

- Empirical risk minimization
- Stochastic approximation of the gradient
  - unbiasedness
  - role of variance
- Lower bound  $\underline{g}_t^\top (\underline{x}_t - \underline{x}^*)$  ?

Wed. 20<sup>th</sup> 18:00 - 19:00 hrs (TA Session)

Recall gradient descent:

$$\underline{x}_{-t+1} = \underline{x}_{-t} - \eta_t \underline{g}_t \quad ; \quad \underline{g}_t \in \partial f(\underline{x})$$

$\nabla f(\underline{x}_t)$

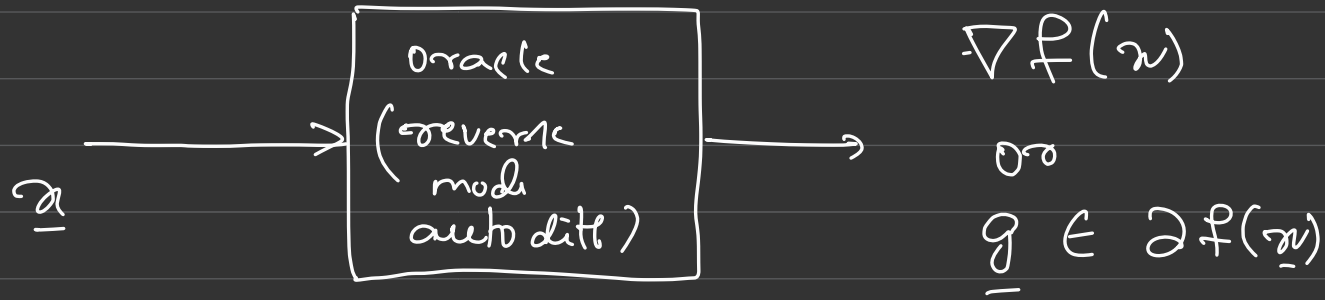
Motivation:

- $\underline{g}_t$  may be expensive to compute
- complete gradient  $\underline{g}_t$  may not be available

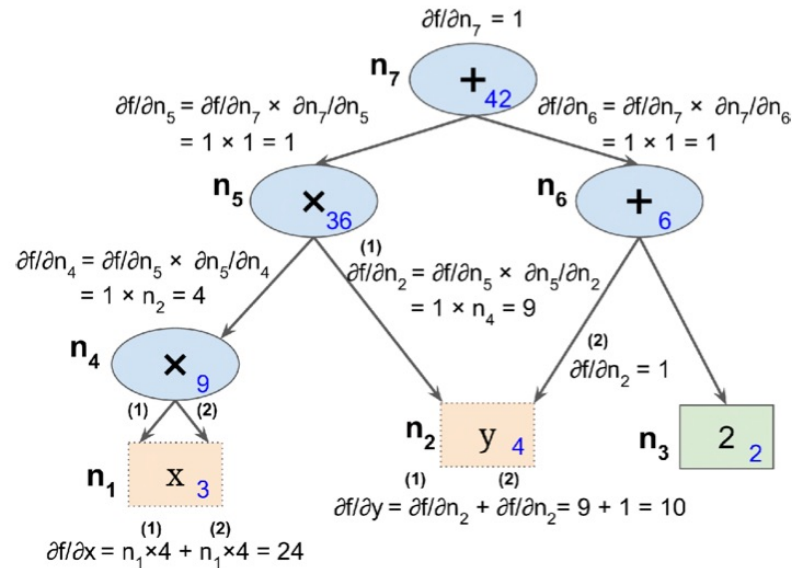
So we use an stochastic or approximate version of  $\underline{g}_t \in \partial f(\underline{x})$

