

Large scale greedy feature-selection for multi-target learning

Antti Airola, Tapio Pahikkala et al.

ECML 2015 BigTargets Workshop

Joint work with many authors

- University of Turku: Antti Airola, Pekka Naula, Tapio Pahikkala, Tapio Salakoski (Multi-target greedy RLS)

- Large scale feature selection for multi-target learning
- Task: select minimal set of common features allowing accurate predictions over target tasks
- Greedy RLS: greedy regularized least-squares
- Linear time ($\#$ inputs, $\#$ features, $\#$ outputs, $\#$ selected)
- Highlights from experiments
 - Broad-DREAM Gene Essentiality Prediction Challenge
 - Outperforms multi-task Lasso for small feature budgets
- Also scales to full Genome Wide Association Studies; thousands of samples, hundreds of thousands of features (recent PhD thesis: Sebastian Okser)

Why feature selection?

- ① Accuracy: regularizing effect, avoiding overfitting leads to better generalization
- ② Interpretability: obtain a small set of features understandable by human expert
- ③ Budget constraints: obtaining features costs time and money

Model sparsity

$$\mathbf{W}_1 = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 3 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 3 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 2 & 0 \end{pmatrix}, \quad \mathbf{W}_2 = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 2 & 3 & -1 & 2 \\ 0 & 0 & 0 & 0 \\ 3 & 1 & 4 & 1 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

- features \times targets coefficient matrices
- \mathbf{W}_1 8 features needed for prediction
- \mathbf{W}_2 2 features needed for prediction

Least-squares formulation

$$\begin{aligned} \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times t}} & \|\mathbf{XW} - \mathbf{Y}\|_F^2 \\ & \text{subject to } C(\mathbf{W}) \end{aligned}$$

Notation

X	data matrix
Y	output matrix
W	model coefficients
$\ \cdot\ _F$	Frobenius norm
$C(\cdot)$	Constraint (regularizer)

Multi-task Lasso (baseline)

Multi-task Lasso (Zhang, 2006)

$$\begin{aligned} & \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times t}} \|\mathbf{XW} - \mathbf{Y}\|_F^2 \\ & \text{subject to } \sum_{i=1}^d \max_j |\mathbf{W}_{i,j}| \leq r \end{aligned}$$

- $L_{1,\infty}$ norm enforces sparsity in the number of features
- $r > 0$ regularization parameter

Greedy RLS (proposed)

$$\begin{aligned} & \arg \min_{\mathbf{W} \in \mathbb{R}^{d \times t}} \|\mathbf{XW} - \mathbf{Y}\|_F^2 \\ & \text{subject to } \|\mathbf{W}\|_F^2 < r \quad \text{and} \\ & |\{i \mid \exists j, \mathbf{W}_{i,j} \neq 0\}| \leq k \end{aligned}$$

- $r > 0$ regularization parameter
- $k > 0$ constraint on the number of features
- heuristics needed to search over the power set of features

- Greedy regularized least-squares (Greedy RLS)
- Starting from empty feature set, at each point add the feature reducing leave-one-out cross-validation error most
- Stop once k features have been selected

Algorithm 1 Multi-target greedy RLS

```
1:  $\mathcal{S} \leftarrow \emptyset$  ▷ selected features common for all tasks
2: while  $|\mathcal{S}| < k$  do ▷ select  $k$  features
3:    $e \leftarrow \infty$ 
4:    $b \leftarrow 0$ 
5:   for  $i \in \{1, \dots, d\} \setminus \mathcal{S}$  do ▷ test all features
6:      $e_{avg} \leftarrow 0$ 
7:     for  $j \in \{1, \dots, t\}$  do
8:        $e_{i,j} \leftarrow \mathcal{L}(\mathbf{X}_{:,S \cup \{i\}}, \mathbf{Y}_{:,j})$  ▷ LOO for task  $j$ 
9:        $e_{avg} \leftarrow e_{avg} + e_{i,j}/t$ 
10:    if  $e_{avg} < e$  then
11:       $e \leftarrow e_{avg}$ 
12:       $b \leftarrow i$ 
13:     $\mathcal{S} \leftarrow \mathcal{S} \cup \{b\}$  ▷ feature with lowest LOO-error
14:  $\mathbf{W} \leftarrow \mathcal{A}(\mathbf{X}_{:, \mathcal{S}}, \mathbf{Y})$  ▷ train final models
15: return  $\mathbf{W}, \mathcal{S}$ 
```

- Greedy RLS could be implemented as a general wrapper code calling a black-box solver
- $\#selected \times \#features \times \#targets \times \#CV\text{-rounds}$ calls for naive implementation!
- Matrix algebraic optimization for feature addition, leave-one-out... (for all targets simultaneously)
- Linear time algorithm ($\#inputs$, $\#features$, $\#outputs$, $\#selected$)
- P. Naula, A. Airola, T. Salakoski and T. Pahikkala. Multi-label learning under feature extraction budgets. *Pattern Recognition Letters*, 2014.

