

## 在线学习

## 场景

在线学习框架，每一次迭代，环境给出一个样本，模型作出预测并产生loss，更新样本

广告场景：用户请求->模型预估->产生loss(用户是否点击)->更新模型

业务实践：用户请求->模型预估->等待一天(或小时级)->产生loss(用户是否点击)->更新模型

For  $t = 1, \dots, T$

- Player chooses  $w_t \in \mathcal{W}$ , where  $\mathcal{W}$  is a *convex* set in  $\mathbb{R}^n$ .
- Environment chooses a *convex* loss function  $f_t : \mathcal{W} \rightarrow \mathbb{R}$ .
- Player incurs a loss  $\ell_t = f_t(w_t) = f_t(w_t; (x_t, y_t))$ .
- Player receives feedback  $f_t$ .

## 在线学习算法评估

$$R(T) = \sum_{t=1}^T f_t(w_t) - \min_{w \in \mathcal{W}} \sum_{t=1}^T f_t(w).$$

解释：评估后悔程度， $w$ 为全局最优解， $w_t$ 为每次更新的权

重,  $R(T)$ 表示了在线学习和批量学习的gap。当 $R(T)$ 随 $T$ 增长低于线性时,效果随迭代而增加。

分析方法：对于在线学习算法，证明其 $R(T) \leq a$ ，得到效果下界。

举例：

最简单的在线学习算法，在线梯度下降， $R(T)$ 分析略。

白话主人明不叫 •

### Algorithm 2. Stochastic Gradient Descent

```

Loop {
  for j=1 to M {
     $W^{(t+1)} = W^{(t)} - \eta^{(t)} \nabla_W \ell(W^{(t)}, Z_j)$ 
  }
}

```

## 关于稀疏性

一般是通过约束区域和损失函数等高线，证明I1比I2更容易获得稀疏解。

另外一种解释：I2在靠近0的附近变化变小，而I1变化量不变。

Consider the vector  $\vec{x} = (1, \varepsilon) \in \mathbb{R}^2$  where  $\varepsilon > 0$  is small. The  $l_1$  and  $l_2$  norms of  $\vec{x}$ , respectively, are given by

$$\|\vec{x}\|_1 = 1 + \varepsilon, \quad \|\vec{x}\|_2^2 = 1 + \varepsilon^2$$

Now say that, as part of some regularization procedure, we are going to reduce the magnitude of one of the elements of  $\vec{x}$  by  $\delta \leq \varepsilon$ . If we change  $x_1$  to  $1 - \delta$ , the resulting norms are

$$\|\vec{x} - (\delta, 0)\|_1 = 1 - \delta + \varepsilon, \quad \|\vec{x} - (\delta, 0)\|_2^2 = 1 - 2\delta + \delta^2 + \varepsilon^2$$

On the other hand, reducing  $x_2$  by  $\delta$  gives norms

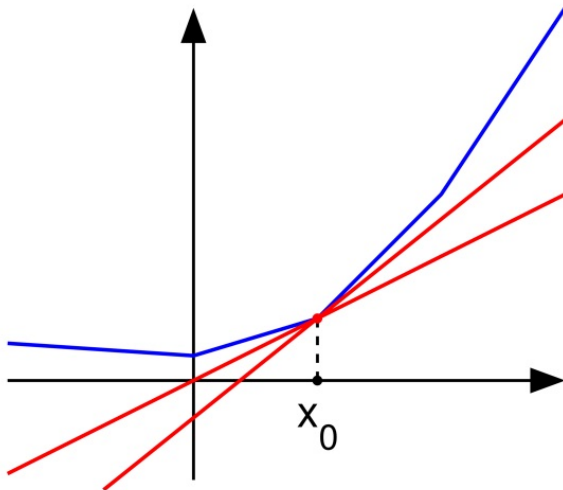
$$\|\vec{x} - (0, \delta)\|_1 = 1 - \delta + \varepsilon, \quad \|\vec{x} - (0, \delta)\|_2^2 = 1 - 2\varepsilon\delta + \delta^2 + \varepsilon^2$$

The thing to notice here is that, for an  $l_2$  penalty, regularizing the larger term  $x_1$  results in a much greater reduction in norm than doing so to the smaller term  $x_2 \approx 0$ . For the  $l_1$  penalty, however, the reduction is the same. Thus, when penalizing a model using the  $l_2$  norm, it is highly unlikely that anything will ever be set to zero, since the reduction in  $l_2$  norm going from  $\varepsilon$  to 0 is almost nonexistent when  $\varepsilon$  is small. On the other hand, the reduction in  $l_1$  norm is always equal to  $\delta$ , regardless of the quantity being penalized.

