

Last time:

Stochastic Approximation Iteration:

$$x_{n+1} = x_n + a(n) [h(x_n) + M_{n+1}]$$

We came to this via : Uen model ($h(x) = p(x) - x$)

We saw the need for two conditions on step sizes:

$$\textcircled{1} \quad \sum_n a(n) = \infty$$

$$\textcircled{2} \quad \sum_n a^2(n) < \infty$$

Today: Example 2: Regression.

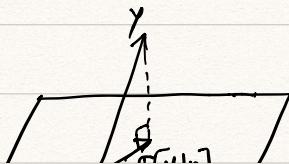
Preliminaries: $(X, Y) \sim P_{XY}$, $X \in \mathbb{R}^n$, $Y \in \mathbb{R}^k$.

Problem: Observe X , estimate Y .

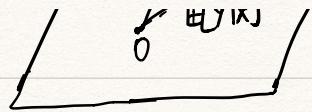
$$\underset{f(\cdot)}{\rightarrow} \min \frac{1}{2} E[\|Y - f(X)\|^2]$$

The answer to this is: $f^*(x) = E[Y|X]$.

Proof (Exercise): Add & sub $E[Y|X]$ & use the orthogonality property:



$$E[(Y - E[Y|X])^T f(X)] = 0$$



Modified Problem: Given a family of parameterized functions $\{f_w : \mathbb{R}^m \rightarrow \mathbb{R}^k, w \in \mathbb{R}^d\}$, find the best fit w .

$$w^* = \underset{w}{\operatorname{arg\min}} \quad g(w) \quad , \quad g(w) = \frac{1}{2} \mathbb{E}[\|Y - f_w(x)\|^2]$$

e.g.: linear fns, frequencies of sines/cosines

$$g(w) = \frac{1}{2} \mathbb{E} \left[Y^T Y - 2 f_w(x)^T Y + f_w(x)^T f_w(x) \right]$$

$$\nabla g(w) = -\mathbb{E} \left[\begin{array}{c} \vdots \\ \frac{\partial}{\partial w_i} \sum_{j=1}^k f_{w,j}(x) Y_j \\ \vdots \end{array} \right] + \mathbb{E} \left[\begin{array}{c} \vdots \\ \frac{1}{2} \sum_{j=1}^k f_{w,j}^2(x) \\ \vdots \end{array} \right]$$

(Assuming we can push $\frac{\partial}{\partial w_i}$ inside $\mathbb{E}[\cdot]$)

$$= -\mathbb{E} \left[\begin{array}{c} \vdots \\ \sum_{j=1}^k \frac{\partial}{\partial w_i} f_{w,j}(x) (Y_j - f_{w,j}(x)) \\ \vdots \end{array} \right]$$

$$= -\mathbb{E} \left[\left[\begin{array}{c} \frac{\partial f_{w,1}(x)}{\partial w_1} \frac{\partial f_{w,2}(x)}{\partial w_1} \dots \frac{\partial f_{w,k}(x)}{\partial w_1} \\ \vdots \quad \ddots \quad \vdots \end{array} \right] \left[\begin{array}{c} Y - f_w(x) \end{array} \right] \right] \quad (*)$$

$$\left[\begin{array}{ccc} \frac{\partial}{\partial w_1} f_{w,l}(x) & \dots & \frac{\partial}{\partial w_d} f_{w,l}(x) \end{array} \right] \quad \left] \right]$$

As an aside, Jacobian: Given $f_w(x) = \begin{bmatrix} f_{w,1}(x) \\ \vdots \\ f_{w,k}(x) \end{bmatrix}$,

$$\text{Jacobian} = \left[\begin{array}{ccc} \frac{\partial}{\partial w_1} f_{w,1}(x) & \dots & \frac{\partial}{\partial w_d} f_{w,1}(x) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial w_1} f_{w,k}(x) & \dots & \frac{\partial}{\partial w_d} f_{w,k}(x) \end{array} \right] =: D_w f_w(x)$$

$$\text{So } (*) = (D_w f_w(x))^T$$

$$\text{We thus have } \nabla g(w) = -E[(D_w f_w(x))^T (y - f_w(x))]$$

So the idea is that: $\nabla g(w)$ gives direction of steepest ascent. So take a step in opp dir.

$$\text{Try } w_{n+1} = w_n - \nabla g(w) \Big|_{w=w_n}$$

$$\Rightarrow w_{n+1} = w_n + E[(D_w f_w(x))^T (y - f_w(x))]$$

Now, we don't know P_{xy} & hence it's difficult for us to evaluate the expectation.

