

# STOCHASTIC OPTIMIZATION

## AN EP VIEWPOINT

SP 23, STUDY GR.

## EXAMPLE I (MLE)

$$Y \sim p_{\theta_0}(x), \quad x \in \mathcal{X}.$$

$p_{\theta_0}$  is a density w.r.t  
a  $\sigma$ -finite measure.

$\theta_0 \in \Theta$  is unknown

Available data:  $Y_1, Y_2, \dots, Y_n$  iid

## EXAMPLE I (MLE) contd...

Notice that  $\forall \theta \in \Theta$

$$\mathbb{E} \log P_{\theta_0}(Y)$$

$$\geq \mathbb{E} \log P_{\theta}(Y)$$

— (1)

Why?

## EXAMPLE I (MLE) contd...

(1) implies that

$$\theta_0 \in \underset{\theta}{\text{argmin}} \mathbb{E} \log P_{\theta}(Y)$$

## EXAMPLE I (MLE) contd...

How to estimate  $\theta_0$ ?

Suppose  $Y_1, Y_2, \dots, Y_n \stackrel{iid}{\sim} P_{\theta_0}$

Solve:

$$\hat{\theta}_n \in \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log P_{\theta}(Y_i)$$

 max. likelihood estimator.

## EXAMPLE I (MLE) contd...

The alternate EP viewpoint.

$$\log P_{\theta}(Y) \equiv F(\theta, Y)$$

Where  $F: \Theta \times \mathcal{Y} \rightarrow \mathbb{R}$

## EXAMPLE I (MLE) contd...

The alternate EP viewpoint.

$$\theta_0 \in \arg \min_{\theta} \mathbb{E} [F(\theta, Y)]$$

$$\hat{\theta}_n \in \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n F(\theta, Y_i)$$

## EXAMPLE I (MLE) contd...

Main Questions:

(Q.1) Does  $\hat{\theta}_n \rightarrow \theta_0$  in some sense?

(Q.2) Is there a "CLT" on  $\hat{\theta}_n$ ?

(Q.3) What if  $\hat{\theta}_n$  can be solved only approximately?



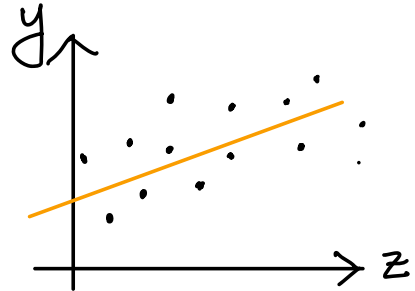
## EXAMPLE I (MLE) contd...

Food for Thought.

- (a) Construct an example where  $\Theta$  is finite dimensional.
- (b) Construct an example where  $\Theta$  is infinite dimensional.

## EXAMPLE II (Least Squares)

$(\tilde{Y}, Z)$



"response"

"covariate"

$$\tilde{Y} \in \mathbb{R} ; Z \in \mathcal{L}.$$

Model:

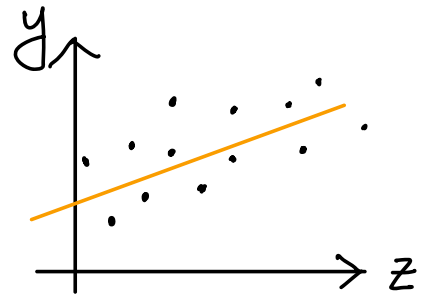
$$\tilde{Y} = \theta_0(Z) + W$$

"error"

"true reg. fn."

## Least Squares ...

Observe data



$$(\tilde{Y}_i, Z_i), i = 1, 2, \dots, n$$

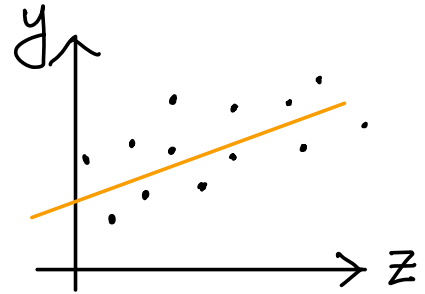
-  $(\tilde{Y}_i, Z_i)$  iid.

