


Chapter 6

- Motivating example: mean estimation.

1. expectation: $E[X] \leq \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$

$$\bar{x} \rightarrow E[X] \text{ as } N \rightarrow \infty$$

2. Calculate:

directly: have to wait for a long time.

$$\text{iteration: } w_{k+1} = \frac{1}{k} \sum_{i=1}^k x_i \quad k = 1, 2, 3, \dots$$

$$w_k = \frac{1}{k-1} \sum_{i=1}^{k-1} x_i \quad k = 1, 2, 3, \dots$$

$$\begin{aligned} w_{k+1} &= \frac{1}{k} \sum_{i=1}^k x_i = \frac{1}{k} \left(\sum_{i=1}^{k-1} x_i + x_k \right) \\ &= \frac{1}{k} ((k-1)w_k + x_k) \\ &= w_k - \frac{1}{k}(w_k - x_k) \end{aligned}$$

Iteration algorithm: $w_{k+1} = w_k - \frac{1}{k}(w_k - x_k)$ 

Furthermore, general expression:

$$w_{k+1} = w_k - \alpha_k(w_k - x_k)$$

where $\alpha_k > 0$ and satisfy some mild condition

2. Robbins - Monro algorithm

1. Stochastic approximation (SA)

① SA refers to a broad of stochastic iterative algorithm solving root finding or optimization problem.

② SA does not require to know the expression of the objective function

③ Stochastic gradient descent (SGD)
is a special form of RM algorithm

2. find root of $g(w) = 0$

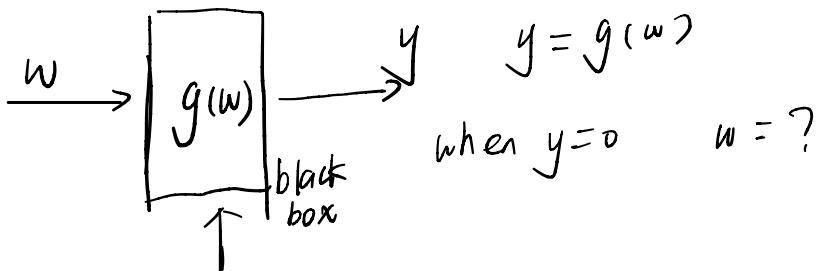
example: $\underbrace{g(w)}_{\text{or } g(w) - c} = 0$

or $g(w) = c \Rightarrow \underbrace{g(w) - c}_{g'(w)} = 0$

If $g(w)$ is known: use numerical algorithm

If $g(w)$ is unknown?

example: such as a artificial network



The Robbins - Monroe algorithm can solve

this problem: $w_{k+1} = w_k - \alpha_k \bar{g}(w_k, \eta_k)$ $k = 1, 2, \dots$

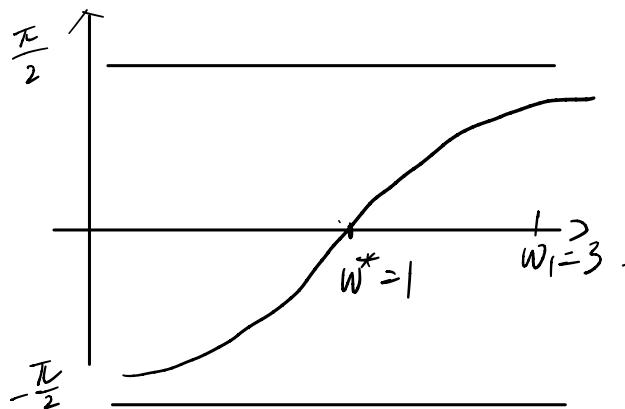
$$\bar{g}(w_k, \eta_k) = g(w_k) + \eta_k \quad \text{noisy observation}$$

need data

3. Convergence properties -

an illustrative example:

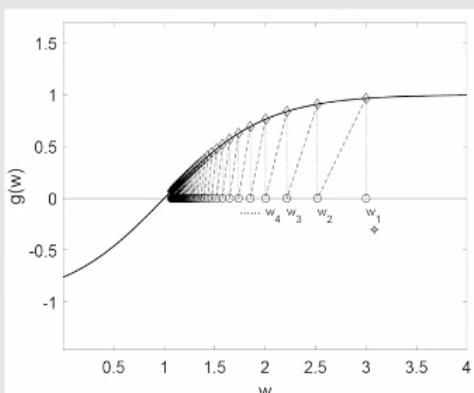
$$\cdot g(w) = \tanh(w-1)$$



$$w_1 = 3 \quad a_k = \frac{1}{k} \quad \eta_k \equiv 0 \text{ (no noise)}$$

$$w_{k+1} = w_k - a_k g(w_k)$$

Result: w_k converges to the true root $w^* = 1$.



不动点迭代
牛顿迭代

Intuition: w_{k+1} is closer to w^* than w_k .

- When $w_k > w^*$, we have $g(w_k) > 0$. Then, $w_{k+1} = w_k - a_k g(w_k) < w_k$ and hence w_{k+1} is closer to w^* than w_k .
- When $w_k < w^*$, we have $\underline{g(w_k)} < 0$. Then, $w_{k+1} = w_k - a_k g(w_k) > w_k$ and w_{k+1} is closer to w^* than w_k .

平滑数学推导：

Theorem (Robbins-Monro Theorem)

1) $0 < c_1 \leq \nabla_w g(w) \leq c_2$ for all w :

2) $\sum_{k=1}^{\infty} a_k = \infty$ and $\sum_{k=1}^{\infty} a_k^2 < \infty$:

3) $E[\eta_k | \mathcal{H}_k] = 0$ and $E[\eta_k^2 | \mathcal{H}_k] < \infty$;

where $\mathcal{H}_k = \{w_k, w_{k-1}, \dots\}$ then w_k converge with probability 1 to the root w^* satisfying $g(w^*) = 0$

依赖率收敛？

Explanation:

$\nabla g(w)$ 有界且大正， $g(w) = 0$ exist and unique.

