

Jose' Luis Preciado Areola

Theory

- Nonparametric algorithms to estimate stochastic production functions & frontiers
 - Bayesian
 - Frequentist
- Bayesian estimation of unobserved variables and functions (MCMC).

Industry

- Senior Data Scientist — ITD Innovation Lab @ Austin
- Revenue opt. — MC algorithms
— Math. Programming
- Tollway optimization modeler.

Now

- Predictive defect detection for Steel Industry
- Survival Analysis related to workforce rotation.
- Econometrics
- Wind Energy

Lecture 1: Algorithms and Inference

Statistics:

- Branch of math
- Related to prob. theory
- Extracting learnings from data.
- Finding a "true" signal from experience (and separate it from noise).

Algorithmic : • Like a pipeline that transforms data into more summarized learnings.

- Following a set of data processing steps to produce estimates.

Inferential : • How "good" is the outcome of the algorithm?
• Measuring the uncertainty around estimates.

Statistical
Fundamentals

- Statistics courses (Grad)

This course!

Statistical
Practice

- Online courses about OS, ML, ISD, etc.

Warmup Example

• The mean $\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$ summarizes data into one number.

• How precise is this number?

Denotes estimator $\hat{\sigma}_{\bar{x}}$ (standard error) = $\left[\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n(n-1)} \right]^{1/2}$
(\hat{se})

sample variance σ^2

$$= \left[\frac{\sigma^2}{n} \right]^{1/2} = \frac{\sigma}{\sqrt{n}}$$

Inference in regression

• A regression is an estimator of a conditional mean. In same way $\bar{y} = E(Y)$

$$\hat{y} = E(Y | X)$$

Height Gender

For the univariate case

Liver
Func.
Index

$$= \beta_0 + \beta_1 \cdot \text{Age (model)}$$

$$\rightarrow \underbrace{L_i}_{\text{LFI for patient } i} = \underbrace{\beta_0}_{\text{Age of patient } i} + \underbrace{\beta_1}_{\text{Age of patient } i} x_i + \underbrace{\epsilon_i}_{\text{error for observation } i}$$

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (\epsilon_i)^2 \quad (\text{Algorithm}) \quad \text{Minimize Least Squares Loss Function.}$$

$$\min_{\beta_0, \beta_1} \sum_{i=1}^n (L_i - \beta_0 - \beta_1 x_i)^2$$

By the way, your data for this problem, looks like:

$$\hat{\sigma}_{\epsilon LS(1)} = \left[\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n(n-2)} \right]^{1/2}$$

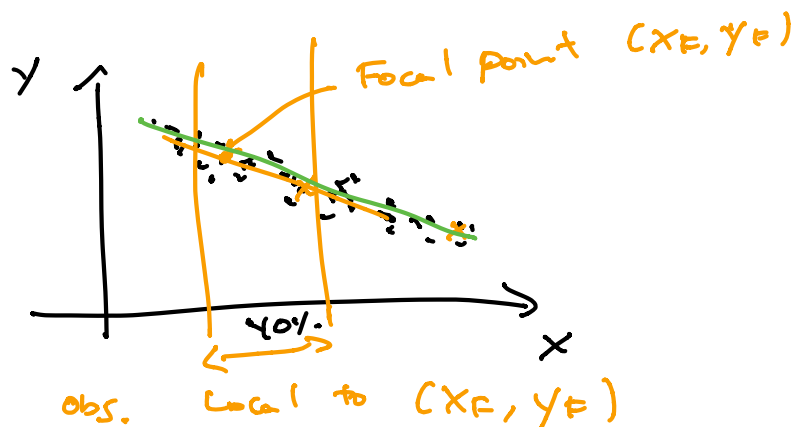
Because we're estimating 2 parameters from data.

Lowess (Local weighted Regression)

- Say we choose $f = 40\%$ of the data to be local to each focal point.

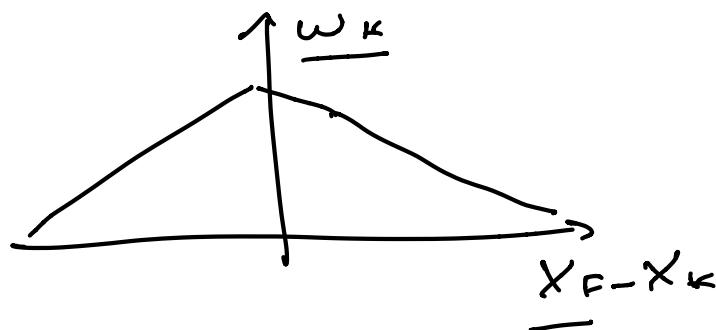
Step 1. Choose the N data points closest
in X to your focal point (x_F, y_F)

Your sample



Step 2. Compute a weighted linear regression for the focal point.

This means that behind the weights for each observation there's a weighting pattern.

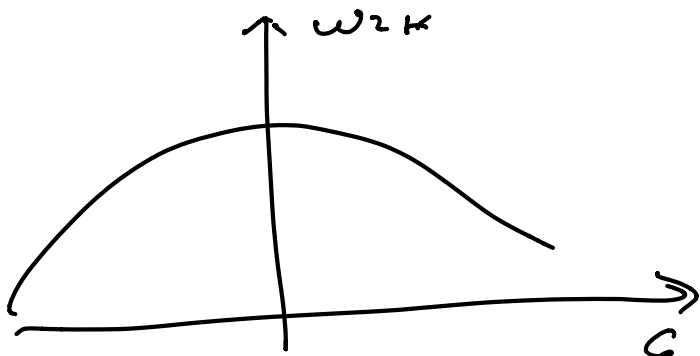


Weighting function

where x_k is any other observation.

Step 3. Solve ^{40%} $\min_{\beta_0, \beta_1} \sum_{k=1}^N (w_k \cdot e_k)^2$ for $i = 1, \dots, N$

Step 4. Apply a second weighting function depending on $e_{Fi} = y_F - \hat{y}_{F(-i)}$, multiply the weights & redo step 3.



"Give more weight to good predictions from previous step".

$$\xi_i(i) = YF - \hat{Y}_F(i)$$

Step 5. Repeat step 4 until satisfied.

Instead of one regression with N observations, we need to compute $N \cdot \boxed{\text{\# of iterations}}$ regressions with $\lceil N \cdot f \rceil$ observations each.

- There is no exact inference for lower estimates ".

How do we get inference then?

→ Bootstrapping → Sampling w/ replacement