

Learning Representations that Support Robust Transfer of Predictors

思路是把模型在不同domain的泛化能力直接作为objective优化。提出了TRM(Transfer Risk Minimization)。所谓的Transfer Risk，也就是将一个环境的最优分类器在其他环境上的泛化能力进行平均得到的risk。可以用下面的loss来衡量：

w 是分类器， Φ 是特征表示， Ω 是环境集合， P 、 Q 是对应环境的empirical distribution。对于给定的 Φ ，TRM所提出的泛化能力衡量标准如下：

$R(\Phi; \Omega) = \sum_{Q \in \Omega} (\sup_{P \in \text{Conv}(\Omega/Q)} E_P[l(w(Q; \Phi) \circ \Phi)])$ 其中 $w(Q; \Phi) = \arg \min_w E_Q[l(w \circ \Phi)]$ 是 Φ 在 Q 上的最优分类器， $\text{Conv}(\Omega/Q)$ 是 Ω 中除去 Q 的环境的凸包。

优化：

1、内层优化：需要对抗地找到让 $w(Q; \Phi)$ 性能最差的环境 P ，要搜索凸包

$\text{Conv}(\Omega/Q) = \{ \sum_{P_i \in \Omega/Q} \alpha_i(Q) P_i \mid \alpha_i(Q) \geq 0, \|\alpha(Q)\|_1 = 1 \}$ 中的元素 P ，于是需要对 α 进行优化。做法是对 α_i 进行 Exponential Gradient Ascent。

2、外层优化：需要通过优化 Φ 来让 $w(Q; \Phi)$ 在 $P(Q)$ 上的loss最低。将 $w(Q; \Phi)$ 在 $P(Q)$ 上的loss

$E_{P(Q)}[l(w(Q; \Phi) \circ \Phi)]$ 记为 $L_P(Q)$ 。对 Φ 的全梯度 $\frac{dL_P(Q)}{d\Phi}$ 可以写成两项，
 $\frac{\partial L_P(Q)}{\partial \Phi} + (\frac{\partial L_P(Q)}{\partial w(Q; \Phi)})^T \frac{dw(Q; \Phi)}{d\Phi}$

第二项 (implicit gradient) 可以改写为 $-\frac{\partial (sg((\frac{\partial L_P(Q)}{\partial w(Q)})^T H_{w(Q)}^{-1})^T (\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q))})}{\partial \Phi}$ 其中

$H_{w(Q)} = \frac{\partial^2 E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q)^2}$ ， sg 为 stop_gradient 操作，即让括号内的函数不对 Φ 求导。

证明：注意到， $w(Q; \Phi)$ 是由方程 $\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q)} = 0$ 确定的隐函数。由隐函数存在定理，

$\frac{dw(Q)}{d\Phi} = -\frac{\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial \Phi}}{\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q)}} = -H_{w(Q)}^{-1} \frac{\partial^2 E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q) \partial \Phi}$ 代入 implicit gradient，即可得证。

将(2)代入(1)并对 Φ 积分，得到外层优化的损失函数：

$E_P[l(w(Q) \circ \Phi)] - (sg((\frac{\partial L_P(Q)}{\partial w(Q)})^T H_{w(Q)}^{-1})^T (\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q)})) + C$

第一项保证环境 Q 上的最优分类器 $w(Q)$ 在 P 上的性能；第二项是环境 Q 和 $w(Q)$ 的最差性能环境 $P(Q)$ 对 $w(Q)$ 梯度的 alignment。

第二项相比于 Fish 直接 align 不同 domain 的梯度： $\mathcal{L}_{idgm} = \mathcal{L}_{erm}(\mathcal{D}_{tr}; \theta) - \gamma \frac{2}{S(S-1)} \sum_{i,j \in S}^{i \neq j} G_i G_j$

其中， $G_i = E_{\mathcal{D}_i} \frac{\partial l(x,y;\theta)}{\partial \theta}$ ，区别在于 TRM 加入了 Hessian Inverse 这一项，而且是对最优分类器的导数。

对 Hessian Inverse 的近似：对 H^{-1} 泰勒展开到前项： $H_j^{-1} = \sum_{i=0}^j (I - H)^i$ 可以用递推式计算 H_j^{-1} ： $H_j^{-1} = I + (I - H) H_{j-1}^{-1}$ 递推式的矩阵乘法用参考 Fast Exact Multiplication by the Hessian，可以在线性时间内完成。

算法：

将 (3) 的损失函数加上 ERM 的 loss $E_Q[l(w_{all} \circ \Phi)]$ 用于更新总的分类器 w_{all} ，在给定 $P(Q)$ 的条件下得到环境 Q 的 loss：

$R(\Phi, w_{all}; Q) = E_Q[l(w_{all} \circ \Phi)] + E_P[l(w(Q) \circ \Phi)] - \lambda \frac{\partial (sg((\frac{\partial L_P(Q)}{\partial w(Q)})^T H_{w(Q)}^{-1})^T (\frac{\partial E_Q[l(w(Q) \circ \Phi)]}{\partial w(Q)}))}{\partial \Phi}$