-  $W_1, W_2, \dots$  independent with  
 $E[W_i] = 0, \text{Var}(W_i) \leq \sigma_0^2 < \infty.$

-  $\theta_0 : \mathcal{L} \rightarrow \mathbb{R}; \theta_0 \in \Theta$

↓  
"regression class"

## Least Squares ...



$$\hat{\theta}_n \in \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \left( \theta(z_i) - \tilde{y}_i \right)^2$$

↓

"least squares estimator"

## EXAMPLE I (MLE) contd...

The alternate EP viewpoint.

$$(\theta(Z) - \tilde{Y})^2 \equiv F(\theta, Y)$$

Where  $F: \theta \times y \rightarrow \mathbb{R}$

$$\theta_0 \in \arg \min_{\theta} \mathbb{E} [F(\theta, Y)]$$

$$\hat{\theta}_n \in \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n F(\theta, Y_i)$$

## Least Squares ...

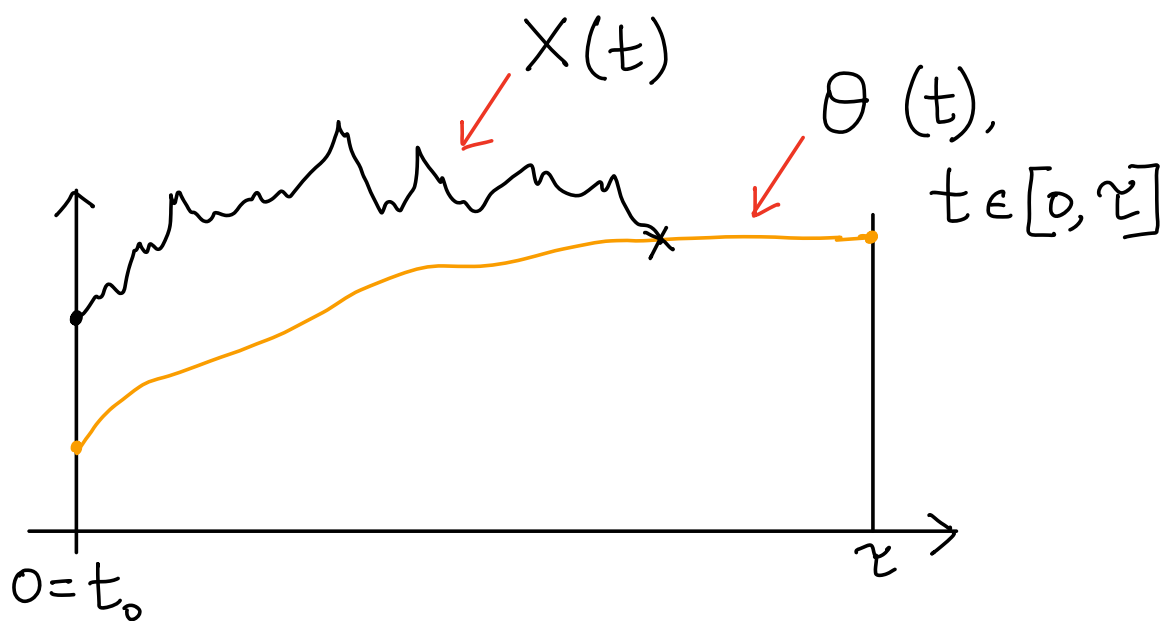
- (Q.1) Does  $\hat{\theta}_n \rightarrow \theta_0$  in some sense?
- (Q.2) Is there a "CLT" on  $\hat{\theta}_n$ ?
- (Q.3) What if  $\hat{\theta}_n$  can be solved only approximately?

Least Squares ...

Food for thought.