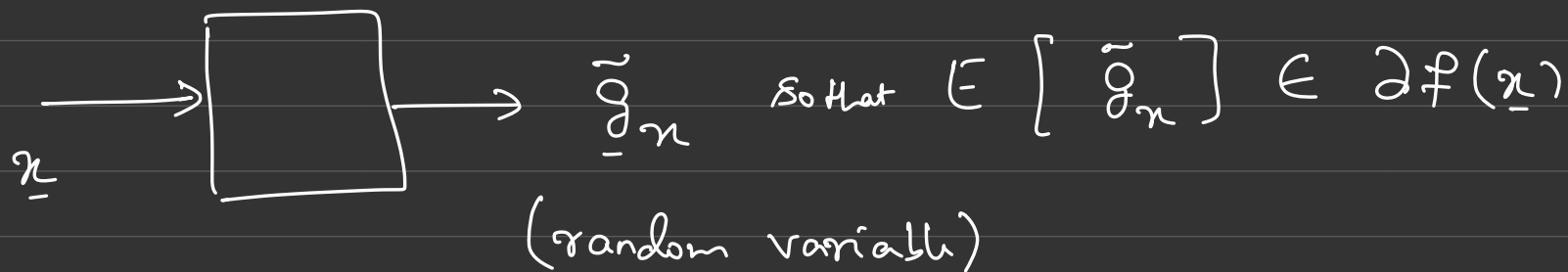
minimize  $f(x)$

s. to  $x \in C$



$$f(x, y) = x^2y + y + 2$$





Example:

①  $\tilde{g}_{\underline{x}} = \nabla f(\underline{x}) + \underline{w}$  ;  $\underline{w}$  is zero-mean noise

$$E[\tilde{g}_{\underline{x}}] = E[\nabla f(\underline{x}) + \underline{w}] = \nabla f(\underline{x})$$

② Random coordinate descent:

$\underline{x} \rightarrow$    $\rightarrow \tilde{g}_{\underline{x}} = \begin{bmatrix} 0 \\ \vdots \\ \partial f / \partial x_j \\ \vdots \\ 0 \end{bmatrix} \cdot d$  ;  $\underline{x} \in \mathbb{R}^d$

$j \sim \text{Unif}(1, \dots, d)$

$\nabla f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_d} \end{bmatrix}$

$$E_j[\tilde{g}_{\underline{x}}] = \sum_{i=1}^d \frac{1}{d} \cdot \begin{bmatrix} 0 \\ \vdots \\ \partial f / \partial x_i \\ \vdots \\ 0 \end{bmatrix} = \nabla f(\underline{x})$$

### ③ Stochastic programming:

$$\underset{\underline{x} \in C}{\text{minimize}} \quad F(\underline{x}) = \underbrace{E [f(\underline{x}; \xi)]}_{\substack{\text{Expected risk} \\ \text{Population risk}}}$$

- $\xi$  : random in the problem
- If  $f(\underline{x}; \xi)$  is convex for every  $\xi$ , then  $F(\underline{x})$  is convex.

## Empirical risk minimization:

Let  $\{\underline{a}_i, y_i\}_{i=1}^n$  be  $n$  random data samples

ERM: minimize  $\underline{x}$   $f(\underline{x}) = \frac{1}{n} \sum_{i=1}^n f(\underline{x}; \{\underline{a}_i, y_i\})$

Empirical risk

Regression problem: (more generally, any supervised learning)

$$f(\underline{x}; \{\underline{a}_i, y_i\}) = (\underline{a}_i^T \underline{x} - y_i)^2$$

prediction/hypothesis:  $\underline{a}_i^T \underline{x}$

minimize expected loss: draw  $j \sim \text{unif}(1, 2, \dots, n)$ , then

$$E_j [f(\underline{x}; \{\underline{a}_j, y_j\})] = \sum_{j=1}^n f(\underline{x}; \{\underline{a}_j, y_j\}) \cdot \frac{1}{n}$$

$$\text{minimize}_{\underline{x} \in \mathcal{C}} f(\underline{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\underline{x})$$

ERM:  $n$  is the number of data points

Gradient descent:

$$\begin{aligned} \underline{x}_{t+1} &= \underline{x}_t - \eta \nabla \left( \frac{1}{n} \sum_{i=1}^n f_i(\underline{x}_t) \right) \\ &= \underline{x}_t - \eta \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(\underline{x}_t) \end{aligned}$$

When  $n$  is large, computing  $\nabla f(\underline{x})$  is expensive (full pass over the data  $\rightarrow$  memory issues)

$$\begin{array}{c} \underline{x} \longrightarrow \boxed{\begin{array}{l} j \sim \text{unif} \\ (1, \dots, n) \\ \text{oracle} \end{array}} \longrightarrow \tilde{g}_{\underline{x}} = \nabla f_j(\underline{x}) \\ E_j \left[ \nabla f_j(\underline{x}) \right] = \sum_{j=1}^n \nabla f_j(\underline{x}) \cdot \frac{1}{n} = \nabla f(\underline{x}) \end{array}$$

