

Given data $(X_1, Y_1), \dots, (X_n, Y_n) \sim F$

Goal: Find β minimizing dist in high dimension

$$R(\beta) = \mathbb{E}[L(r_\beta(X), Y)] = \int L(r_\beta(x), y) dF$$

where $L(\hat{y}, y)$ is loss function

Since we don't know F , replace w. empirical CDF F_n

$$R_n(\beta) = \int L(r_\beta(x), y) dF_n = \frac{1}{n} \sum_{i=1}^n L(r_\beta(X_i), Y_i)$$

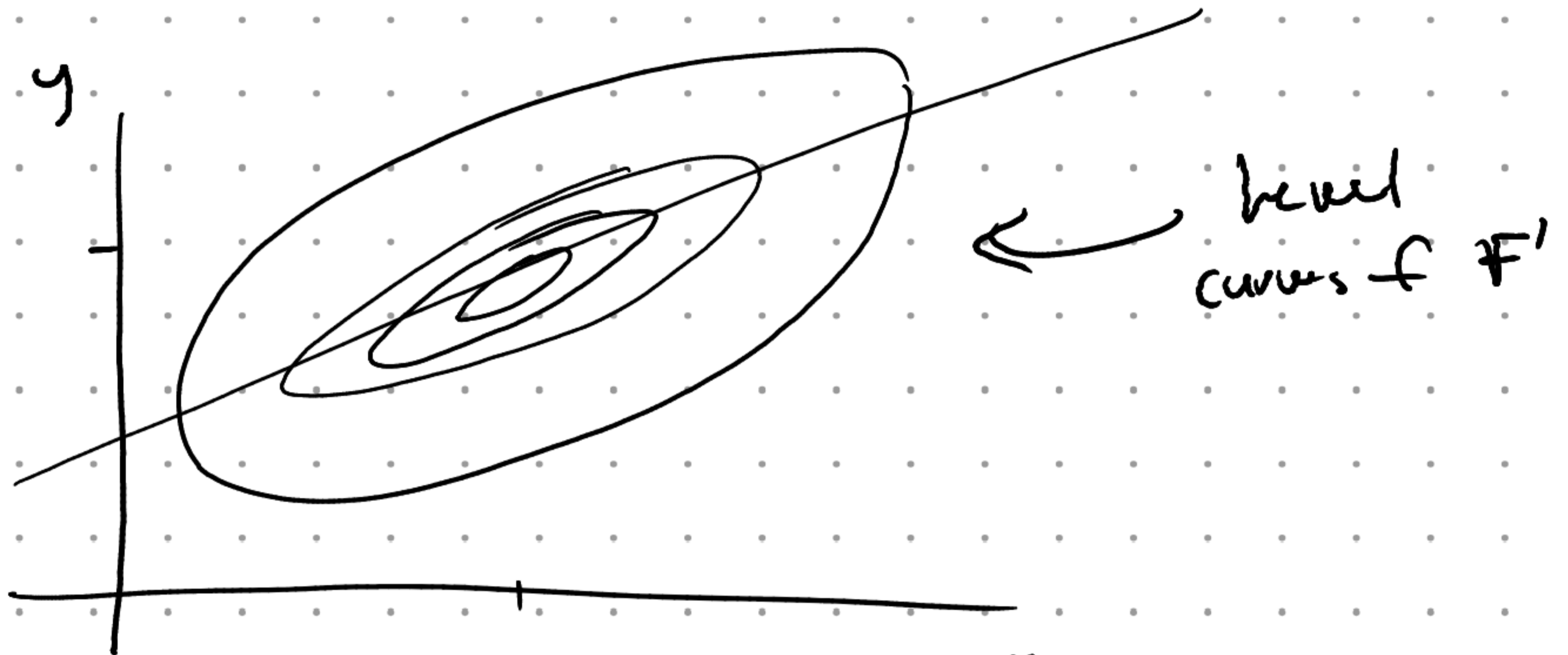
Now minimize $R_n(\beta)$ to find $\hat{\beta}_n$

How good is this estimate? IE what is $R(\hat{\beta}_n)$?

