

GAN: Meeting #3

Due on May , 07 2019 at 11:00pm

Student Zhaokun Zhou Section A

Zhaokun Zhou

Problem 1

一. 给定四个公式的证明

符号说明:

$l-1$ 层有 k 个神经元, l 层有 j 个神经元, 设 $l+1$ 层有 m 个神经元

$W_{jk}^{(l)}$ 是第 $l-1$ 层的第 k 个神经元到第 l 层第 j 个神经元的权值

$b_j^{(l)}$ 是第 l 层的第 j 个神经元的偏置 (在此次推导中设每层偏置量相同)

$a_j^{(l)}$ 是第 l 层的第 j 个神经元的激活函数输出值

$\delta_j^{(l)}$ 定义为第 l 层的第 j 个神经元的误差值

$z_j^{(l)}$ 定义为第 l 层的第 j 个神经元的输入

J 为代价函数

$$a_j^{(l)} = \sigma \left(\sum_k W_{jk}^{(l)} a_k^{(l-1)} + b_j^{(l)} \right) \quad (1)$$

$$z_j^{(l)} = \sum_k W_{jk}^{(l)} a_k^{(l-1)} + b_j^{(l)} \quad (2)$$

将式1转化成矩阵形式:

$$a^{(l)} = \sigma \left(\sum_k W^{(l)} a^{(l-1)} + b^{(l)} \right) \quad (3)$$

其中 $a^{(l)}$ 是 j 行1列向量, $W^{(l)}$ 是 j 行 k 列矩阵, $a^{(l-1)}$ 是 k 行1列向量, $b^{(l)}$ 是 j 行1列向量。
将式2转化成矩阵形式:

$$z^{(l)} = W^{(l)} a^{(l-1)} + b^{(l)} \quad (4)$$

结合函数链式求导法则与上述公式, 得到以下四个公式:
公式一:

$$\delta_j^{(l)} = \frac{\partial J}{\partial z_j^{(l)}} = \frac{\partial J}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} = \frac{\partial J}{\partial a_j^{(l)}} \cdot \sigma' \left(z_j^{(l)} \right) \quad (5)$$

将式5转化成矩阵形式:

$$\delta^{(l)} = \nabla_a J \odot \sigma' \left(z^{(l)} \right) \quad (6)$$

其中 $\delta^{(l)}$ 为 j 行1列, $\nabla_a J$ 为 j 行1列, $\sigma' \left(z^{(l)} \right)$ 为 j 行1列。

公式二:

已知 $\delta^{(l+1)}$ 求 $\delta^{(l)}$:

$$\delta^{(l)} = (W^{(l+1)})^T \delta^{(l+1)} \odot \sigma' \left(z^{(l)} \right) \quad (7)$$

设 $l+1$ 层有 m 个神经元, l 层有 j 个神经元, 推导如下:
第 $l+1$ 层第 m 个神经元的偏差 $\delta_j^{(l+1)}$ 对第 l 层第 j 个神经元的偏差 $\delta_j^{(l)}$ 的影响:

$$\delta_j^{(l)}(m) = \frac{\partial J}{\partial z_m^{(l+1)}} \cdot \frac{\partial z_m^{(l+1)}}{\partial z_j^{(l)}} = \delta_j^{(l+1)} \cdot W_{mj}^{(l+1)} \cdot \sigma' \left(z_j^{(l)} \right) \quad (8)$$

整个第 $l+1$ 层对第 l 层第 j 个神经元的偏差 $\delta_j^{(l)}$ 的影响:

$$\delta_j^{(l)} = \sum_m \frac{\partial J}{\partial z_m^{(l+1)}} \cdot \frac{\partial z_m^{(l+1)}}{\partial z_j^{(l)}} = \delta_j^{(l+1)} \cdot W_{mj}^{(l+1)} \cdot \sigma' \left(z_j^{(l)} \right) \quad (9)$$

把上式归纳成整个 l 层的偏差:

$$\delta^{(l)} = \sum_j \sum_m \frac{\partial J}{\partial z_m^{(l+1)}} \cdot \sum_j \frac{\partial z_m^{(l+1)}}{\partial z_j^{(l)}} = (W^{l+1})^T \delta^{(l+1)} \cdot \sigma' \left(z^{(l)} \right) \quad (10)$$

公式三:

基于以上公式, 对权重求偏导得到:

$$\frac{\partial J(W, b)}{\partial W_{jk}^{(l)}} = \delta_j^{(l)} a_k^{(l-1)} \quad (11)$$

公式四:

