

8.3 二阶方法

8.4 随机梯度下降

纪一仑

Harbin Institute of Technology, ShenZhen

Group Meeting, January 2023

1 8.3 二阶方法

- 8.3.1 牛顿法
- 8.3.2 BFGS 和其他类牛顿法
- 8.3.3 Trust region method

2 8.4 随机梯度下降

- 8.4.1 有限和问题的应用
- 8.4.2 实例: SGD 用于线性回归

1 8.3 二阶方法

- 8.3.1 牛顿法
- 8.3.2 BFGS 和其他类牛顿法
- 8.3.3 Trust region method

2 8.4 随机梯度下降

- 8.4.1 有限和问题的应用
- 8.4.2 实例: SGD 用于线性回归

8.3.1 牛顿法

牛顿法

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \mathbf{H}_t^{-1} \mathbf{g}_t$$

where

$$\mathbf{H}_t \triangleq \nabla^2 \mathcal{L}(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_t} = \nabla^2 \mathcal{L}(\boldsymbol{\theta}_t) = \mathbf{H}(\boldsymbol{\theta}_t)$$

8.3.1 牛顿法

推导

考虑用二阶泰勒级数近似在 $\boldsymbol{\theta}_t$ 周围的 $\mathcal{L}(\boldsymbol{\theta})$

$$\mathcal{L}_{quad}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_t)^\top \mathbf{H}_t (\boldsymbol{\theta} - \boldsymbol{\theta}_t)$$

最小值点是 $\boldsymbol{\theta} = \boldsymbol{\theta}_t - \mathbf{H}_t^{-1} \mathbf{g}_t$, 需要 \mathbf{H}_t 正定

线性回归应用

$$\begin{aligned} \text{RSS}(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \left(y_n - \mathbf{w}^\top \mathbf{x}_n \right)^2 = \frac{1}{2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^\top (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= \frac{1}{2} \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\mathbf{w} + \frac{1}{2} \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w} \end{aligned}$$

损失函数就是二阶泰勒级数, 使用牛顿法可以直接得到结果

8.3.2 BFGS 和其他类牛顿法

计算 Hessian 矩阵 \mathbf{H}_t 的代价比较大，为了减小这个代价，类牛顿法使用 Pseudo Hessian 矩阵 \mathbf{B}_t 来近似 \mathbf{H}_t 。

Pseudo Hessian 的要求

- 1. \mathbf{B}_t 正定对称
- 2. \mathbf{B}_{t+1} 由之前迭代步骤的梯度/更新方向得到
- 3. \mathbf{B}_{t+1} 需要逼近 \mathbf{B}_t ，来保证收敛
- 4. \mathbf{B}_{t+1} 的更新开销需要尽可能小
- 5. $\mathbf{B}_{t+1}(\boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}) = \mathbf{g}_t - \mathbf{g}_{t-1}$

8.3.2 BFGS 和其他类牛顿法

BFGS 公式的推导

求如下的优化问题

$$\min_{\mathbf{B}_+} \|\mathbf{B}_+ - \mathbf{B}\|_F$$

$$\text{s. t. } \mathbf{B}_+ \boldsymbol{\delta} = \boldsymbol{\gamma} \text{ and } \mathbf{B}_+ = \mathbf{B}_+^\top \text{ and } \mathbf{B}_+ \succ 0$$

参考链接: <https://zhuanlan.zhihu.com/p/573703008>

BFGS 公式

$$\mathbf{B}_{t+1} = \mathbf{B}_t + \frac{\mathbf{y}_t \mathbf{y}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} - \frac{(\mathbf{B}_t \mathbf{s}_t) (\mathbf{B}_t \mathbf{s}_t)^\top}{\mathbf{s}_t^\top \mathbf{B}_t \mathbf{s}_t}$$

$$\mathbf{s}_t = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1}$$

$$\mathbf{y}_t = \mathbf{g}_t - \mathbf{g}_{t-1}$$

8.3.2 BFGS 和其他类牛顿法

\mathbf{B}_t 保持正定需要满足下面的条件

Wolfe Conditions

- $\mathcal{L}(\boldsymbol{\theta}_t + \eta \mathbf{d}_t) \leq \mathcal{L}(\boldsymbol{\theta}_t) + c\eta \mathbf{d}_t^\top \nabla \mathcal{L}(\boldsymbol{\theta}_t)$
- $-\mathbf{d}_t^\top \nabla \mathcal{L}(\boldsymbol{\theta}_t + \eta \mathbf{d}_t) \leq -c_2 \mathbf{d}_t^\top \nabla \mathcal{L}(\boldsymbol{\theta}_t)$