Algorithm 2 Multi-target greedy RLS

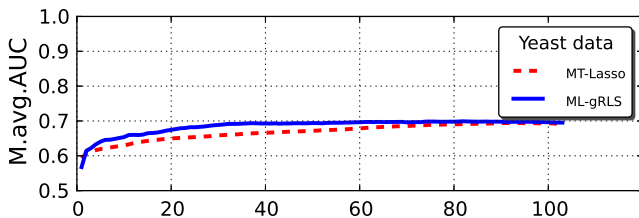
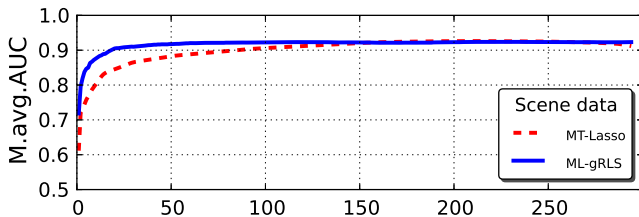
```
A  $\leftarrow \lambda^{-1} \mathbf{Y}$ 
g  $\leftarrow \lambda^{-1} \mathbf{1}$ 
C  $\leftarrow \lambda^{-1} \mathbf{X}$ 
S  $\leftarrow \emptyset$ 
while  $|S| < k$  do
  e  $\leftarrow \infty$ 
  b  $\leftarrow 0$ 
  for  $i \in \{1, \dots, d\} \setminus S$  do
    u  $\leftarrow \mathbf{C}_{:,i} (1 + (\mathbf{X}_{:,i})^T \mathbf{C}_{:,i})^{-1}$ 
    ei  $\leftarrow 0$ 
    A  $\leftarrow \mathbf{A} - \mathbf{u} ((\mathbf{X}_{:,i})^T \mathbf{A})$ 
    for  $h \in \{1, \dots, t\}$  do
      for  $j \in \{1, \dots, n\}$  do
        gj  $\leftarrow \mathbf{g}_j - \mathbf{u}_j \mathbf{C}_{j,i}$ 
        ei  $\leftarrow e_i + (\mathbf{g}_j)^{-2} (\tilde{\mathbf{A}}_{j,h})^2$ 
      if  $e_i < e$  then
        e  $\leftarrow e_i$ 
        b  $\leftarrow i$ 
  u  $\leftarrow \mathbf{C}_{:,b} (1 + (\mathbf{X}_{:,b})^T \mathbf{C}_{:,b})^{-1}$ 
  A  $\leftarrow \mathbf{A} - \mathbf{u} ((\mathbf{X}_{:,b})^T \mathbf{A})$ 
  for  $j \in \{1, \dots, n\}$  do
    gj  $\leftarrow \mathbf{g}_j - \mathbf{u}_j \mathbf{C}_{j,b}$ 
  C  $\leftarrow \mathbf{C} - \mathbf{u} ((\mathbf{X}_{:,b})^T \mathbf{C})$ 
  S  $\leftarrow S \cup \{b\}$ 
W  $\leftarrow (\mathbf{X}_{:,S})^T \mathbf{A}$ 
```

Benchmarking greedy RLS and multi-task Lasso

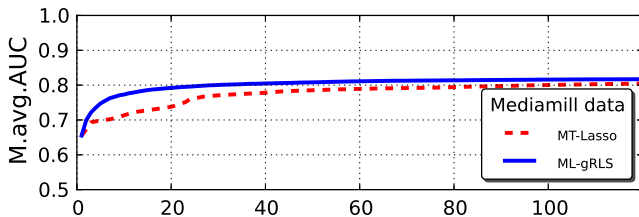
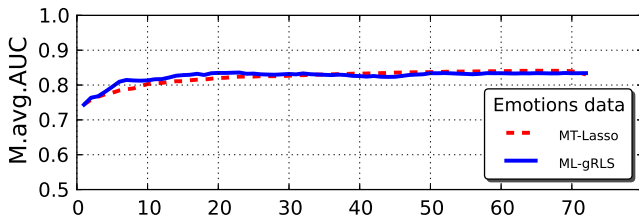
Table: Mulan datasets (Tsoumakas et al. 2011).

Data sets	domain	labels	features	instances
Scene	image	6	294	2407
Yeast	biology	14	103	2417
Emotions	music	6	72	593
Mediamill*	text	9	120	41583
Delicious	text	983	500	16105
Tmc2007	text	22	49060	28596

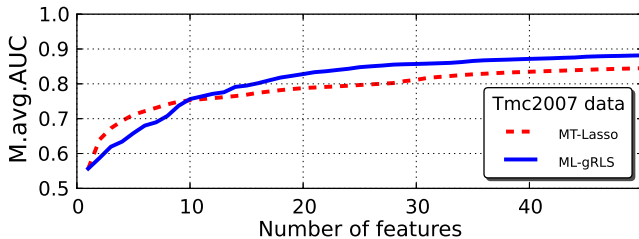
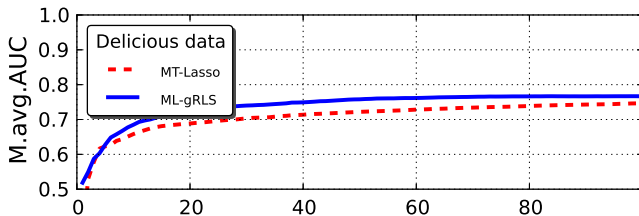
Greedy RLS vs. Lasso



Greedy RLS vs. Lasso



Greedy RLS vs. Lasso



- Greedy RLS: linear time algorithm for (multi-target) feature selection
- Selects joint features for the target tasks
- Competitive, when number of features to be selected small
- Applications on Genome-Wide Association Studies
- RLScore open source implementation at <https://github.com/aatapa/RLScore>