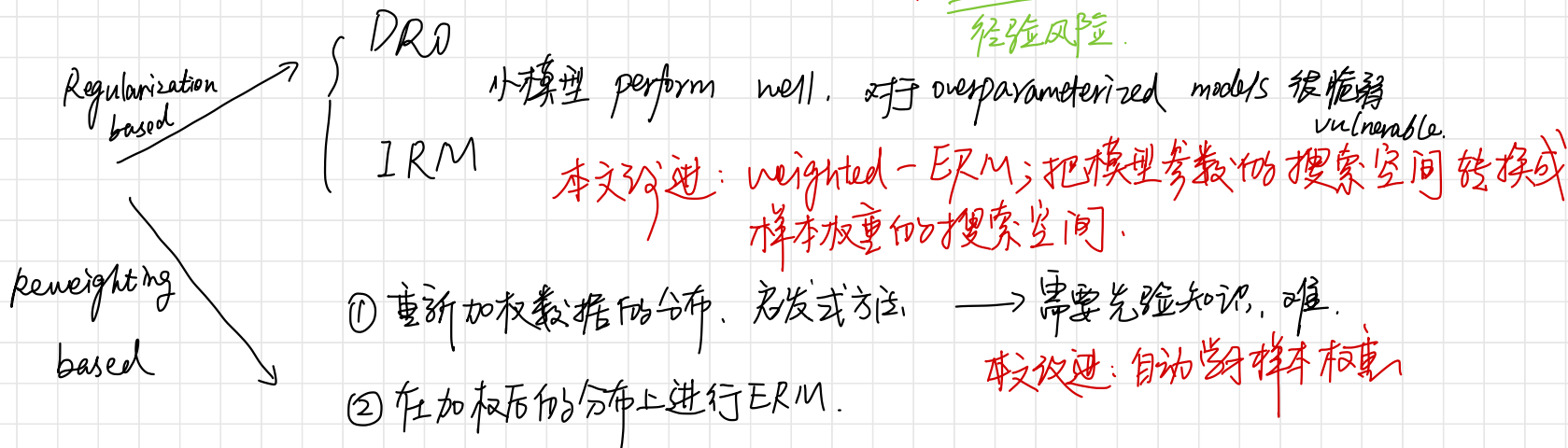


# MAPLE

key idea: 找到一个 reweighting method. 在加权训练数据上对大 model 进行 ERM 可以得到好的 OOD 泛化性能.

在已知训练数据上衡量其性能.  
经验风险.



本文: 用 weighted-ERM 解决正则化方法过参数化的问题.  $\rightarrow$  (过拟合)

自动学习样本权重, 解决重新加权方法的限制.

model agnostic 的两层优化:

内循环: 在加权训练样本上训练 DNN. (ERM 训练), 得到模型  $\theta$ .

外循环: 在验证集上评估的 OOD 标准作为外部目标来学习样本权重

自动权重学习

不容易过拟合的原因: 搜索权重, 而不是模型参数

举例: CIFAR-10 : 50K 训练样本

ResNet-18 : 11.4M 参数

• 在 waterbirds 数据集中比 G-ORO 得分更高.

Question: 权重空间为什么与样本数作为参考?

有可能是因为样本的权重, 对样本操作

数据集  $D := \{(x_i, y_i)\}_{i=1}^n, (x_i, y_i) \in X \times Y$

加权经验损失  $L(D, \theta; w) := \frac{1}{n} \sum_{i=1}^n w_i l(f(x_i; \theta), y_i)$

$L(D, \theta)$ : 无偏损失  $= L(D, \theta; \mathbf{1})$

$Z_c$ : bias-conflict (core)

$Z_s$ : bias-aligned (spurious)

Validation  $\longrightarrow$  相比于正则化方法: 使用验证集缓解训练集的过拟合

$\min_{w \in C} R(D_v, \theta^*(w))$

s.t.  $\theta^*(w) \in \underset{\theta}{\operatorname{argmin}} L(D_{\text{tr}}, \theta; w)$   $C = \{w: w \geq 0, \|w\| \leq k\}$

$w \leftarrow \operatorname{proj}_C(w - \eta \nabla_w R|_{\theta^*} \frac{\partial L}{\partial w} \Big|_{\theta^*})$

ERM loss  $\hat{\theta}_{\text{ERM}} := \underset{\theta \in \Theta}{\operatorname{argmin}} \mathbb{E}_{(x,y) \sim p} [l(\theta; (x,y))]$

训练集上的经验分布.

sparse. 减少 cost

$\min_{(w,s) \in C'} \Phi(w,s) = \mathbb{E}_{p(m|s)} R(D_v, \theta^*(w,m))$

$C' = \{(w,s): w \geq 0, 0 \leq s \leq 1, \|s\| \leq k\}$

s.t.  $\theta^*(w,m) \in \underset{\theta}{\operatorname{argmin}} L(D_{\text{tr}}, \theta; w \circ m)$

$m_i$ : Bernoulli Random Variable.

$\nabla_{w,s} \Phi \approx \nabla_{w,s} R(\theta^*(w, \mathbf{1}(\log(\frac{s}{1-s}) + g, -g_0 \geq 0)))$

Gumbel(0,1) 中随机选取两个随机变量

$p(m_i=1)=s_i, p(m_i=0)=1-s_i$   
 $p(m|s) = \prod_{i=1}^m (s_i)^{m_i} (1-s_i)^{(1-m_i)}$

$(w,s) \leftarrow \operatorname{proj}_C(w - \eta \nabla_w \Phi, s - \eta \nabla_s \Phi)$  projected gradient descent.

Datasets: Colored MNIST, CIFAR10, MNIST, Waterbirds, CelebA

validate MAPLE on DRO

Baseline: ERM, IRMv1, REx, MRM, Sparse IRM, Bayesian IRM  $\rightarrow$  IRM 系列

---

ERM, CVaRDRO, LfF, JTT, Up Weighting, Group DRO  $\rightarrow$  DRO 系列.

---

**Algorithm 1** Model Agnostic Sample Reweighting (MAPLE)

---

**Input:** a network  $\theta$ , remaining training sample size  $K$ , training set  $\mathcal{D}_{tr}$  and validation set  $\mathcal{D}_v$ .

- 1: Initialize sample weights  $w = \mathbf{1}$  and probabilities  $s = \frac{K}{|\mathcal{D}_{tr}|} \mathbf{1}$ .
- 2: **for** training iteration  $i = 1, 2 \dots I$  **do**
- 3:   Sample mask  $m$  according to the probability distribution  $p(m|s) = \prod_{i=1}^n (s_i)^{m_i} (1 - s_i)^{(1-m_i)}$ . After Sparse.
- 4:   Train the inner loop to converge:  $\theta^*(w, m) \leftarrow \arg \min_{\theta} \mathcal{L}(\mathcal{D}_{tr}, \theta; w, m)$  started from randomly initialized  $\theta$ .
- 5:   Estimate  $\nabla_s \Phi(w, s)$  and  $\nabla_w \Phi(w, s)$  by Straight-through Gumbel-softmax and 1-step truncated backpropagation.
- 6:   Perform projected gradient descent:  $(w, s) \leftarrow \text{proj}_{\mathcal{C}}(w - \eta \nabla_w \Phi(w, s), s - \eta \nabla_s \Phi(w, s))$
- 7: **end for**

**output** The weighted set  $\{(\mathbf{x}_i, \mathbf{y}_i, w_i) : m_i \neq 0, (\mathbf{x}_i, \mathbf{y}_i) \in \mathcal{D}_{tr}\}$  with  $m$  sampled from  $p(m|s)$

---