参考链接: <https://www.zhihu.com/question/49600881>

其中 $0 < c < c_2 < 1$

第一个条件确保函数单调下降。

第二个条件确保函数的梯度也下降。

8.3.2 BFGS 和其他类牛顿法

其他需要注意的事情

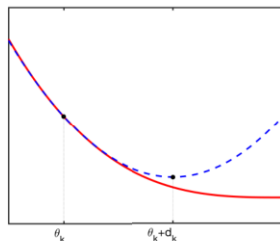
- 可以使用 SWM 公式直接维护 \mathbf{B}_{t+1} 的逆 \mathbf{C}_{t+1}

$$\mathbf{C}_{t+1} = \left(\mathbf{I} - \frac{\mathbf{s}_t \mathbf{y}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} \right) \mathbf{C}_t \left(\mathbf{I} - \frac{\mathbf{y}_t \mathbf{s}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} \right) + \frac{\mathbf{s}_t \mathbf{s}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t}$$

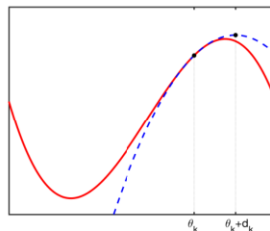
- 存储 \mathbf{B}_t 需要 $O(D^2)$ 的空间, 对于很大的问题, 可以只存储最近的 M 对 $(\mathbf{s}_t, \mathbf{y}_t)$, 然后用这些 $(\mathbf{s}_t, \mathbf{y}_t)$ 直接计算 \mathbf{B}_t , 这样只需要 $O(MD)$ 的空间。这就是 limited memory BFMS, 即 L-BFMS。

8.3.3 Trust region method

牛顿法得到的是多元二次函数的驻点，所以当 \mathbf{H}_t 不正定时， $\mathbf{d}_t = -\mathbf{H}_t^{-1} \mathbf{g}_t$ 可能不是函数下降的方向。



(a)



(b)

Figure: 牛顿法不一定向着函数下降的方向

8.3.3 Trust region method

但是，这个多元二次函数在一定区域内肯定有最小值。

Trust-region optimization

$$\delta^* = \arg \min_{\delta \in \mathcal{R}_t} M_t(\delta)$$

其中 $\delta = \theta - \theta_t$ ， \mathcal{R}_t 是所取区域， $M(\delta)$ 是目标函数或者它的估计。

8.3.3 Trust region method

Tikhonov damping

通常情况下，假设

$$M_t(\boldsymbol{\delta}) = \mathcal{L}_{quad}(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^\top \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^\top \mathbf{H}_t \boldsymbol{\delta}$$

$$\mathcal{R}_t = \{\boldsymbol{\delta} : \|\boldsymbol{\delta}\|_2 \leq r\}$$

使用拉格朗日乘数法将有约束问题转化为无约束问题

$$\boldsymbol{\delta}^* = \arg \min_{\boldsymbol{\delta}} M(\boldsymbol{\delta}) + \lambda \|\boldsymbol{\delta}\|_2^2 = \arg \min_{\boldsymbol{\delta}} \mathbf{g}^\top \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^\top (\mathbf{H} + \lambda \mathbf{I}) \boldsymbol{\delta}$$

当 λ 足够大，这个函数的 Hessian，即 $\mathbf{H} + \lambda \mathbf{I}$ 肯定是正定的，于是有最小值点

$$\boldsymbol{\delta} = -(\mathbf{H} + \lambda \mathbf{I})^{-1} \mathbf{g}$$

8.3.3 Trust region method

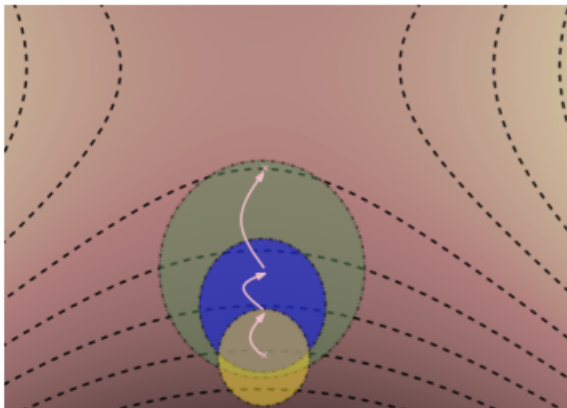


Figure: trust region method

1 8.3 二阶方法

- 8.3.1 牛顿法
- 8.3.2 BFGS 和其他类牛顿法
- 8.3.3 Trust region method

2 8.4 随机梯度下降

- 8.4.1 有限和问题的应用
- 8.4.2 实例: SGD 用于线性回归

8.4 随机梯度下降

目标函数:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z})}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{z})]$$

