

# Lecture Note 12: Non-parametric Methods

NOT ON FINAL.

parametric vs. non-parametric

↳ ① estimating the density

↳ ② estimating the reg. function

Linear relationship:  $Y = \beta_0 + \beta_1 X + U$

↓  
 $E[Y|X] = \beta_0 + \beta_1 X$

Non-linear relationship?  $E[Y|X] = g(X)$

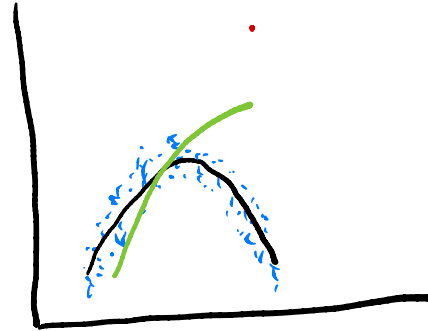
Want to estimate:  $Y = \underbrace{g(X)}_{\text{2 methods}} + U$

- 1) polynomial regression
- 2) local linear regression

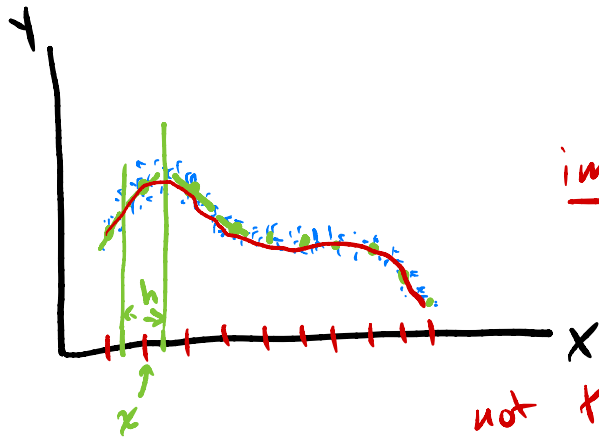
## Polynomial regression

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_k X^k + U$$

- Calculus: As  $k \rightarrow \infty$ , polynomial converges to  $g(x)$
- bias vs. variance: high  $k \rightarrow$  less bias, more variance  
low  $k \rightarrow$  more bias, less variance
- unappealing: outlier sensitivity



# Local regression



→ local linear regression

→ bandwidth:  $h$  (half of window)

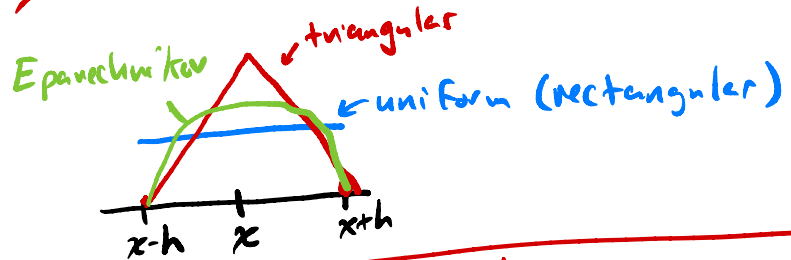
→ bias vs. variance → high  $h$ : more bias  
less variance

low  $h$ : less bias  
more variance

important

→ kernel function:  $K\left(\frac{x_i - x}{h}\right)$

not that important



→ for each  $x$ , estimate:

$$\min_{b_0, b_1} \sum_i K\left(\frac{x_i - x}{h}\right) (y_i - b_0 - b_1 x_i)^2$$

Stats: `lpol`

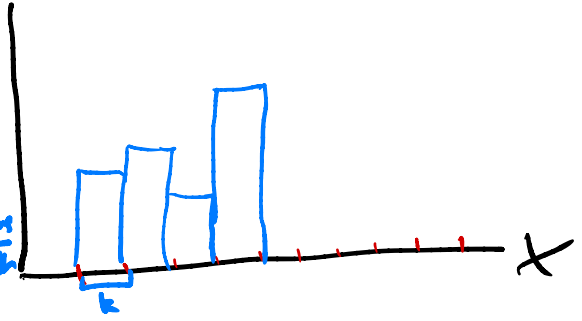
R: `kernSmooth::locpoly()`

ggplot: `geom-smooth()`

# Estimating the density

→ Histogram:

→ vertical axis: → count / frequency:  $N_k$   
→ share of obs:  $N_k/N$   
→ density:  $\hat{f}(x) = \frac{N_k}{N} \cdot \frac{\text{num of obs}}{\text{bin width}}$



→ Stata: hist

R: hist() or ggplot::geom-histogram()

→ Centered histogram

↪  
→ kernel density estimation

$$\hat{f}(x) = \frac{1}{Nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right)$$

Stata: kdensity

R: ggplot::geom-density()