对偏置求偏导得到:

$$\frac{\partial J(W, b)}{\partial b_j^{(l)}} = \delta_j^{(l)} \quad (12)$$

将上述11式重写为矩阵形式:

$$\nabla_{W^{(l)}} J(W, b) = \delta^{(l)} a^{(l-1)T} \quad (13)$$

将上述12式重写为矩阵形式:

$$\nabla_{b^{(l)}} J(W, b) = \delta^{(l)} \quad (14)$$

完整的迭代小批量随机梯度下降算法:

1. 设定 $\Delta W^{(l)} := 0, \Delta b^{(l)} := 0$ 为误差积累项

2. for $i = 1:m$

2a. 使用BP算法计算 $\nabla_{W^{(l)}} J(W, b), \nabla_{b^{(l)}} J(W, b)$

2b. 设置 $\nabla W^{(l)} := \Delta W^{(l)} + \nabla_{W^{(l)}} J(W, b)$

2c. 设置 $\Delta b^{(l)} := \Delta b^{(l)} + \nabla_{b^{(l)}} J(W, b)$

3. 更新参数

$$W^{(l)} := W^{(l)} - \alpha \left[\left(\frac{1}{m} \right) \Delta W^{(l)} + \lambda W^{(l)} \right]$$

$$b^{(l)} := b^{(l)} - \alpha \left[\left(\frac{1}{m} \right) \Delta b^{(l)} \right]$$

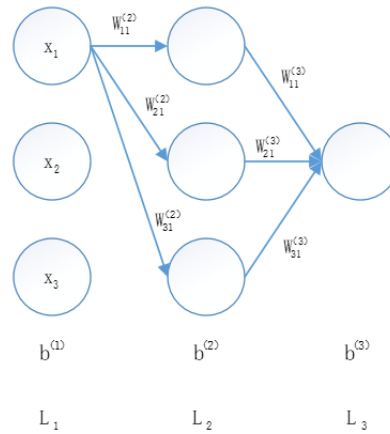


Figure 1: 3-layer MLP

其中 α 为学习率。

代价函数：

$$J(W, b) = \left[\frac{1}{m} \sum_1^m \left(\frac{1}{2} \|h_{W,b}(x^{(i)} - y^{(i)})\|^2 \right) \right] + \frac{\lambda}{2} \sum_{l=1}^2 \sum_{j=1}^{S_l} \sum_{k=1}^{S_{l+1}} W_{jk}^{(l)2} \quad (15)$$

实例计算：

已知数值的变量有： $a_j^{(l)}, x_i, y_i, W_{jk}^{(l)}$
 设激活函数为logistic sigmoid函数，用 $\sigma(x)$ 表示

$$\sigma'(z_j^l) = \frac{\exp^x}{(1 + \exp^x)^2} = \frac{1}{(1 + \exp^x)} \left(1 - \frac{1}{(1 + \exp^x)} \right) = a_i^l (1 - a_i^l)$$

那么第三层权值参数的梯度计算为：

$$\frac{\partial J(W, b)}{\partial W_{11}^{(3)}} = \left[\frac{1}{m} \sum_1^m \frac{\partial J(W, b; x^i, y^i)}{\partial W_{11}^{(3)}} \right] + \lambda W_{11}^{(3)}$$

$$\frac{\partial J(W, b; x^i, y^i)}{\partial W_{11}^{(3)}} = \delta_1^{(3)} \cdot a_1^2 = \frac{\partial J}{\partial a_1^{(3)}} \cdot \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \cdot a_1^2 = (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_1^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{11}^{(3)}} = (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_1^{(2)} + \lambda W_{11}^{(3)}$$

同理：

$$\frac{\partial J(W, b)}{\partial W_{12}^{(3)}} = (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_2^{(2)} + \lambda W_{12}^{(3)}$$

$$\frac{\partial J(W, b)}{\partial W_{13}^{(3)}} = (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_3^{(2)} + \lambda W_{13}^{(3)}$$

第二层权值参数:

$$\begin{aligned}\frac{\partial J(W, b)}{\partial W_{11}^{(2)}} &= \left[\frac{1}{m} \sum_1^m \frac{\partial J(W, b; x^i, y^i)}{W_{11}^{(2)}} \right] + \lambda W_{11}^{(3)} = W_{11}^{(3)} \delta_1^{(3)} \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot x_1 \\ &= W_{11}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot x_1 + \lambda W_{11}^{(2)}\end{aligned}$$