Suppose we have access to samples from the joint distⁿ
 $(x_n, y_n)_{n \geq 1} \stackrel{iid}{\sim} P_{XY}$

Try: "Plug-in estimate" for the gradient:

$$w_{n+1} = w_n + D_w f_w(x_n)^T (y_n - f_w(x_n)) \Big|_{w=w_n}$$

Problem: like a random walk, can wander all over.

Idea: smoothen the dynamics.

Smoothed dynamics:

$$\begin{aligned} w_{n+1} &= (1 - a(n)) w_n + a(n) \left[w_n + D_w f_w(x_n)^T (y_n - f_w(x_n)) \right] \Big|_{w=w_n} \\ &= w_n + a(n) \left[D_w f_w(x_n)^T (y_n - f_w(x_n)) \right] \Big|_{w=w_n} \end{aligned}$$

As time progresses, $a(n) \rightarrow 0$

Add & sub. mean.

$$\Rightarrow w_{n+1} = w_n + a(n) \left[-\nabla g(w_n) + D_w f_w(x_n)^T (y_n - f_w(x_n)) \right] \Big|_{w=w_n} - (-\nabla g(w_n)) \underbrace{\}_{M_{n+1}, \text{ noise}}$$

$$\Rightarrow w_{n+1} = w_n + a(n) \left[-\nabla g(w_n) + M_{n+1} \right] \quad (\text{Stoch. grad. descent})$$

This is a noisy discretization of

$$\dot{w}(t) = -\nabla g(w(t)) \rightarrow \text{Gradient dynamics}$$

We anticipate that $w_n \xrightarrow{n \rightarrow \infty} H = \{w : \nabla g(w) = 0\}$.

If "noise is rich enough", we can show convergence to minima, possibly local.

$$\sum_{n=0}^{\infty} a(n) = \infty, \sum_{n=0}^{\infty} a^2(n) < \infty$$

Chapter 2: Preliminaries: ODEs.

$$\dot{x}(t) = h(x(t)), \quad x(0) = \bar{x} \in \mathbb{R}^d.$$

Defn: The ode is "well-posed" if for any $\bar{x} \in \mathbb{R}^d$, the solution is unique, $x(\cdot) \in C([0, \infty); \mathbb{R}^d)$,

$\xrightarrow[\substack{\text{cont.} \\ \forall t \in [0, \infty)}}]{} x(\cdot) \in \mathbb{R}^d$

space of continuous fns

and further the mapping $\bar{x} \xrightarrow{\mathbb{R}^d} x(\cdot) \in C([0, \infty); \mathbb{R}^d)$

is continuous.

Remarks:

- $C([0, \infty); \mathbb{R}^d)$: space of continuous fns for all positive time,

taking values in \mathbb{R}^d .

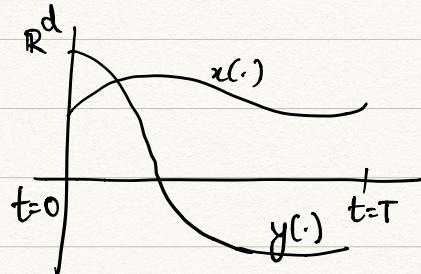
$x(\cdot) \in C([0, \infty); \mathbb{R}^d)$ is a function $x: [0, \infty) \rightarrow \mathbb{R}^d$.

- $C([0, \infty); \mathbb{R}^d)$ is a metric space $\rho(x(\cdot), y(\cdot))$

- This induces a topology \rightarrow collection of open sets generated by open balls.

- A mapping is continuous, if inverse images of open sets are open.

- To understand ρ , let us consider first $C([0, T]; \mathbb{R}^d)$, $T > 0$, finite.



$$\|x(\cdot) - y(\cdot)\|_T = \sup_{t \in [0, T]} \|x(t) - y(t)\|$$

= width of the smallest tube around $x(\cdot)$ that
swallows $y(\cdot)$

(Exercise: Verify that this is a metric).

- (Back to ρ):

$$\rho(x(\cdot), y(\cdot)) = \sum_{T=1}^{\infty} \left(\left(\|x(\cdot)|_T - y(\cdot)|_T\right) \wedge 1 \right) 2^{-T}$$