Lecture Ten

High-Dimensional hinear Regression

1 Linear Regression: y = B + A+ + E+ y = B+ B+ B+2+ &+ Y== B+ B con 2xt + B sin =xt+Ex

2 Non Linean Regression:

y = B + B cos 27 + B cin 27 + E. RSS(f) = min = (4---)²

3 High-Dimension Linear Regression:

7 =

Change of Slope Model

California yearly popular at time year 1900 to 2024 y = log (cA population at time t)

100xB: percent rate

y = B + B + B (t - c) + Et.

the contraction of slope model:

$$t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta_2 \\
t \leq c : \text{ slope } \beta_1 + \beta$$

- (1) How to choose K?
- 2 Estimation of parameters (??)
- 3 (x large then there will be overtitting.)

We create a high-dimensional model by using all possible choices of C1, C21.

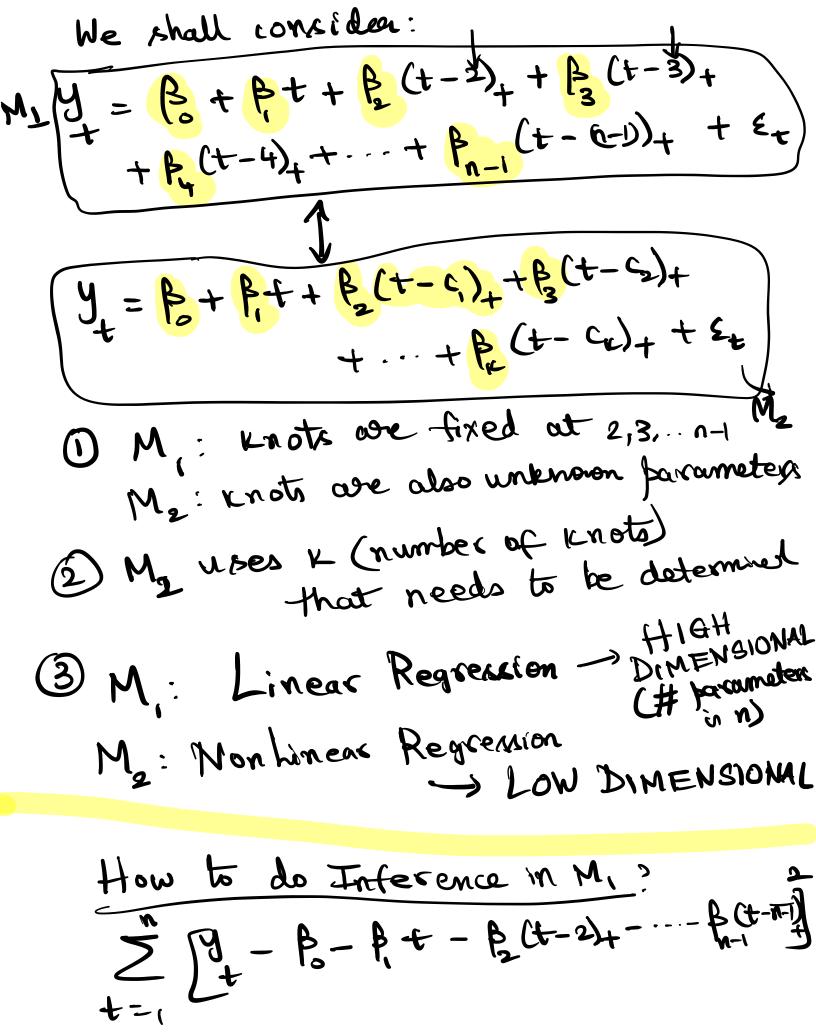
$$y_{+} = \beta_{3} + \beta_{1} + \beta_{2} + \beta_{3} + \beta_{4} + \beta_{5} + \beta_{5$$

$$(t-5.7)_{+} = (0.3)(t-5)_{+} + (0.7)(t-6)_{+}$$

for all $t \ge 5$

& all +>6 If +36, Ths = (0.3) (t-5) +(0.7) lhs = t-5.7, = +-5.7

(+-6)



Slight change in notation for M,:

B + B(t-1) + B (t-2) + + B(t-3) + ...+ B(t-1) + B = Y_{4} , t=1,...n $(t=1) \Rightarrow \beta_0 = Y_1$ More generally, $B_{i} = (y_{i+1} - y_{i}) - (y_{i} - y_{i-1})$ (izz) $y_t = \log P_t cA$ > population to Bo = Population in year 1 on log scale $P_{1} = y_{2} - y_{1} = ky \frac{P_{2}}{P_{1}} \approx \frac{P_{2} - P_{1}}{P_{1}}$ 100 x B = percent change in Population from year 1 to year 2 B2 = (3- 42) - (3- 41)

100x B = charge in percent charge between years 2 to 3 Use this model but change extimation 1) Frequentist -> Regularized
MLE >2 Bayesian -> Change the prior. An ièd Unif (-C, C) hange this