we have continuous random variable  $X$  and discrete random variables  $Z_1$  and  $Z_2$ .

# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship

Possible to visualize relationships in **lattice**:

Suppose we have continuous random variable  $X$  and discrete random variables  $Z_1$  and  $Z_2$ .

Then we might estimate



# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship

Possible to visualize relationships in **lattice**:

Suppose we have continuous random variable  $X$  and discrete random variables  $Z_1$  and  $Z_2$ .

Then we might estimate

$$f(x, z_1, z_2) = E[Y|X = x, Z_1 = z_1, Z_2 = z_2]$$

# Loess and Conditioning

Many models: condition on **many** factors in order to learn relationship

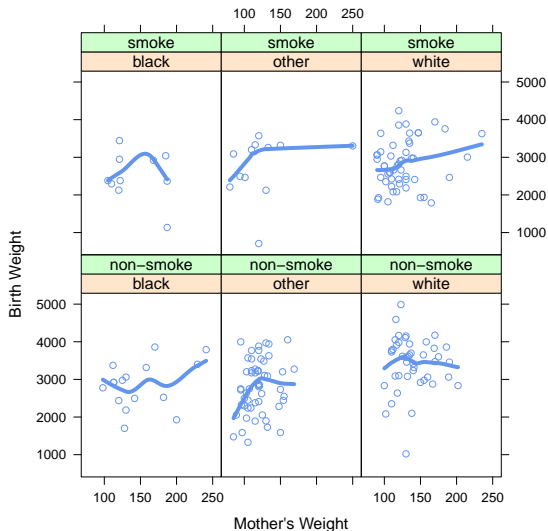
Possible to visualize relationships in **lattice**:

Suppose we have continuous random variable  $X$  and discrete random variables  $Z_1$  and  $Z_2$ .

Then we might estimate

$$\begin{aligned} f(x, z_1, z_2) &= E[Y|X = x, Z_1 = z_1, Z_2 = z_2] \\ &= \int_{-\infty}^{\infty} y f(y|X = x, Z_1 = z_1, Z_2 = z_2) dy \end{aligned}$$

# Loess and Conditioning



< R Code >

```
xyplot(bwt~lwt|race*smoke,
       lwd=4,
       col='cornflowerblue',
       xlab='Mother's Weight',
       ylab='Birth Weight',
       panel= function(x, y,
       ...){
         panel.xyplot(x, y,
         ...);
         panel.loess(x,y,span=.8,
         degree=1, ...) } )
```

# Densities

# Histogram, Probability Mass Function

Simplest case: probability mass functions

# Histogram, Probability Mass Function

Simplest case: probability mass functions

- Suppose we observe  $x_1, x_2, \dots, x_N$

# Histogram, Probability Mass Function

Simplest case: **probability mass functions**

- Suppose we observe  $x_1, x_2, \dots, x_N$
- If  $X$ 's takes on few unique values, say  $K \ll N$ , then an estimator for the pmf is,

# Histogram, Probability Mass Function

Simplest case: **probability mass functions**

- Suppose we observe  $x_1, x_2, \dots, x_N$
- If  $X$ 's takes on few unique values, say  $K \ll N$ , then an estimator for the pmf is,

$$\hat{p}(X = x) \equiv \hat{p}(x) = \frac{\sum_{i=1}^N I(x_i = x)}{N}$$



# Histogram, Probability Mass Function

Simplest case: **probability mass functions**

- Suppose we observe  $x_1, x_2, \dots, x_N$
- If  $X$ 's takes on few unique values, say  $K \ll N$ , then an estimator for the pmf is,

$$\hat{p}(X = x) \equiv \hat{p}(x) = \frac{\sum_{i=1}^N I(x_i = x)}{N}$$

Will converge on the probability mass function (by weak law of large numbers)

# Histogram, Probability Density Function

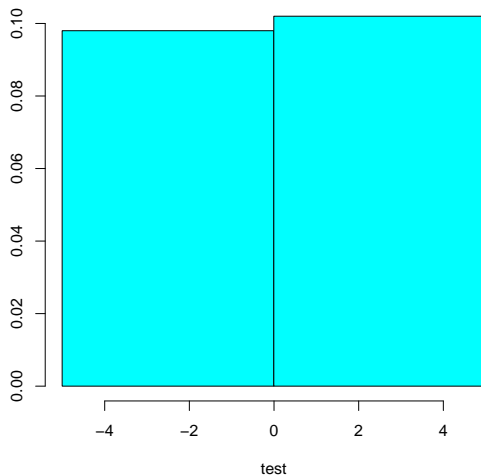
## Harder case: probability density functions