# Stochastic gradient descent / Stochastic approximation:

$$\underline{x}_{t+1} = \underline{x}_t - \eta_t \tilde{g}(\underline{x}_t; \xi_t)$$

where  $\tilde{g}(\underline{x}_t; \xi_t)$  is an unbiased estimate of  $\nabla F(\underline{x}_t)$ , i.e.,

$$E[\tilde{g}(\underline{x}_t; \xi_t)] = \nabla F(\underline{x}_t)$$

- Stochastic algorithm for finding a critical point  $\underline{x}$  that obeys  $\nabla F(\underline{x}) = 0$

or for finding roots of  $G(\underline{x}) = E[\tilde{g}(\underline{x}, \xi)]$



# SGD for ERM

minimize  $\underline{x}$   $f(\underline{x}) = \frac{1}{n} \sum_{i=1}^n f(\underline{x}; \{a_i, y_i\})$

For  $t = 0, 1, \dots$  do

Choose  $i_t$  uniformly at random

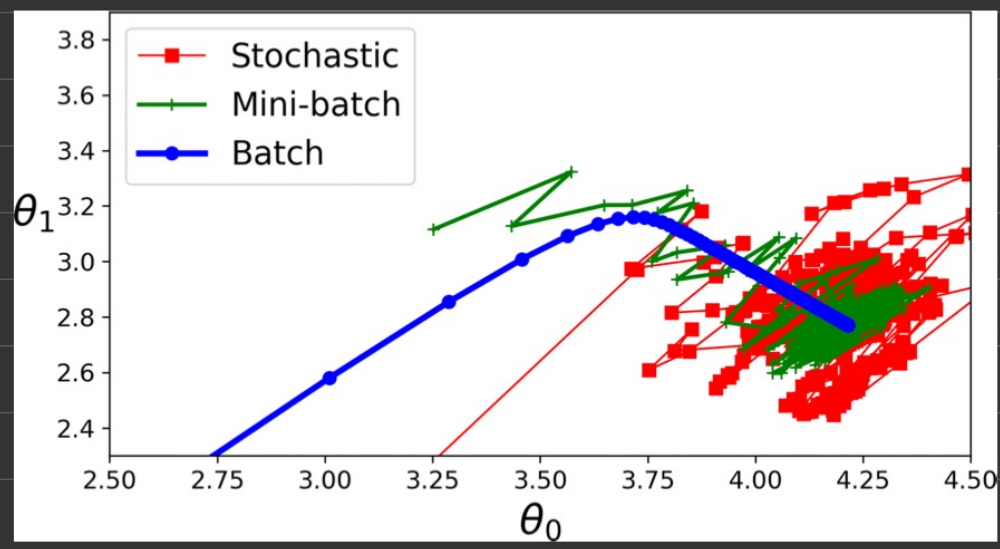
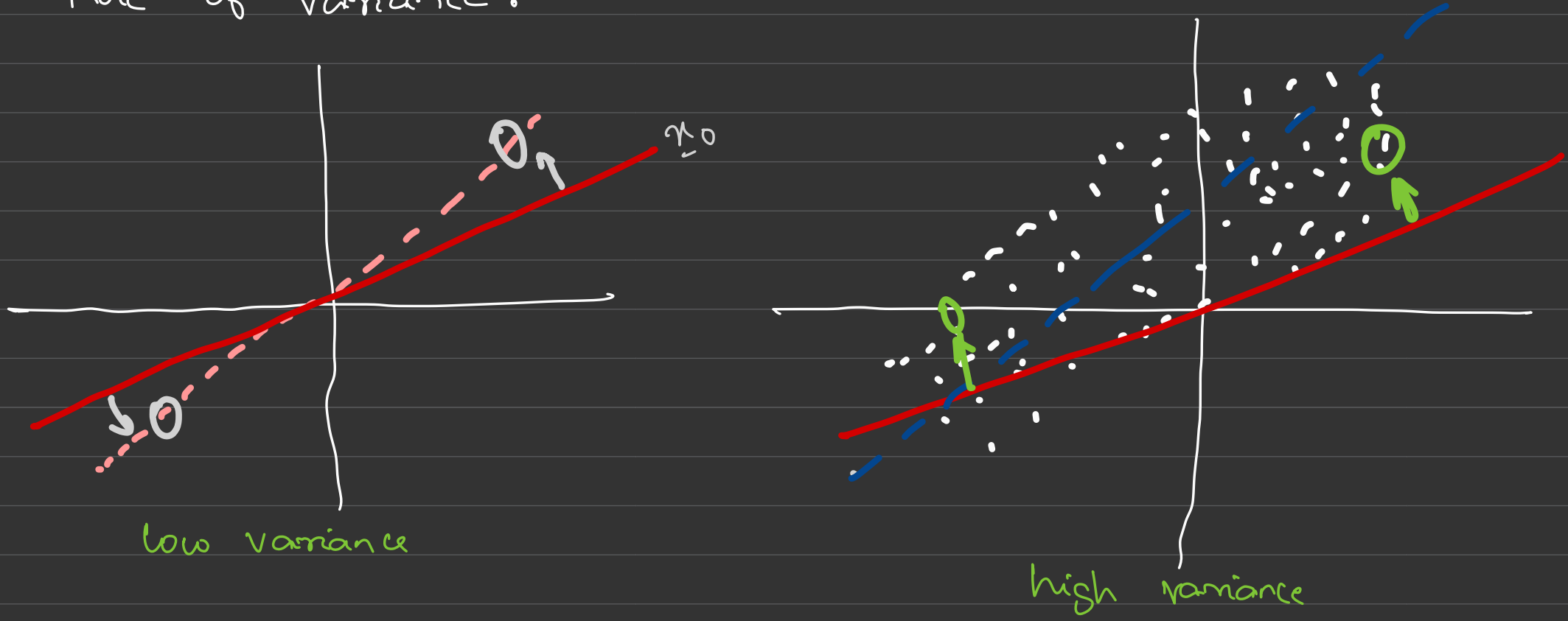
$$\underline{x}_{t+1} = \underline{x}_t - \eta_t \nabla_{\underline{x}} f_{i_t}(\underline{x}; \{a_i, y_i\})$$