Linear regression, logistic regression, isotonic regression, deep learning fall in this framework. Convince yourself...

## EXAMPLE III (Limit Order Book)



$\{X(t), t \in [0, \infty)\}$  "asset price process"

$\theta(\cdot)$  : "execution boundary"

$[0, \tau]$  : "horizon"



LOB ...

$$T := \inf \{ t : \theta(t) \geq X(t) \}$$

Purchase price:

$$C^*(\theta) := \begin{cases} \theta(T) & T \leq \tau \\ m(\tau) & T > \tau \end{cases}$$

↙  
"market price"

LOB ...

$$\min. \mathbb{E} [c^*(\theta)]$$

$$\text{s.t. } \text{Var}(c^*(\theta)) \leq \gamma.$$

$$\theta \in \Theta$$

notation.

$$\hat{\theta}_n \in \arg \min \frac{1}{n} \sum_{j=1}^n c^*(\theta, Y_j)$$

$$\text{s.t. } \frac{1}{n} \sum_{j=1}^n (c^*(\theta, Y_j) - \bar{c}^*(\theta))^2$$

$$\theta \in \Theta.$$

## LOB ...

- (Q.1) Does  $\hat{\theta}_n \rightarrow \theta_0$  in some sense?
- (Q.2) Is there a "CLT" on  $\hat{\theta}_n$ ?
- (Q.3) What if  $\hat{\theta}_n$  can be solved only approximately?

# Stochastic Optimization (unconstrained)

$$\min_{\theta_0} f(\theta) = \mathbb{E}[F(\theta, Y)]$$

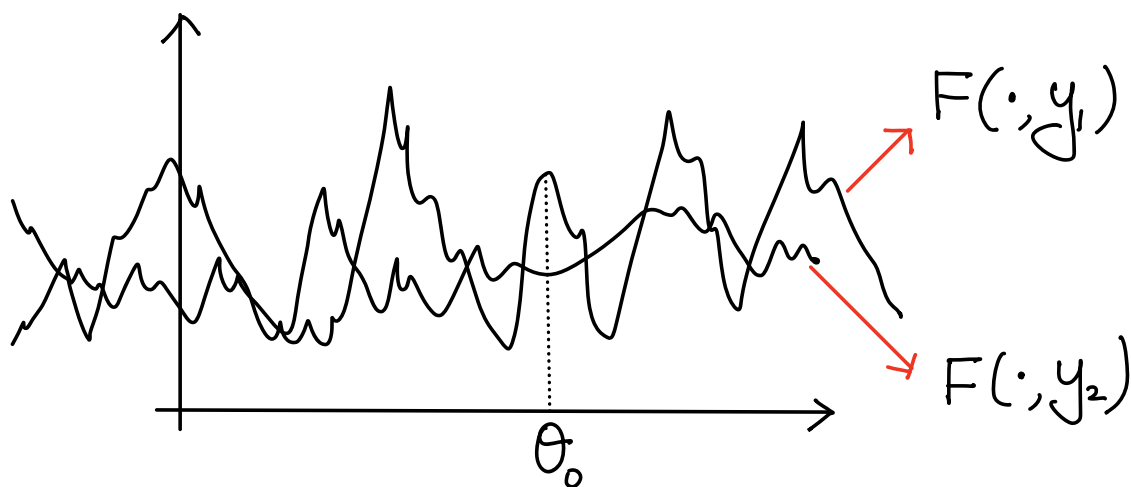
$$\text{s.t. } \theta \in \Theta.$$

$$f_0: \Theta \rightarrow \mathbb{R}, \quad F: \Theta \times \mathcal{Y} \rightarrow \mathbb{R}.$$

$F$  and possibly some of its "derivatives" are observable.

"Function Class"  $\mathcal{F}$ .

This is the class of "random variables" labeled by  $\theta$ .



$$\Theta \equiv \mathbb{R}$$

$$\mathcal{F} := \{ F(\theta, Y), \theta \in \Theta \}.$$

We will see that the  
"richness" of  $\mathcal{F}$  will play  
a very important role.

