

Week 2: Tutorial Hand Note

Unit Vector \hat{v} /单位向量

定义:

$$\vec{v} \text{ is } \hat{v} \text{ iff. } \|\vec{v}\| = 1$$

即长度为 1 的向量为单位向量。

Norm 2

$\|\vec{v}\|$ 即向量 v 的 2-范数 (Norm 2, 即欧几里得范数, 或者说是几何距离), 其定义为

$$\|\vec{v}\| = \sqrt{\sum_{i=1}^d v_i^2}$$

巧思: 这个也同样是 L2 Loss 的表达形式

Normalisation

对于任意向量 (除 0 以外), 其总能获得与其方向 (Direction) 相同的单位向量。其公式为:

$$\|\hat{v}\| = \frac{\vec{v}}{\|\vec{v}\|}$$

我们称上述公式为 Normalisation。

Taylor Polynomials/泰勒级数

单变量公式: $\mathbb{R} \mapsto \mathbb{R}$

对于函数 $E(x)$ 于 w_0 处展开, 获得 w 处的近似值, 其公式被定义为:

$$T(w) = \sum_{k=0}^n \left[\frac{E^{(k)}(w_0)}{k!} (w - w_0)^k \right]$$

多变量公式: $\mathbb{R}^d \mapsto \mathbb{R}$

$$T(\vec{w}) = \sum_{k=0}^n \left[\frac{E_{\vec{w}}^{(k)}(\vec{w}_0)}{k!} (\vec{w} - \vec{w}_0)^k \right]$$

其中 $E_{\vec{w}}^{(k)}$, (k) 表示 k 阶导数, 而 \vec{w} 表示是对函数 $E(w)$ 的全导数 (total derivative)。

$$E_{\vec{w}}^{(0)}(\vec{w}_0) = E(\vec{w}_0)$$

$$E_{\vec{w}}^{(1)}(\vec{w}_0) = \nabla E(\vec{w}_0)$$

$$E_{\vec{w}}^{(2)}(\vec{w}_0) = \mathbf{H}(\vec{w}_0)$$

$$(\vec{w} - \vec{w}_0)^0 = 1$$

$$1$$

$$(\vec{w} - \vec{w}_0)^1 = (\vec{w} - \vec{w}_0)$$

$$d \times 1$$

$$(\vec{w} - \vec{w}_0)^2 = (\vec{w} - \vec{w}_0) \cdot (\vec{w} - \vec{w}_0)^T$$

$$(d \times 1) \times (1 \times d) = d \times d$$

$$E_{\vec{w}}^{(k)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^k = ?$$

$$E_{\vec{w}}^{(0)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^0 = E(\vec{w}_0)$$

$$E_{\vec{w}}^{(1)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^1 = \nabla E(\vec{w}_0)^T (\vec{w} - \vec{w}_0)$$

$$E_{\vec{w}}^{(2)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^2 = (\vec{w} - \vec{w}_0)^T \mathbf{H}(E(\vec{w}_0)) (\vec{w} - \vec{w}_0)$$

其中:

$$\begin{aligned} \nabla E(\mathbf{w}) &= \begin{pmatrix} \frac{\partial E(\mathbf{w})}{\partial w_1} \\ \frac{\partial E(\mathbf{w})}{\partial w_2} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w_{d-1}} \\ \frac{\partial E(\mathbf{w})}{\partial w_d} \end{pmatrix} = \mathbf{H}_E \\ &= \begin{bmatrix} \frac{\partial^2 E(\mathbf{w})}{\partial w_1^2} & \frac{\partial^2 E(\mathbf{w})}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 E(\mathbf{w})}{\partial w_1 \partial w_d} \\ \frac{\partial^2 E(\mathbf{w})}{\partial w_2 \partial w_1} & \frac{\partial^2 E(\mathbf{w})}{\partial w_2^2} & \cdots & \frac{\partial^2 E(\mathbf{w})}{\partial w_2 \partial w_d} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\mathbf{w})}{\partial w_d \partial w_1} & \frac{\partial^2 E(\mathbf{w})}{\partial w_d \partial w_2} & \cdots & \frac{\partial^2 E(\mathbf{w})}{\partial w_d^2} \end{bmatrix} \\ &\quad (\mathbf{H}_E)_{i,j} = \frac{\partial^2 E(\mathbf{w})}{\partial w_i \partial w_j} \end{aligned}$$

How to prove that GD really takes a step in the direction of the steepest descent?

如何证明 "梯度下降法" 真的向最陡下降的方向迈出了一步?

可以将 GD 的公式记为

$$\vec{w} = \vec{w}_0 - \eta \nabla E(\vec{w}_0)$$

假设其的公式为

$$\vec{w} = \vec{w}_0 + \eta \hat{v}$$

其中 η 是学习率, 也同样是学习的步长 (每一次学习, 学习多长)

而 \hat{v} 则是学习的方向 (向哪里学习?)