同理:

$$\frac{\partial J(W, b)}{\partial W_{21}^{(2)}} = W_{12}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_2^{(2)} (1 - a_2^{(2)}) \cdot x_1 + \lambda W_{21}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{31}^{(2)}} = W_{13}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_3^{(2)} (1 - a_3^{(2)}) \cdot x_1 + \lambda W_{31}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{12}^{(2)}} = W_{11}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot x_2 + \lambda W_{12}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{22}^{(2)}} = W_{12}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_2^{(2)} (1 - a_2^{(2)}) \cdot x_2 + \lambda W_{22}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{32}^{(2)}} = W_{13}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_3^{(2)} (1 - a_3^{(2)}) \cdot x_2 + \lambda W_{32}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{13}^{(2)}} = W_{11}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_1^{(2)} (1 - a_1^{(2)}) \cdot x_3 + \lambda W_{13}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{23}^{(2)}} = W_{12}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_2^{(2)} (1 - a_2^{(2)}) \cdot x_3 + \lambda W_{23}^{(2)}$$

$$\frac{\partial J(W, b)}{\partial W_{33}^{(2)}} = W_{13}^{(3)} (a_1^{(3)} - y^i) \cdot a_1^{(3)} (1 - a_1^{(3)}) \cdot a_3^{(2)} (1 - a_3^{(2)}) \cdot x_3 + \lambda W_{33}^{(2)}$$

Problem 2

WGAN论文理解:

此文章与作者第一篇文章《Towards Principled Methods for Training Generative Adversarial Networks》内容联系紧密。相比于原始GAN改进了四处:

1. 判别器最后一层去掉sigmoid
2. 生成器和判别器的loss不取log
3. 每次更新判别器的参数之后把它们的绝对值截断到不超过一个固定常数c
4. 不用基于动量的优化算法 (包括momentum和Adam), 推荐RMSProp, SGD亦可

原始GAN存在的问题:

根据GAN论文中的判别器的loss, 求得最有判别器的形式; 而在最有判别器的条件下, 将生成器loss等式变换为 P_r 与 P_g 之间的JS散度, 在训练Discriminator接近最有过程中, 最小化Generator的loss也会近似于最小化 P_r 与 P_g 之间的JS散度。

$$G_{loss} = 2JS(P_r || P_g) - 2\log 2 \quad (12)$$

当 P_r 与 P_g 的支撑集是高维空间中的低维流形时, P_r 与 P_g 重叠部分测度为0的概率为1。在近似最优D下, 优化G的loss等价于最小化 P_r 与 P_g 的JS散度, 几乎不可能存在不可忽略的重叠, 因此DCGAN与MLPGAN大部分时间的JS散度都是 $\log 2$, 不仅无法指示训练过程, 最终还会导致生成器的梯度近似为0, 梯度消失。

判别器训练得太好, 生成器梯度消失, 生成器loss降不下去; 判别器训练得不好, 生成器梯度不准, 无法收敛。只有判别器训练得不好不坏才行, 但是这个火候又很难把握, 甚至在同一轮训练的前后不同阶段的火候都可能不一样, 所以GAN难训练。

对生成样本和真实样本加噪声, 让其产生重叠, 可以解决训练不稳定的问题。

EM距离的性质: 平滑, 本质是一个路径规划的最低成本问题。

弱对偶、强对偶解决Wasserstein的求解问题

W距离的定义:

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} E_{(x,y) \in \gamma} [\|x - y\|] \quad 13$$

分析上式, 在联合分布函数连续的情况下将公式转变为下式:

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} E_{(x,y) \in \gamma} \int \int \gamma(x, y) d(x, y) dx dy \quad 14$$

分析上式, $d(x, y)$ 式成本函数, 形式为范数, 在此论文中选取1范数, 根据范数等价相容性知范数形式可变。

联合分布函数 $\gamma(x, y)$, 根据联合分布性质有 $\int \gamma(x, y) dy = P_r(x)$ 且 $\int \gamma(x, y) dx = P_g(y)$

原始优化为找出总成本最低的W距离即可。即MIN13式, 在满足如下约束条件下:

$$\int \gamma(x, y) dy = p(x), \int \gamma(x, y) dx = q(y), \gamma(x, y) \geq 0$$

将上式离散化变成内积形式:

$$\mathbf{\Gamma} = \begin{pmatrix} \gamma(\mathbf{x}_1, \mathbf{y}_1) \\ \gamma(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_2, \mathbf{y}_1) \\ \gamma(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ \gamma(\mathbf{x}_n, \mathbf{y}_1) \\ \gamma(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} d(\mathbf{x}_1, \mathbf{y}_1) \\ d(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ d(\mathbf{x}_2, \mathbf{y}_1) \\ d(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ d(\mathbf{x}_n, \mathbf{y}_1) \\ d(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}$$