- Suppose we observe  $x_1, x_2, \dots, x_N$ , with  $X$ 's continuous
- Goal: estimate the pdf  $f(x)$ .
- Define a set of bins  $b_1, b_2, \dots, b_M$  such that
  - $b_1 \leq x_1$  and  $b_M \geq x_N$
- Then, we can define an estimator for the pdf  $f(x)$  as

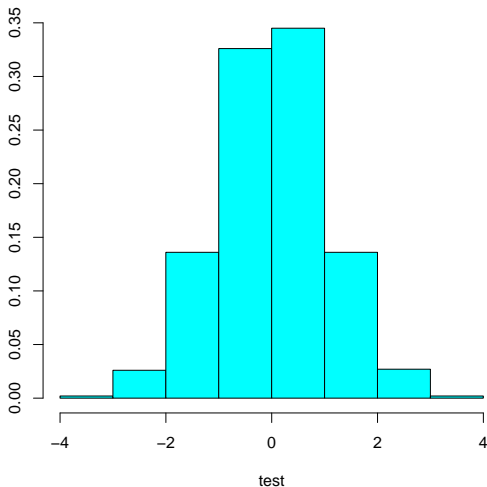
$$\widehat{f(x)} = \frac{\sum_{i=1}^N I(x_i > b_z \text{ and } x_i \leq b_{z+1})}{N} \times I(x \in (b_z, b_{z+1}])$$

In words:  $\widehat{f(x)}$  is equal to proportion of observations in the bandwidth that  $x$  resides in

