BINARY INSTANCE-SPECIFIC COST LEARNING

- Assumes instance-specific costs for every observation: $\mathcal{D}^{(n)} = \{(\mathbf{x}^{(i)}, \mathbf{c}^{(i)})\}_{i=1}^{n}$, where $(\mathbf{x}^{(i)}, \mathbf{c}^{(i)}) \in \mathbb{R}^{p} \times \mathbb{R}^{2}$.
- Define "true class" as cost minimal class
- Define observation weights: $|\mathbf{c}^{(i)}[1] \mathbf{c}^{(i)}[0]|$

	c ⁽ⁱ⁾ [0]	c ⁽ⁱ⁾ [1]	y ⁽ⁱ⁾	w ⁽ⁱ⁾
x ⁽¹⁾	1	1	0	0
x ⁽²⁾	1	2	0	1
x ⁽³⁾	7	3	1	4

× 0 0 × × ×

• Now solve weighted ERM:

$$\mathcal{R}_{emp}(\boldsymbol{\theta}) = \sum_{i=1}^{n} w^{(i)} L\left(y^{(i)}, f\left(\mathbf{x}^{(i)} \mid \boldsymbol{\theta}\right)\right)$$

• NB: Instances with equal costs are effectively ignored.

MULTICLASS COSTS

 Consider g > 2. Vanilla CSL is special case of instance specific, use c⁽ⁱ⁾ same for all x⁽ⁱ⁾ of the same class



• For two $\mathbf{x}^{(i)}$ with y = 2 and y = 3:

	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]	y ⁽ⁱ⁾
x ⁽¹⁾	1	0	1	2
x ⁽²⁾	3	1	0	3
x ⁽³⁾	1	0	1	2

• Set $\mathbf{c}^{(i)}[y^{(i)}] = 0$, i.e. zero-cost for correct prediction.

CSOVO LIN ET AL. 2014

• Let
$$\mathcal{D}^{(n)} = \{(\mathbf{x}^{(i)}, \mathbf{c}^{(i)})\}_{i=1}^n, (\mathbf{x}^{(i)}, \mathbf{c}^{(i)}) \in \mathbb{R}^p \times \mathbb{R}^g$$
.

• Example:

	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]
x ⁽¹⁾	0	2	3
x ⁽²⁾	1	0	1
x ⁽³⁾	2	0	3

× 0 0 × × ×

- Idea: Reduction principle to binary case (weighted fit) by one-versus-one (OVO).
- For class *j* vs. *k*:
 - How to deal with the label $y^{(i)}$? $y^{(i)}$ can be neither *j* nor *k*.
 - How to deal with the costs $\mathbf{c}^{(i)}[j]$ and $\mathbf{c}^{(i)}[k]$?

CSOVO

- When training a binary classifier $f^{(j,k)}$ for class *j* vs. *k*,
 - Choose cost min class from pair arg min<sub>*l*∈{*j*,*k*} c^(*i*)[*l*] as ground truth
 </sub>
 - Sample weight is simply diff between the 2 costs
 |c⁽ⁱ⁾[j] c⁽ⁱ⁾[k]|
- Example continued:

Example continued.							
	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]	c ⁽ⁱ⁾ [1 vs 2]	$ ilde{y}^{(i)}$ [1 vs 2]	w ⁽ⁱ⁾ [1 vs 2]	
x ⁽¹⁾	0	2	3	0/2	1	2	
x ⁽²⁾	1	0	1	1/0	2	1	
x ⁽³⁾	2	0	3	2/0	2	2	
		-	Ũ	=, •	-	_	
	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]	c ⁽ⁱ⁾ [2 vs 3]	<i>ỹ</i> ⁽ⁱ⁾ [2 vs 3]	w ⁽ⁱ⁾ [2 vs 3]	
x ⁽¹⁾	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]	c ⁽ⁱ⁾ [2 vs 3] 2/3	<i>ỹ</i> ⁽ⁱ⁾ [2 vs 3] 2	<i>w</i> ^(<i>i</i>) [2 vs 3] 1	
x ⁽¹⁾ x ⁽²⁾	c ^(<i>i</i>) [1] 0 1	c ⁽ⁱ⁾ [2] 2 0	c ⁽ⁱ⁾ [3] 3 1	c ⁽ⁱ⁾ [2 vs 3] 2/3 0/1	$\tilde{y}^{(i)}$ [2 vs 3] 2 2	w ⁽ⁱ⁾ [2 vs 3] 1 1	

× 0 0 × × ×

CSOVO

• Example continued

	c ⁽ⁱ⁾ [1]	c ⁽ⁱ⁾ [2]	c ⁽ⁱ⁾ [3]	c ⁽ⁱ⁾ [1 vs 3]	$ ilde{y}^{(i)}$ [1 vs 3]	<i>w</i> ⁽ⁱ⁾ [1 vs 3]
x ⁽¹⁾	0	2	3	0/3	1	3
x ⁽²⁾	1	0	1	-/-	-	0
x ⁽³⁾	2	0	3	2/3	1	1

- Wrap everything up:
 - For class *j* vs. *k*, transform all $(\mathbf{x}^{(i)}, \mathbf{c}^{(i)})$ to $(\mathbf{x}^{(i)}, \arg\min_{l \in \{j,k\}} \mathbf{c}^{(i)}[l])$ with sample-wise weight $|\mathbf{c}^{(i)}[j] \mathbf{c}^{(i)}[k]|$.
 - 2 Train a weighted binary classifier $f^{(j,k)}$ using the above
 - Solution Repeat step 1 and 2 for different (j, k).
 - Predict using the votes from all $f^{(j,k)}$.
- Theoretical guarantee:

test costs of final classifier $\leq 2 \sum_{j < k}$ test cost of $f^{(j,k)}$.

0 0 X X 0 X X