→ 因此我们需要找到一个 \hat{v} 使 $E(\vec{w}) - E(\vec{w}_0)$ 尽可能的负 (**Negative**)

解释: 因为是学习, 所以我们期望学习后的 $E(\vec{w})$ 尽可能小 (当然, 需要小于 $E(\vec{w}_0)$)

所以我们期望 $E(\vec{w}) - E(\vec{w}_0) \leq 0$

而如果我们学习方向足够正确, 那么一定会让 $E(\vec{w})$ 足够小, 也就意味着

$$(E(\vec{w}) \downarrow -E(\vec{w}_0)) \downarrow$$

也就是期望其尽可能的负。

→ 因此我们需要证明 $-\nabla E(\vec{w}_0)$ 是尽可能的负

假设: 学习率 η 足够小

假设原因: 我们需要使用泰勒级数进行展开

我们对 $E(\vec{w})$ 进行 1 阶泰勒展开

$$E(\vec{w}) = \frac{E_{\vec{w}}^{(0)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^0}{0!} + \frac{E_{\vec{w}}^{(1)}(\vec{w}_0)(\vec{w} - \vec{w}_0)^1}{1!} + \dots (\text{ignored})$$

$$E(\vec{w}) = E(\vec{w}_0) + \nabla E(\vec{w}_0)^T (\vec{w} - \vec{w}_0)$$

$$E(\vec{w}) - E(\vec{w}_0) = \nabla E(\vec{w}_0)^T (\vec{w} - \vec{w}_0)$$

代入 $\vec{w} = \vec{w}_0 + \eta \hat{v}$

$$E(\vec{w}_0 - \eta \hat{v}) - E(\vec{w}_0) = \nabla E(\vec{w}_0)^T (\vec{w}_0 + \eta \hat{v} - \vec{w}_0)$$

$$= \eta \nabla E(\vec{w}_0)^T \hat{v}$$

End of the Tutorial (Thursday, Week 2), TBC.

Following Content is Week2, Friday, the Normal Lecture.

整理上式，我们可以得到：

$$E(\vec{w}_0 - \eta \vec{v}) - E(\vec{w}_0) = \eta \nabla E(\vec{w}_0)^T \vec{v}$$

并且我们期望 $E(\vec{w}_0 - \eta \vec{v}) - E(\vec{w}_0)$ 尽可能的负。

→ 因此我们期望 $\eta \nabla E(\vec{w}_0)^T \vec{w}_0$ 尽可能的负

→ 不考虑常数项 η ，因此我们期望 $\nabla E(\vec{w}_0)^T \vec{w}_0$ 尽可能的负

而对于 $\nabla E(\vec{w}_0)^T \vec{w}_0$ ，我们可以将其看作向量 $\nabla E(\vec{w}_0)^T$ 与向量 \vec{w}_0 的点乘 (Dot Product)。因此我们可以将上式写成

$$\begin{aligned} \nabla E(\vec{w}_0)^T \vec{v} &= \|\nabla E(\vec{w}_0)\| \|\vec{v}\| \cos \theta \\ &= \|\nabla E(\vec{w}_0)\| \cos \theta \end{aligned}$$

其中 θ 表示向量 \vec{v} 与梯度 $\nabla E(\vec{w}_0)$ 的夹角。

因此我们可以改变上式得到：

$$\cos \theta = \frac{\nabla E(\vec{w}_0)^T \hat{v}}{\|\nabla E(\vec{w}_0)\|} \in [-1, 1]$$

我们知道 $\cos \theta \in [-1, 1]$ ，因此

$$(\theta = 180^\circ) - 1 \leq \cos \theta = \frac{\nabla E(\vec{w}_0)^T \hat{v}}{\|\nabla E(\vec{w}_0)\|} \leq 1 \quad (\theta = 0^\circ)$$

$$(\theta = 180^\circ) - \|\nabla E(\vec{w}_0)\| \leq \nabla E(\vec{w}_0)^T \hat{v} \leq \|\nabla E(\vec{w}_0)\| \quad (\theta = 0^\circ)$$

因此

$$\begin{aligned} E(\vec{w}_0 - \eta \vec{v}) - E(\vec{w}_0) &= \eta \nabla E(\vec{w}_0)^T \vec{v} \\ &\geq -\eta \|\nabla E(\vec{w}_0)\| \end{aligned}$$

因为我们期望 $\nabla E(\vec{w}_0)^T \vec{w}_0$ 尽可能的负，因此我们使其取最小值 $-\|\nabla E(\vec{w}_0)\|$

而此时，表示方向的单位向量 \hat{v} 与梯度 $\nabla E(\vec{w}_0)$ 的夹角 θ 为 180° 。

因此我们学习时候的 \vec{v} 的方向与 $-\nabla E(\vec{w}_0)$ 一致。

这证明了"梯度下降法"真的向最陡下降的方向迈出了一步。