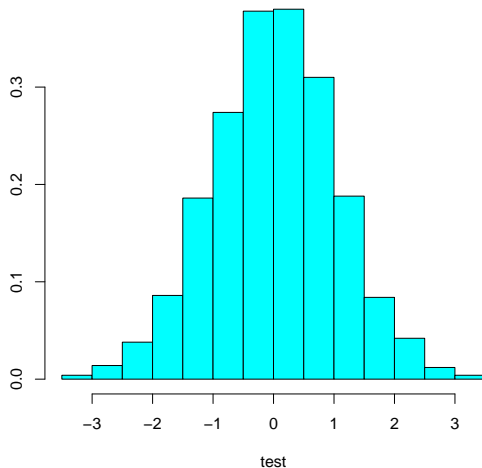
# Bias Variance Tradeoff on Bin Size



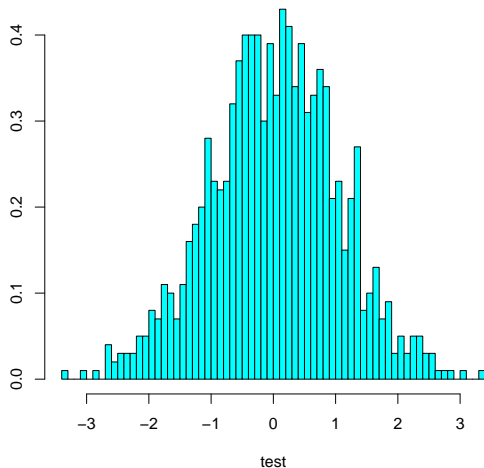
# Bias Variance Tradeoff on Bin Size



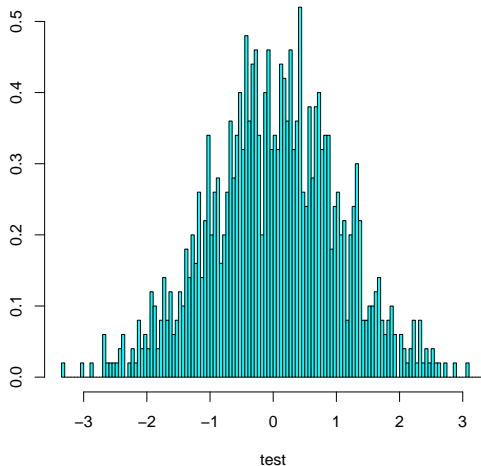
# Bias Variance Tradeoff on Bin Size



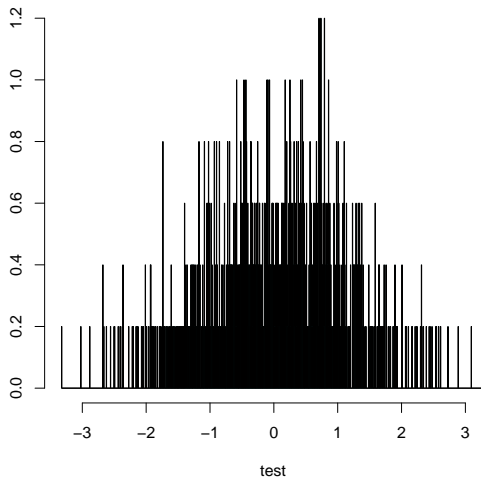
# Bias Variance Tradeoff on Bin Size



# Bias Variance Tradeoff on Bin Size



# Bias Variance Tradeoff on Bin Size

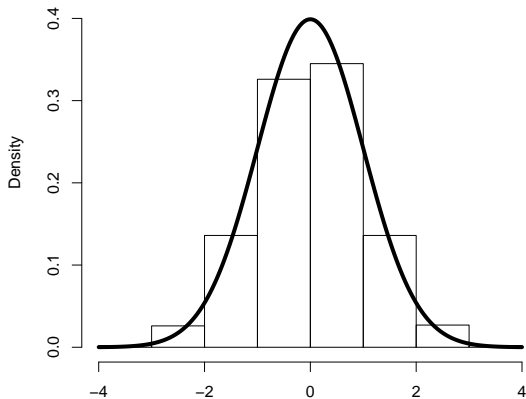




# Bias-Variance Tradeoff/Guidance

Classic problem:

**Bias-Variance** Tradeoff

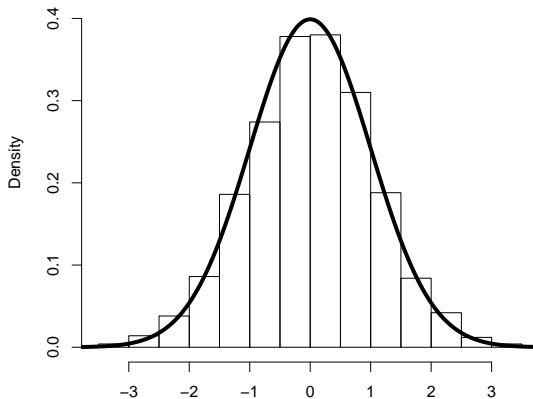


test

# Bias-Variance Tradeoff/Guidance

Classic problem:

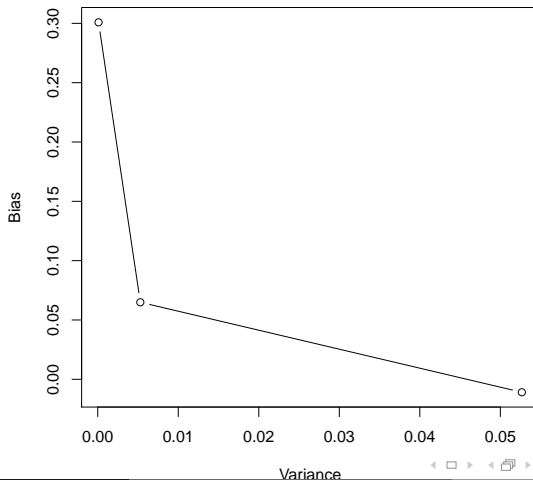
**Bias-Variance** Tradeoff



# Bias-Variance Tradeoff/Guidance

Classic problem:

**Bias-Variance** Tradeoff



# Navigating the Bias Variance Tradeoff

Two formulas for navigating bias-variance tradeoff for data  $x$

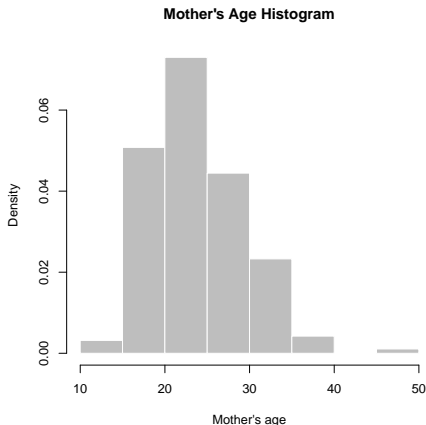
1)  $\text{width} = \frac{3.5 \times \text{sd}(x)}{n^{1/3}}$  (Scott 1979)