Figure 2: 离散化内积

$$\underbrace{\begin{pmatrix} 1 & 1 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ 0 & 0 & \dots & 1 & 1 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 1 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ \hline 1 & 0 & \dots & 1 & 0 & \dots & \dots & 1 & 0 & \dots & \dots \\ 0 & 1 & \dots & 0 & 1 & \dots & \dots & 0 & 1 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \\ 0 & 0 & \dots & 0 & 0 & \dots & \dots & 0 & 0 & \dots & \dots \\ \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \ddots & \vdots & \vdots & \ddots & \ddots \end{pmatrix}}_{\mathbf{A}} \underbrace{\begin{pmatrix} \gamma(\mathbf{x}_1, \mathbf{y}_1) \\ \gamma(\mathbf{x}_1, \mathbf{y}_2) \\ \vdots \\ \gamma(\mathbf{x}_2, \mathbf{y}_1) \\ \gamma(\mathbf{x}_2, \mathbf{y}_2) \\ \vdots \\ \vdots \\ \gamma(\mathbf{x}_n, \mathbf{y}_1) \\ \gamma(\mathbf{x}_n, \mathbf{y}_2) \\ \vdots \\ \vdots \end{pmatrix}}_{\mathbf{\Gamma}} = \underbrace{\begin{pmatrix} p(\mathbf{x}_1) \\ p(\mathbf{x}_2) \\ \vdots \\ p(\mathbf{x}_n) \\ \vdots \\ \hline q(\mathbf{y}_1) \\ q(\mathbf{y}_2) \\ \vdots \\ q(\mathbf{y}_n) \\ \vdots \end{pmatrix}}_{\mathbf{b}}$$

Figure 3: 离散化约束

其中 $\mathbf{\Gamma} \geq 0$ W 距离可以用下式描述:

$$\min_{\Gamma} \{ \langle \Gamma, D \rangle \mid A\Gamma = b, \Gamma > 0 \}$$

重写上式

$$\min_x \{ c^T x \mid Ax = b, x \geq 0 \}$$

设存在最小值 x^* 使得上式成立, 则 $Ax^* = b$ 也成立, 在等式两边乘 $y^T \in R^m$ 得到 $y^T Ax^* = y^T b$, 此时设 $y^T A \leq c^T$ 则 $y^T Ax^* \leq c^T x^*$, 也有 $y^T b \leq c^T x^*$ 。整理上述公式变为语言描述: 在条件 $y^T A \leq c^T$ 下, 任意 $y^T b$ 总是不大于 $\min_x \{ c^T x \mid Ax = b, x \geq 0 \}$ 。

所以有:

$$\max_y \{ b^T y \mid A^T y \leq c \} \leq \min_x \{ c^T x \mid Ax = b, x \geq 0 \}$$

再利用Farkas引理证明上式的强对偶形式, 取得等号。

$$\max_y \{ b^T y \mid A^T y \leq c \} = \min_x \{ c^T x \mid Ax = b, x \geq 0 \}$$

通过强对偶为W距离找对偶表达式:

$$\min_{\Gamma} \{ \langle \Gamma, D \rangle \mid A\Gamma = b, \Gamma > 0 \} = \max_F \{ \langle b, F \rangle \mid A^T F \geq D \}$$

处理问题是依然将其转化为离散形式:

$$\langle b, F \rangle = \sum_n p(x_n) f(x_n) + \sum_n q(x_n) g(x_n)$$

约束条件是 $A^T F \leq D$

由约束条件得到:

$$\forall i, j, f(x_i) + g(y_i) \leq d(x_i, y_i)$$

$$\mathbf{F} = \begin{pmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \\ \vdots \\ g(\mathbf{y}_1) \\ g(\mathbf{y}_2) \\ \vdots \\ g(\mathbf{y}_n) \\ \vdots \end{pmatrix}$$

Figure 4: 神经网络拟合函数

进而得到:

$$W\{p, q\} = \max_{f, g} \left\{ \int [p(x)f(x) + q(x)g(x)]dx \mid f(x) + g(y) \leq d(x, y) \right\}$$

由 $d(x, x) = 0$ 得到 $f(x) + g(x) \leq d(x, x) + 0$, 即 $g(x) \leq f(x)$, 所以有:

$$p(x)f(x) + q(x)g(x) \leq p(x)f(x) + q(x)(-f(x)) = p(x)f(x) - q(x)f(x)$$

再次重写W距离:

$$W\{p, q\} = \max_f \left[\int [p(x)f(x) + q(x)(-f(x))]dx \mid f(x) - f(y) \leq d(x, y) \right]$$

由上得到W的对偶形式, 约束条件重写为 $\|f\|_L \leq 1$, 称为利普希茨约束。将上式重写为采样形式就得到WGAN论文中的具体优化公式, 1则为文章中假设的值。同样K作为利普希茨常数亦可。

算法流程图中关于梯度迭代的D为加, G迭代为减。

WGAN: 先最大化去拟合Wasserstein距离, 再去最小化L, 来优化生成器。D也就是Critice不是为了判别, 正如其意思, 是为了评判, 所以拿掉了二分类, 变成了回归问题。目标是优化G。