Another way to think of it: it's not so much that  $l_1$  penalties encourage sparsity, but that  $l_2$  penalties in some sense **discourage** sparsity by yielding diminishing returns as elements are moved closer to zero.

优化问题没有解析解时，即使使用l1也很难获得稀疏解。l1只能保证快速逼近到0附近。在线学习只使用l1不容易获得稀疏解。需要额外手段来保证稀疏。

## 在线梯度下降l1正则化



## 次梯度

以 $k$ 为窗口，当 $t/k$ 不为整数时采用标准的SGD进行迭代，当 $t/k$ 为整数时，采用如下权重更新方式：

$$W^{(t+1)} = T_0(W^{(t)} - \eta^{(t)} G^{(t)}, \theta) \quad (3-1-2)$$

$$T_0(v_i, \theta) = \begin{cases} 0 & \text{if } |v_i| \leq \theta \\ v_i & \text{otherwise} \end{cases}$$

注意，这里面 $\theta \in \mathbb{R}$ 是一个标量，且 $\theta \geq 0$ ；如果 $V = [v_1, v_2, \dots, v_N] \in \mathbb{R}^N$ 是一个向量， $v_i$ 是向量的一个维度，那么有 $T_0(V, \theta) = [T_0(v_1, \theta), T_0(v_2, \theta), \dots, T_0(v_N, \theta)] \in \mathbb{R}^N$ 。

## 简单截断法

$$W^{(t+1)} = T_1(W^{(t)} - \eta^{(t)} G^{(t)}, \eta^{(t)} \lambda^{(t)}, \theta) \quad (3-1-3)$$

$$T_1(v_i, \alpha, \theta) = \begin{cases} \max\{0, v_i - \alpha\} & \text{if } v_i \in [0, \theta] \\ \min\{0, v_i + \alpha\} & \text{if } v_i \in [-\theta, 0] \\ v_i & \text{otherwise} \end{cases}$$

直观的获得稀疏性的方法，当权重小于阈值时设置为0

## 截断梯度法

### Algorithm 3. Truncated Gradient

---

```

1  input  $\theta$ 
2  initial  $W \in \mathbb{R}^N$ 
3  for  $t = 1, 2, 3 \dots$  do
4     $G = \nabla_W \ell(W, X^{(t)}, y^{(t)})$ 
5    refresh  $W$  according to
      
$$w_i = \begin{cases} \max\{0, w_i - \eta^{(t)} g_i - \eta^{(t)} \lambda^{(t)}\} & \text{if } (w_i - \eta^{(t)} g_i) \in [0, \theta] \\ \max\{0, w_i - \eta^{(t)} g_i + \eta^{(t)} \lambda^{(t)}\} & \text{if } (w_i - \eta^{(t)} g_i) \in [-\theta, 0] \\ w_i - \eta^{(t)} g_i & \text{otherwise} \end{cases}$$

6  end
7  return  $W$ 

```

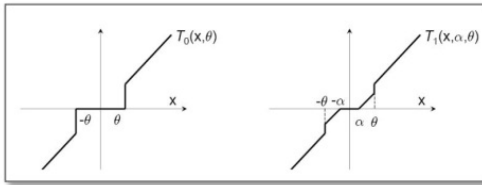
---

$$w_i^{(t+1)} = \begin{cases} \text{Trnc}\left((w_i^{(t)} - \eta^{(t)} g_i^{(t)}), \lambda_{TG}^{(t)}, \theta\right) & \text{if } \text{mod}(t, k) = 0 \\ w_i^{(t)} - \eta^{(t)} g_i^{(t)} & \text{otherwise} \end{cases}$$

$$\lambda_{TG}^{(t)} = \eta^{(t)} \lambda k \quad (3-1-4)$$

$$\text{Trnc}(w, \lambda_{TG}^{(t)}, \theta) = \begin{cases} 0 & \text{if } |w| \leq \lambda_{TG}^{(t)} \\ w - \lambda_{TG}^{(t)} \text{sgn}(w) & \text{if } \lambda_{TG}^{(t)} \leq |w| \leq \theta \\ w & \text{otherwise} \end{cases}$$

l1正则化和简单截断法都是梯度截断的特殊形式。



$$W^{(t+\frac{1}{2})} = W^{(t)} - \eta^{(t)} G^{(t)}$$

$$W^{(t+1)} = \underset{W}{\operatorname{argmin}} \left\{ \frac{1}{2} \|W - W^{(t+\frac{1}{2})}\|^2 + \eta^{(t+\frac{1}{2})} \Psi(W) \right\}$$

## FOBOS

权重更新分两个步骤：

$$W^{(t+1)} = \underset{W}{\operatorname{argmin}} \left\{ \frac{1}{2} \|W - W^{(t)} + \eta^{(t)} G^{(t)}\|^2 + \eta^{(t+\frac{1}{2})} \Psi(W) \right\}$$