2)  $\text{width} = \frac{2(Q_3 - Q_1)}{n^{1/3}}$  (Freedman and Diaconis (1981) )

# Histogram Creation

< R Code >

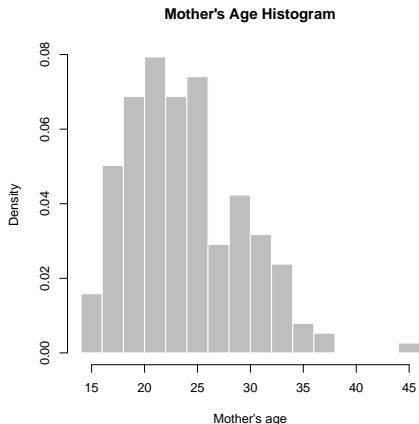
```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white')
```



# Histogram Creation

< R Code >

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='fd')
```



# Histogram Creation

< R Code >

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='scott')
```



# Density Estimation

Histograms make continuous data discrete

**Density** estimators “smooth” out histograms.

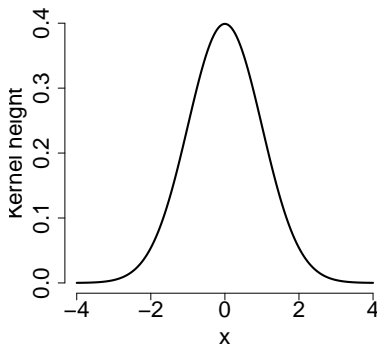
We will use Kernels again



# Gaussian Kernel

Define

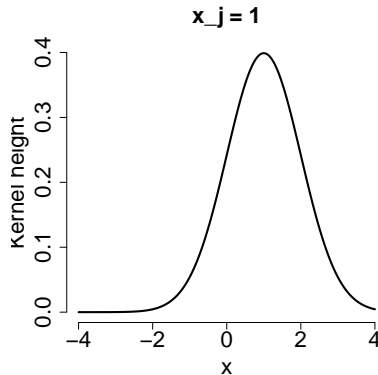
$$K(x) = \frac{\exp[-x^2]}{\sqrt{2\pi}}$$



# Gaussian Kernel

Define

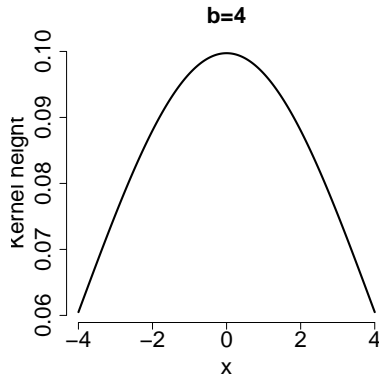
$$K\left(\frac{x - x_j}{b}\right) = \frac{\exp\left[-\left(\frac{x - x_j}{b}\right)^2\right]}{\sqrt{2\pi}}$$



# Gaussian Kernel

Define

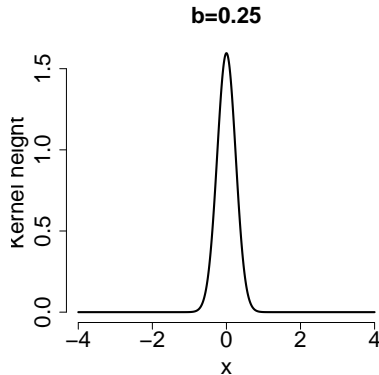
$$K\left(\frac{x - x_j}{b}\right) = \frac{\exp\left[-\left(\frac{x - x_j}{b}\right)^2\right]}{\sqrt{2\pi}}$$



# Gaussian Kernel

Define

$$K\left(\frac{x - x_j}{b}\right) = \frac{\exp\left[-\left(\frac{x - x_j}{b}\right)^2\right]}{\sqrt{2\pi}}$$



# Density Estimation

For all values  $x \in \mathfrak{R}$  define the kernel density estimator as,

$$\hat{f}(x) = \frac{1}{nb} \sum_{j=1}^n K\left(\frac{x - x_j}{b}\right)$$

In words: at each point  $x$ , calculate weighted average of points, where weight is given by kernel