$$y_i = a_i^T x$$

+ Exploits data more efficiently than batch methods

+ Fast initial improvement with low per-iteration cost (data usually has a lot of redundancy)

# Role of variance?



## Unbiasedness and the vanilla analysis:

Recall: In gradient descent, we could lower bound

$$\underline{g}_t^\top (\underline{x}_t - \underline{x}^*) \geq f(\underline{x}_t) - f(\underline{x}^*)$$

but now we cannot as  $\tilde{\underline{g}}_t$  may be far from being the true gradient.

• So inequality  $f(\underline{x}_t) - f(\underline{x}^*) \leq \tilde{\underline{g}}_t^\top (\underline{x}_t - \underline{x}^*)$

(from convexity) may not hold.

We have 
$$\mathbb{E} \left[ \underline{g}_t \mid \underline{x}_t = \underline{x} \right] = \frac{1}{n} \sum_{j=1}^n \nabla f_j(\underline{x}) = \nabla f(\underline{x})$$

Conditional expectation of  $\underline{g}_t$   
given the event  $\{\underline{x} = \underline{x}_t\}$ .

$$\forall \underline{x} \in \mathbb{R}^d$$

$$\Rightarrow \mathbb{E} \left[ g_t^\top (\underline{x} - \underline{x}^*) \mid \underline{x}_t = \underline{x} \right] =$$

$$\mathbb{E} \left[ g_t^\top \mid \underline{x}_t = \underline{x} \right] (\underline{x} - \underline{x}^*) = \nabla f^\top(\underline{x}) (\underline{x} - \underline{x}^*)$$

- $\{ \underline{x}_t = \underline{x} \}$  can occur only for  $\underline{x}$  in finite set  $X$

$$\mathbb{E} \left[ g_t^\top (\underline{x}_t - \underline{x}^*) \right]$$

$$= \sum_{\underline{x} \in X} \mathbb{E} \left[ g_t^\top (\underline{x} - \underline{x}^*) \mid \underline{x}_t = \underline{x} \right] \text{prob}(\underline{x}_t = \underline{x})$$

$$= \sum_{\underline{x} \in X} \nabla f^\top(\underline{x}) (\underline{x} - \underline{x}^*) \text{prob}(\underline{x}_t = \underline{x})$$

$$= \mathbb{E} \left[ \nabla f^\top(\underline{x}_t) (\underline{x}_t - \underline{x}^*) \right]$$

$$\begin{aligned}\Rightarrow \mathbb{E} \left[ \tilde{g}_t^\top (\underline{x}_t - \underline{x}^*) \right] &= \mathbb{E} \left[ \nabla f^\top (\underline{x}_t) (\underline{x}_t - \underline{x}^*) \right] \\ &\geq \mathbb{E} \left[ f(\underline{x}_t) - f(\underline{x}^*) \right]\end{aligned}$$

So the lower bound holds in expectation.