第一步为梯度下降，第二步在第一步梯度下降得到权重附近加入正则进行最优优化。  
合并公式有：

$$W^{(t+1)} = \underset{W}{\operatorname{argmin}} \sum_{i=1}^N \left( \frac{1}{2} (w_i - v_i)^2 + \tilde{\lambda} |w_i| \right)$$

考虑正则为l1的FOBOS：

$$\begin{aligned} w_i^{(t+1)} &= \text{sgn}(v_i) \max(0, |v_i| - \tilde{\lambda}) \\ &= \text{sgn}(w_i^{(t)} - \eta^{(t)} g_i^{(t)}) \max\{0, |w_i^{(t)} - \eta^{(t)} g_i^{(t)}| - \eta^{(t+\frac{1}{2})} \lambda\} \end{aligned}$$

使用技巧求解，得到

$$w_i^{(t+1)} = \begin{cases} 0 & \text{if } |w_i^{(t)} - \eta^{(t)} g_i^{(t)}| \leq \eta^{(t+\frac{1}{2})} \lambda \\ (w_i^{(t)} - \eta^{(t)} g_i^{(t)}) - \eta^{(t+\frac{1}{2})} \lambda \text{sgn}(w_i^{(t)} - \eta^{(t)} g_i^{(t)}) & \text{otherwise} \end{cases}$$

写成梯度截断的形式：

Follow-The-Leader
$w_1$ is set arbitrarily for $t = 1, 2, \dots, T$ $w_t = \operatorname{argmin}_{w \in \mathcal{W}} \sum_{s=1}^{t-1} f_s(w)$

显式的获取稀疏性。

## FTL

$$w_{t+1} = \operatorname{argmin}_{w \in W} (f_{1:t}(w) + R(w)).$$

在线学习中每一轮将w更新为在全部历史轮数上的最有解。

## FTRL

$$x^{(t+1)} = \operatorname{argmin}_x \{ (x - (x^{(t)} - \eta \nabla l(x^{(t)})))^2 + \lambda |x| \}$$

展开并忽略其中一个常数项  $(\eta \nabla l(x^{(t)}))^2$  后得到等价的形式

$$x^{(t+1)} = \operatorname{argmin}_x \{ 2\eta \nabla l(x^{(t)})(x - x^{(t)}) + (x - x^{(t)})^2 + \lambda |x| \}$$

## FTRL-Proximal

考虑l1-FOBOS的更新公式：

$$W^{(t+1)} = \operatorname{argmin}_W \left\{ G^{(1:t)} \cdot W + \lambda_1 \|W\|_1 + \lambda_2 \frac{1}{2} \|W\|_2^2 + \frac{1}{2} \sum_{s=1}^t \sigma^{(s)} \|W - W^{(s)}\|_2^2 \right\}$$

定义FTRL损失函数  $f_{-}(t+1)(x) = g_{-}(x-x_t) + (1/n_t-1/n_{-}(t-1))(x-x_t)^2$

套用FTRL优化方法： $\operatorname{argmin} \sum f_{-}(t)+R$  得到FTRL的优化函数

$$W^{(t+1)} = \operatorname{argmin}_W \left\{ \left( G^{(1:t)} - \sum_{s=1}^t \sigma^{(s)} W^{(s)} \right) \cdot W + \lambda_1 \|W\|_1 + \frac{1}{2} \left( \lambda_2 + \sum_{s=1}^t \sigma^{(s)} \right) \|W\|_2^2 + \frac{1}{2} \sum_{s=1}^t \sigma^{(s)} \|W^{(s)}\|_2^2 \right\}$$

由于  $\frac{1}{2} \sum_{s=1}^t \sigma^{(s)} \|W^{(s)}\|_2^2$  相对于  $W$  来说是一个常数，并且令  $Z^{(t)} = G^{(1:t)} - \sum_{s=1}^t \sigma^{(s)} W^{(s)}$ ，上

式等价于：

$$W^{(t+1)} = \operatorname{argmin}_W \left\{ Z^{(t)} \cdot W + \lambda_1 \|W\|_1 + \frac{1}{2} \left( \lambda_2 + \sum_{s=1}^t \sigma^{(s)} \right) \|W\|_2^2 \right\}$$

将最后一项展开，可以得到

$$w_i^{(t+1)} = \begin{cases} 0 & \text{if } |z_i^{(t)}| < \lambda_1 \\ - \left( \lambda_2 + \sum_{s=1}^t \sigma^{(s)} \right)^{-1} (z_i^{(t)} - \lambda_1 \operatorname{sgn}(z_i^{(t)})) & \text{otherwise} \end{cases}$$