- As  $b$  increases, weight on nearby points is more evenly distributed (redistributed away from  $x$ )
- As  $b$  decreases, weight on nearby points is more concentrated (redistributed towards  $x$ )

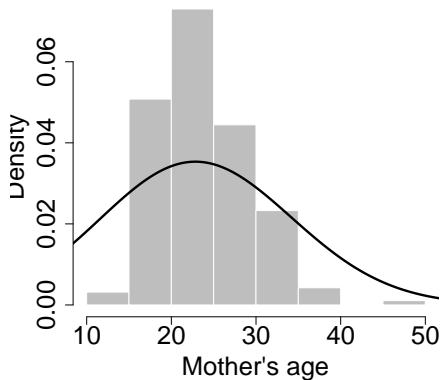
Another bias variance tradeoff!

# Density Estimation

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='scott')
```

```
lines(density(birthwt$age, width=40), lwd=3)
```

**Mother's Age Histogram**

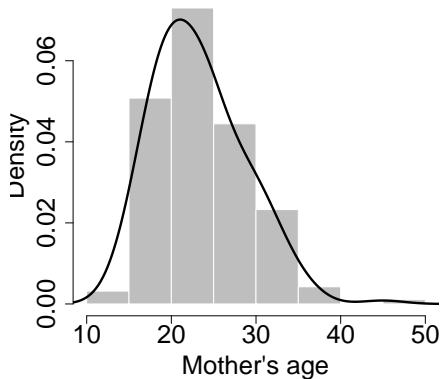


# Density Estimation

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='scott')
```

```
lines(density(birthwt$age, width=10), lwd=3)
```

**Mother's Age Histogram**

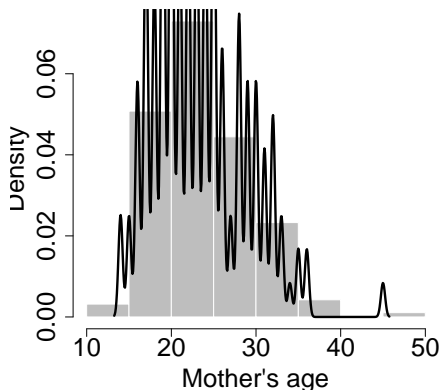


# Density Estimation

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='scott')
```

```
lines(density(birthwt$age, width=1), lwd=3)
```

**Mother's Age Histogram**

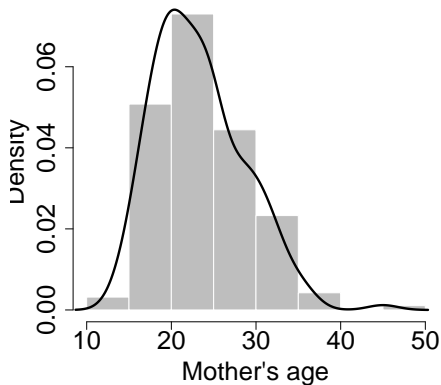




# Density Estimation

```
truehist(birthwt$age, xlab='Mother's age',  
ylab='Density', main='Mother's Age Histogram', col='gray',  
border='white', nbins='scott')  
lines(density(birthwt$age, width='SJ-dpi'), lwd=3)
```

**Mother's Age Histogram**



# Back to the Birth Weight

An initial inference:

# Back to the Birth Weight

An initial inference:  
Compare (contrast):

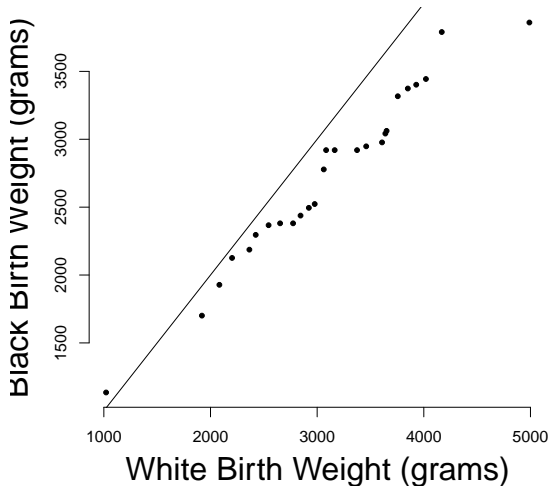
# Back to the Birth Weight

An initial inference:  
Compare (contrast):

$$\hat{p}(y|\text{white}) \quad \text{and} \quad \hat{p}(y|\text{black})$$

# Back to the Birth Weight

Comparing Birthweight across Racial Groups