Recall Stochastic Optimization  
(unconstrained)

$$\min. f(\theta) = \mathbb{E}[F(\theta, Y)]$$

$$\text{s.t. } \theta \in \Theta.$$

$$f : \Theta \rightarrow \mathbb{R}, \quad F : \Theta \times \mathcal{Y} \rightarrow \mathbb{R}.$$

$F$  and possibly some of its  
"derivatives" are observable.

## M-Estimation

$$\min. \frac{1}{n} \sum_{j=1}^n F(\theta, Y_j)$$

$$\text{s.t. } \theta \in \Theta$$

$$\min E[F(\theta, Y)]$$

$$\text{s.t. } \theta \in \Theta$$

## Z-Estimation

solve:

$$\frac{1}{n} \sum_{j=1}^n G_1(\theta, Y_j) = 0$$

solve:

$$E[G_1(\theta, Y)] = 0$$



A solution  $\hat{\theta}_n$  to the  
M-Estimation problem is  
an M-Estimator.

Many, many statistical  
estimators are M-Estimators.

Let's look at examples.

# M-Estimation EXAMPLE I

Recall MLE.

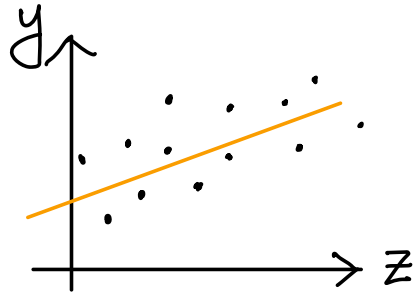
$$X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} P_{\theta^*}.$$

$$\min_{\theta \in \Theta} -\frac{1}{n} \sum_{j=1}^n \underbrace{\log P_{\theta}(X_j)}_{F(\theta, Y_j)}$$

# M-Estimation EXAMPLE II

---

Regression:



Data:  $(X_i, R_i), i=1, 2, \dots, n$

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{j=1}^n (\theta(X_i) - R_i)^2$$

↓  
"space of reg.  
functions"

↓  
 $F(\theta, Y_i)$

# M-Estimation EXAMPLE III

Deep Learning

## Z-Estimation EXAMPLE IV

Mean Estimation.

set  $G(\theta, Y) \equiv Y - \theta$

solve:

$$\frac{1}{n} \sum_{j=1}^n (Y_j - \theta) = 0$$

$$\Rightarrow \hat{\theta}_n = \bar{Y}$$

$$\theta^* = ?$$

## Z-Estimation EXAMPLE $\nabla$

Mean-variance estimation

set  $G_1(\theta, Y) = \begin{pmatrix} Y - \theta_1 \\ (Y - \theta_1)^2 - \theta_2 \end{pmatrix}$

solve:

$$\frac{1}{n} \sum_{j=1}^n \begin{pmatrix} Y_j - \theta_1 \\ (Y_j - \theta_1)^2 - \theta_2 \end{pmatrix} = 0$$

$$\Rightarrow \hat{\theta}_n = \begin{pmatrix} \bar{Y} \\ \hat{\sigma}_n^2 \end{pmatrix} \cdot \theta^* = ?$$

## Z-Estimation EXAMPLE VI

Ratio Estimation.

set  $G_1(\theta, Y) \equiv Y_1 - \theta Y_2$

solve:

$$\frac{1}{n} \sum_{j=1}^n (Y_{1j} - \theta Y_{2j}) = 0$$

$$\Rightarrow \hat{\theta}_n = \bar{Y}_1 / \bar{Y}_2$$

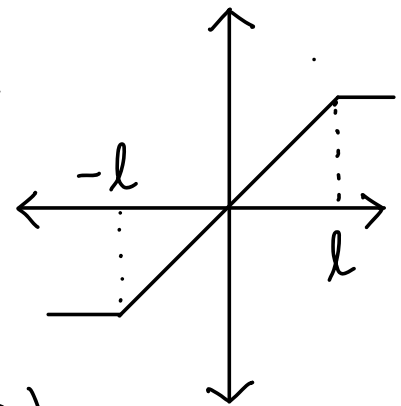
$$\theta^* = ? \quad \theta^* = \frac{E[Y_1]}{E[Y_2]}$$