The gradient is bounded from the above.

$$\min_w J(w)$$

$$g(w) = \nabla_w J(w) = 0$$

$$0 < c_1 \leq \nabla_w g(w) \leq c_2 \text{ for all } w.$$

$J(w)$ 是一个凸函数
convex.

$$\sum_{k=1}^{\infty} a_k = \infty \text{ and } \sum_{k=1}^{\infty} a_k^2 < \infty$$

The condition of $\sum_{k=1}^{\infty} a_k^2 < \infty$ ensure that a_k converges to zero as $k \rightarrow \infty$

The condition of $\sum_{k=1}^{\infty} a_k = \infty$ ensure that a_k do not converge to zero too fast.

$$E[\eta_k | H_k] = 0 \text{ and } E[\eta_k^2 | H_k] < \infty$$

A special yet common case is that $\{\eta_k\}$ is an iid stochastic sequence satisfying $E[\eta_k] = 0$ and $E[\eta_k^2] < \infty$

The observation error η_k is not required to be Gaussian.



$$\textcircled{1} \quad \sum_{k=1}^{\infty} a_k^2 < \infty \Rightarrow a_k \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$w_{k+1} - w_k = - \underbrace{a_k \tilde{g}(w_k, \eta_k)}_{\leftarrow k \rightarrow \infty} \quad w_{k+1} - w_k \rightarrow 0$$

$$\textcircled{2} \quad \sum_{k=1}^{\infty} a_k = \infty$$

$$\begin{cases} w_2 = w_1 - a_1 \tilde{g}(w_1, \eta_1) \\ w_3 = w_2 - a_2 \tilde{g}(w_2, \eta_2) \\ \vdots \\ w_{k+1} = w_k - a_k \tilde{g}(w_k, \eta_k) \end{cases}$$

$$\text{从 } w_2, w_1 - w_\infty = \sum_{k=1}^{\infty} a_k \tilde{g}(w_k, \eta_k)$$

$$w_\infty = w_1 - \sum_{k=1}^{\infty} a_k \tilde{g}(w_k, \eta_k)$$

$$w_\infty \rightarrow w^*$$

$$\sum_{k=1}^{\infty} a_k < \infty \text{ then } \sum_{k=1}^{\infty} a_k \tilde{g}(w_k, \eta_k) \text{ bounded.}$$

One typical sequence: $a_k = \frac{1}{k}$

It holds that: $\lim_{n \rightarrow \infty} \left(\sum_{k=1}^n \frac{1}{k} - \ln n \right) = \underline{\text{Euler-Mascheroni constant.}}$

Euler-Mascheroni

It is notable that:

$$\sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6} < \infty$$

Consider a Function:

$$g(w) = w - E[X]$$

The observation:

$$g(w, x) = w - x$$

$$\begin{aligned}\hat{g}(w, \eta) &= w - x = w - x + E[X] - E[X] \\ &= (w - E[X]) + \frac{(E[X] - x)}{\eta} \\ &= g(w) + \eta\end{aligned}$$

The RM for $\hat{g}(x) = 0$

$$w_{k+1} = w_k - \alpha_k \hat{g}(w_k, \eta_k) = w_k - \alpha_k (w_k - x_k)$$

三、SGD (Stochastic gradient descent)

梯度下降法

1. Algorithm description:

solve the following optimization problem:

$$\min_w J(w) = E[f(w, x)]$$

w is parameter to be optimized

x is random variable.

① method 1: gradient descent (GD)

$$w_{k+1} = w_k - \alpha_k \nabla_w E[f(w_k, x)] = w_k - \alpha_k \nabla_w E[\nabla_w f(w_k, x)]$$

draw-back: the expected value is difficult to obtain.

② method 2: batch gradient descent (BGD)

$$E[\nabla_w f(w_k, x)] \approx \frac{1}{n} \sum_{i=1}^n \nabla_w f(w_k, x_i)$$

$$w_{k+1} = w_k - \alpha_k \frac{1}{n} \sum_{i=1}^n \nabla_w f(w_k, x_i)$$

draw-back: requires many samples in each iteration.

③ SGD

$$w_{k+1} = w_k - \alpha_k \underbrace{\nabla_w f(w_k, x_k)}_{\text{stochastic gradient}}$$

BGD let n=1 = SGD

2. Example:

$$\min_w J(w) = E[f(w, X)] = E\left[\frac{1}{2}\|w - X\|^2\right]$$

$$f(w, X) = \|w - X\|^2/2 \quad \nabla_w f(w, X) = w - X$$

optimal solution: $\nabla_w J(w) = 0 \Rightarrow E[\nabla_w f(w, X)] = 0$

$$\Rightarrow E[w - X] = 0 \Rightarrow w = E[X]$$

① GD: $w_{k+1} = w_k - \alpha_k \nabla_w J(w_k) = w_k - \alpha_k E[\nabla_w f(w_k, X)] = w_k - \alpha_k E[w_k - X]$

② SGD: $w_{k+1} = w_k - \alpha_k \nabla_w f(w_k, X_k) = w_k - \alpha_k (w_k - x_k)$

stochastic gradient = true gradient

3. SGD properties

$$\delta_k \leq \frac{\nabla_w f(w_k, x_k) - E[\nabla_w f(w_k, X)]}{C |w_k - w^*|}$$

$w_k \xrightarrow{\delta_k} w^*$ 从远到近 SGD & GD

