

# Lecture Note 12: Non-parametric Methods

Not on final

parametric vs. non-parametric

↳ ① estimating the density

↳ ② estimating the regression function

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

Conditional expectation:  $E[Y|X] = \beta_0 + \beta_1 X$

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

← flexible function

← regression function

Want to estimate  $Y = g(X) + U$

↳ 2 ways to estimate: ① polynomial regression  
② local regression

## Polynomial regression

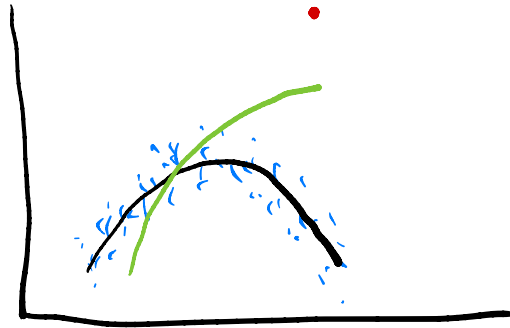
$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \dots + \beta_k X^k + U$$

→ calculus: as  $k \rightarrow \infty$ , polynomial converges to  $g(x)$

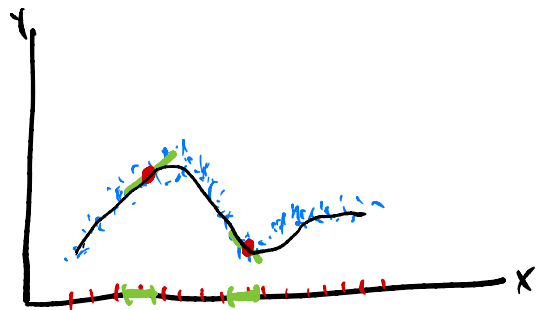
→ bias vs. variance: large  $k \rightarrow$  less bias, more variance

small  $k \rightarrow$  more bias, less variance

→ unappealing feature:



# Local regression



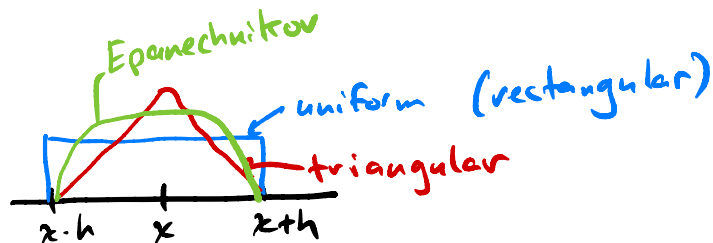
local linear regression

→ bandwidth:  $h$  (half of window)

→ bias vs. variance: large  $h$ : more bias, less var

small  $h$ : less bias, more var

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



→ local linear regression:

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

→ Stata: `lpol`

→ R: `kernSmooth::locpol()`

`ggplot::geom-smooth()`

# Density estimation

Histogram:

→ vertical axis: → number of obs:  $N_k$

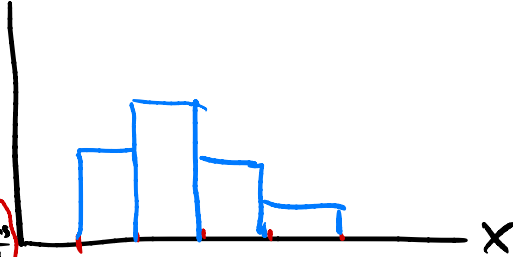
→ share of obs:  $N_k/N$

→ Stata: hist

R: hist()

ggplot::geom\_histogram()

→ density:  $\hat{f}(x) = \frac{N_k}{N} \times \frac{\text{number}}{\text{width}}$



Centered histogram

kernel density estimator

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

Stata: kdensity

R: ggplot::geom\_density()