```
qqplot(  
  bwt[which(race2=='white')],  
  bwt[which(race2=='black')],  
  xlab='White Birth Weight  
  (grams)', ylab='Black Birth  
  Weight (grams)',  
  frame.plot=F,  
  main='Comparing Birth  
  Weight across Racial Groups',  
  pch=20)  
arrows(0,0, 1e7, 1e7)
```

# Box Plot

# Box Plot

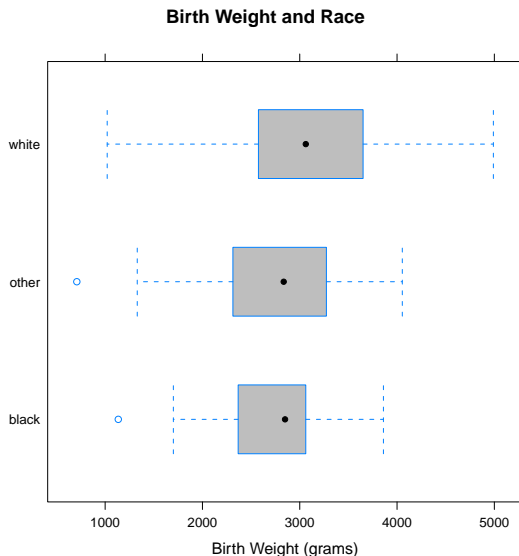
- QQ-plots compare two distributions directly

# Box Plot

- QQ-plots compare two distributions directly
- Often want to compare many distributions (and readers are uncomfortable with qq-plots)



# Box Plot



< R Code >

```
library(lattice)
bwplot(race2 ~bwt,
xlab='Birth Weight
(grams)',
frame.plot=F,main =
'Birthweight and Race',
panel=function(...){
panel.bwplot(...,
fill='gray', lty=1,
lwd=2, pch=20) } )
```

# Violin Plot

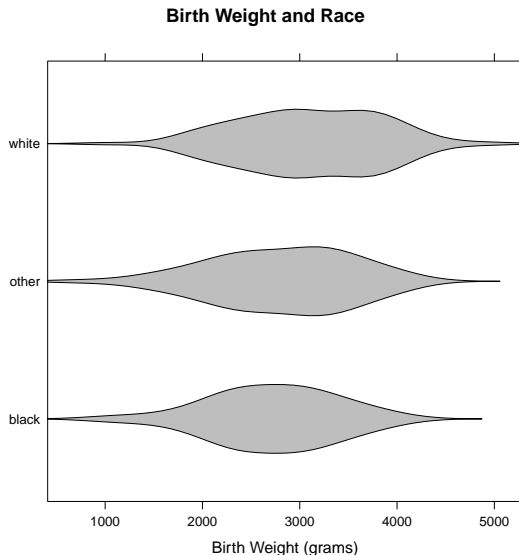
# Violin Plot

- Box plots can obscure variation in distributions [histogram]

# Violin Plot

- Box plots can obscure variation in distributions [histogram]
- **Violin plots** allow smooth estimates of data distribution [density]

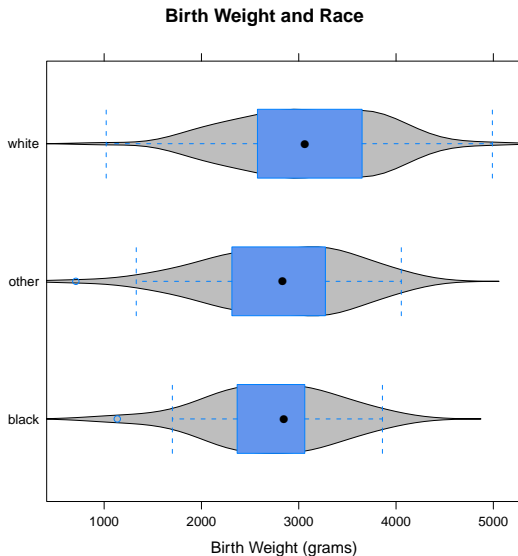
# Violin Plot



< R Code >