这里 \mathbf{z} 是随机输入, 且 $\mathbf{z} \sim q$ 。

在每次迭代中, 假设我们观测到 $\mathcal{L}_t(\boldsymbol{\theta}) = \mathcal{L}(\boldsymbol{\theta}, \mathbf{z}_t)$, 如果分布 $q(\mathbf{z})$ 独立于我们优化的参数 $\boldsymbol{\theta}$, 我们可以使用 $\mathbf{g}_t = \nabla_{\boldsymbol{\theta}} \mathcal{L}_t(\boldsymbol{\theta}_t)$ 来近似 $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}_t)$

于是, 迭代公式为:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t \nabla \mathcal{L}(\boldsymbol{\theta}_t, \mathbf{z}_t) = \boldsymbol{\theta}_t - \eta_t \mathbf{g}_t$$

8.4.1 有限和问题的应用

经验风险最小化的损失函数:

$$\mathcal{L}(\boldsymbol{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{y}_n, f(\mathbf{x}_n; \boldsymbol{\theta}_t)) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\theta}_t)$$

此损失函数的梯度:

$$\mathbf{g}_t = \frac{1}{N} \sum_{n=1}^N \nabla_{\boldsymbol{\theta}} \mathcal{L}_n(\boldsymbol{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \nabla_{\boldsymbol{\theta}} \ell(\mathbf{y}_n, f(\mathbf{x}_n; \boldsymbol{\theta}_t))$$

可以看到, 想要精确计算出梯度, 需要在所有 N 个样本点上计算。

8.4.1 有限和问题的应用

当 N 很大的时候，计算梯度会很慢。幸运的是，我们可以用少量样本上的梯度估计所有样本上的梯度。

在所有 N 个样本中随机取 $B \ll N$ 个样本，作为一个 minibatch，计算 minibatch 上的梯度：

$$\mathbf{g}_t \approx \frac{1}{|\mathcal{B}_t|} \sum_{n \in \mathcal{B}_t} \nabla_{\boldsymbol{\theta}} \mathcal{L}_n(\boldsymbol{\theta}_t) = \frac{1}{|\mathcal{B}_t|} \sum_{n \in \mathcal{B}_t} \nabla_{\boldsymbol{\theta}} \ell(\mathbf{y}_n, f(\mathbf{x}_n; \boldsymbol{\theta}_t))$$

因为 minibatch 是随机取的，所以这个是无偏估计。

8.4.2 实例：SGD 用于线性回归

线性回归的损失函数：

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{2N} \sum_{n=1}^N \left(\mathbf{x}_n^\top \boldsymbol{\theta} - y_n \right)^2 = \frac{1}{2N} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2$$

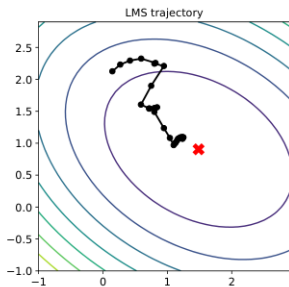
梯度：

$$\mathbf{g}_t = \frac{1}{N} \sum_{n=1}^N (\boldsymbol{\theta}_t^\top \mathbf{x}_n - y_n) \mathbf{x}_n$$

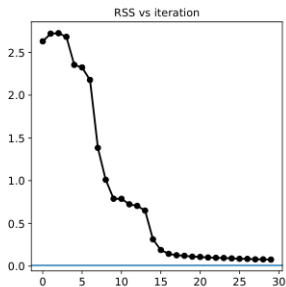
使用 SGD，minibatch 的大小取 $B = 1$ ，迭代公式为：

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta_t (\boldsymbol{\theta}_t^\top \mathbf{x}_n - y_n) \mathbf{x}_n$$

8.4.2 实例：SGD 用于线性回归



(a)



(b)

Figure: SGD for fitting linear regression