Ex

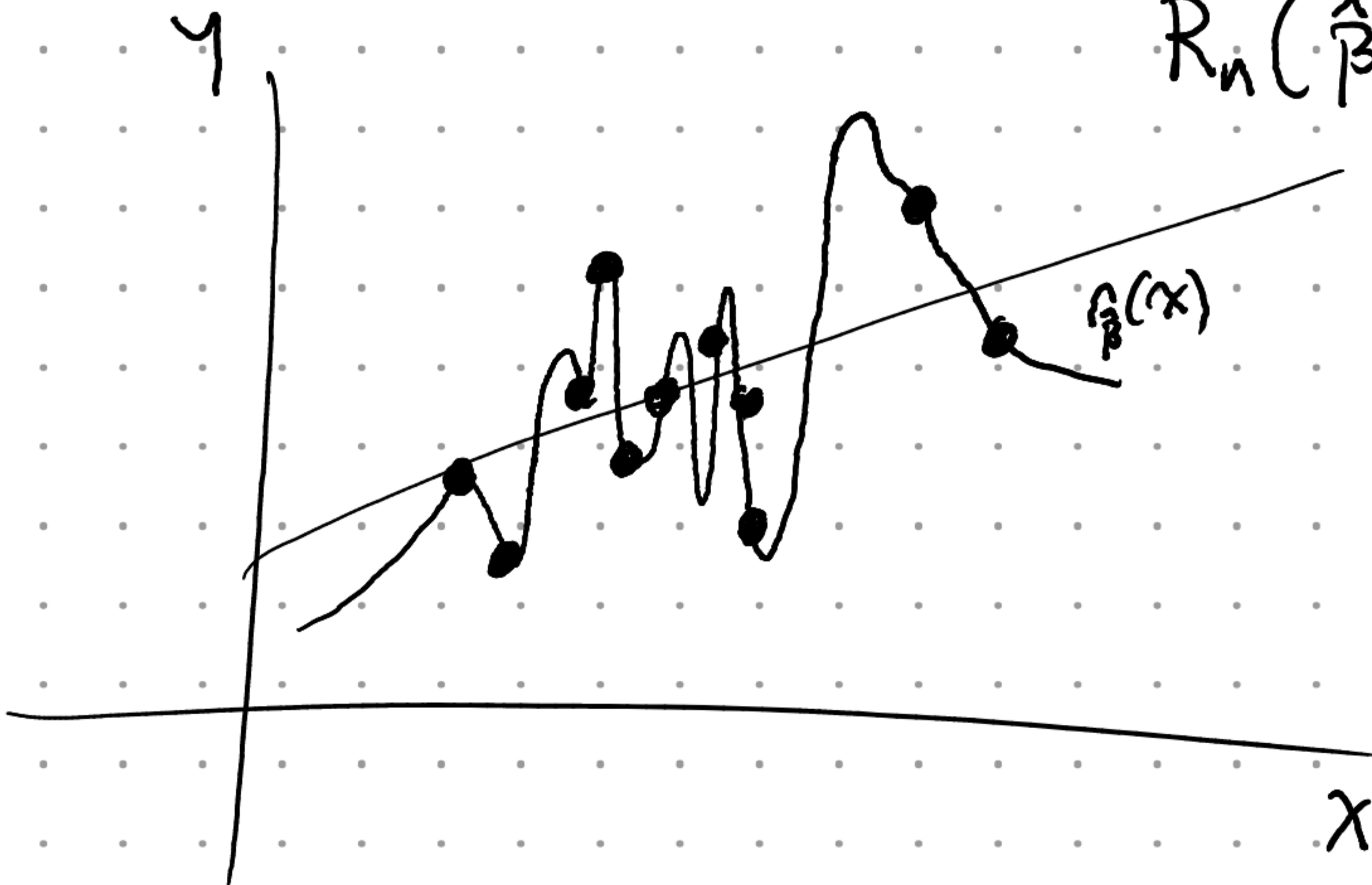
Suppose X, Y joint gaussian

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N \left(\begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{xy} & \sigma_y^2 \end{bmatrix} \right)$$



$$r(x) = \mathbb{E}[Y|X] = \min_r \mathbb{E}[(r(x) - Y)^2 | X]$$

$$R_n(\hat{\beta}) = 0$$



Intuitively, our model was too complex
so we overfit the data. If
we restricted to a simpler model,
then we would be less likely to
overfit. But simpler models may
not be expressive enough to represent
trends in the data.

Here, $R_n(\hat{\beta}) < R(\beta)$

This is not uncommon outcome.

To estimate $R(\beta)$, we can

use new data $(\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n, \tilde{Y}_n)$

and compute $\tilde{R}_n(\hat{\beta})$

New data often not available, so

split original data into two sets,

one for learning $\hat{\beta}$, one for evaluating
the quality.

Least Squares as MLE

Suppose X_1, \dots, X_n fixed set of points

and $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{iid}$$

Then $Y_i \sim N(\mu_i, \sigma^2)$, $\mu_i = \beta_0 + \beta_1 X_i$

$$L_n(\beta) \propto \prod_{i=1}^n f_{\beta}(Y_i) \quad \text{density for } Y_i$$

$$\propto \prod_{i=1}^n \exp\left(-\frac{(Y_i - \mu_i)^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - \mu_i)^2\right)$$

$$\max_{\beta} L_n(\beta) \leftrightarrow \min_{\beta} \sum_{i=1}^n (Y_i - \mu_i)^2$$

$$= \min_{\beta} \sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 X_i))^2$$