## Z-Estimation EXAMPLE VII

Robust Estimation (Huber, 1964)

set

$$h_l(x) = \begin{cases} -l & x < -l \\ x & |x| \leq l \\ l & x > l \end{cases}$$



$$G_1(\theta, Y) \equiv h(Y - \theta)$$

solve:

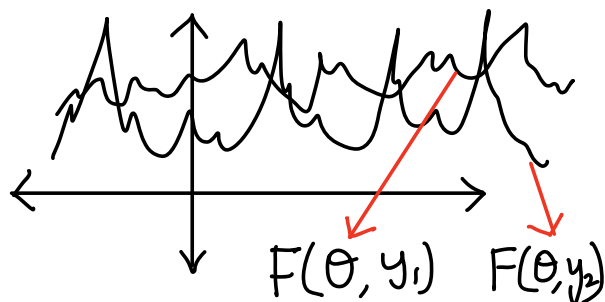
$$\frac{1}{n} \sum_{j=1}^n G_1(\theta, Y_j) = 0$$

( $\theta^*$  satisfies  $\mathbb{E}[h'(Y - \theta)] = 0$ .)



So What? (Heuristic Preview)

So What?



EP formalizes and vastly generalizes by looking at

$$\{ F(\theta, Y), \theta \in \Theta \}$$

(or  $\{ G(\theta, Y), \theta \in \Theta \}$ ).

Let's begin our formal  
study ...

# NOTATION

Suppose  $Y_1, Y_2, \dots$  are  
independent copies of  $Y$   
having distribution  $P$   
on  $(X, \mathcal{A})$ .

$P_n \rightarrow$  "empirical distbn."

$$(P_n(A) = \frac{1}{n} \sum_{j=1}^n \mathbb{I}_A(Y_j))$$


$$\begin{aligned} \mathbb{P}_n F_\theta &\equiv P_n F(\theta, \cdot) = \frac{1}{n} \sum_{j=1}^n F(\theta, Y_j) \\ &\quad \searrow \qquad \qquad \qquad \swarrow \\ &\qquad \qquad \qquad =: \int F_\theta dP_n. \end{aligned}$$

"expectation of  $F(\theta, \cdot)$  under  
the empirical measure"

Similarly,

$$\begin{aligned} PF_{\theta} &\equiv PF(\theta, \cdot) \\ &= E[F(\theta, Y)]. \end{aligned}$$

$$\mathcal{F} := \{F(\theta, Y), \theta \in \Theta\}$$



class of measurable functions  
labeled by  $\theta$ .

## Uniform Norm

$$\|P_n - P\|_{\infty} =$$

$$\sup_{\theta \in \Theta} |P_n F_{\theta} - P F_{\theta}|$$



Uniform Law of Large Numbers  
(ULLN) w.r.t  $\mathcal{F}$ .

$$\lim_{n \rightarrow \infty} \left\| P_n - P \right\|_{\mathcal{F}} = 0 \text{ P a.s.}$$

$\mathcal{F}$  will be called the Glivenko-Cantelli class.

We will write:

$$\int F_{\theta} d(P_n - P)$$

to mean

$$P_n F_{\theta} - P F_{\theta}$$

$$= \frac{1}{n} \sum_{j=1}^n F(\theta, Y_j) - \mathbb{E}[F(\theta, Y)]$$

$$\left\{ \sqrt{n} (P_n F_\theta - P F_\theta), \theta \in \Theta \right\}$$



"empirical process"