```
library(lattice)
bwplot(race2~bwt,
xlab='Birth Weight
(grams)', frame.plot=F,
main='Birth Weight and
Race',
panel=function(...){
panel.violin(...,
col='gray',
bw='SJ-dpi') } )
```

# Violin Plot



```
< R Code >  
library(lattice)  
bwplot(race2~bwt,  
xlab='Birth Weight  
(grams)', frame.plot=F,  
main='Birth Weight and  
Race',  
panel=function(...){  
  panel.violin(...,  
    col='gray',  
    bw='SJ-dpi') } )
```

# Box (Violin) Plot with Conditioning

- Comparing distributions within **strata** (buckets)

# Box (Violin) Plot with Conditioning

- Comparing distributions within **strata** (buckets)
  - Here: compare distribution of birth weights across race, by smoker and non-smoker



# Box (Violin) Plot with Conditioning

- Comparing distributions within **strata** (buckets)
  - Here: compare distribution of birth weights across race, by smoker and non-smoker
- Formally:

# Box (Violin) Plot with Conditioning

- Comparing distributions within **strata** (buckets)
  - Here: compare distribution of birth weights across race, by smoker and non-smoker
- Formally:

$$\hat{p}(y|\text{black, smoke})$$

$$\hat{p}(y|\text{white, smoke})$$

$$\hat{p}(y|\text{other, smoke})$$

# Box (Violin) Plot with Conditioning

- Comparing distributions within **strata** (buckets)
  - Here: compare distribution of birth weights across race, by smoker and non-smoker
- Formally:

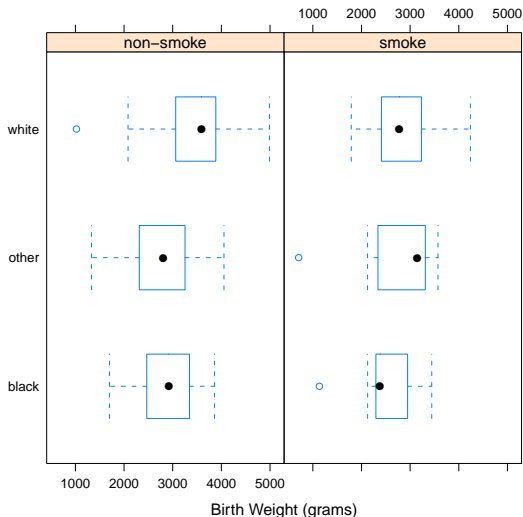
$$\hat{p}(y|\text{black, smoke}) \quad ; \quad \hat{p}(y|\text{black, non-smoke})$$

$$\hat{p}(y|\text{white, smoke}) \quad ; \quad \hat{p}(y|\text{white, non-smoke})$$

$$\hat{p}(y|\text{other, smoke}) \quad ; \quad \hat{p}(y|\text{other, non-smoke})$$

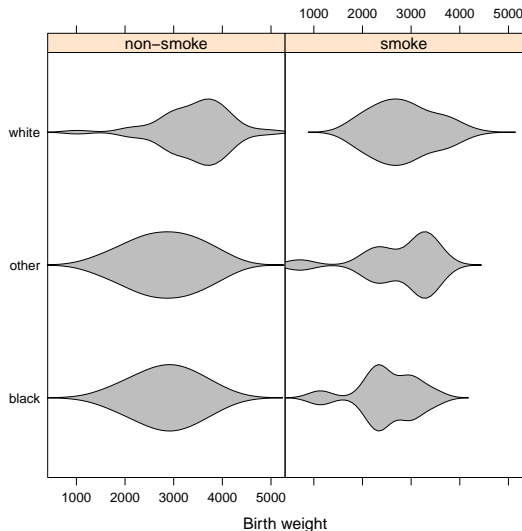
# Box (Violin) Plot with Conditioning

Birth Weight and Race, Given Smoking



```
library(lattice)
bwplot(race2~bwt|smoke,
xlab='Birth Weight
(grams)',
frame.plot=F,main =
'Birth Weight and Race,
Given Smoking')
```

# Box (Violin) Plot with Conditioning



```
library(lattice)
bwplot(race2~bwt|smoke,
       xlab='Birth weight',
       panel=function(...){
         panel.violin(...,
                       col='gray',bw='SJ-dpi')
       })
```

## Nonparametrics: Bias/Variance Tradeoff