

# KRONECKER KERNEL RIDGE REGRESSION

- In MTP with target features, we often use kernel methods.
- Consider the following pairwise model representation in the primal:

$$f(\mathbf{x}, \mathbf{t}) = \omega^\top (\phi(\mathbf{x}) \otimes \psi(\mathbf{t})),$$

## Multi-Target Prediction: Methods Part 2

where  $\phi$  is feature mapping for features and  $\psi$  is feature mapping for target (features) and  $\otimes$  is Kronecker product.

- This yields Kronecker product pairwise kernel in the dual:

$$f(\mathbf{x}, \mathbf{t}) = \sum_{(\mathbf{x}', \mathbf{t}') \in \mathcal{D}} \alpha_{(\mathbf{x}', \mathbf{t}')} \cdot k(\mathbf{x}, \mathbf{x}') \cdot g(\mathbf{t}, \mathbf{t}')$$

### Learning goals

- Kronecker kernel ridge regression
- Graph relations in targets
- Graph relations in features
- Low-rank approximations

where  $k$  is kernel for feature map  $\phi$ ,  $g$  kernel for feature map  $\psi$  and  $\alpha_{(\mathbf{x}', \mathbf{t}')}$  are dual parameters determined by:

$$\min_{\alpha} \|\Gamma \alpha - \mathbf{z}\|_2^2 + \lambda \alpha^\top \Gamma \alpha, \text{ where } \mathbf{z} = \text{vec}(Y)$$

- Commonly used in zero-shot learning.



# EXPLOITING RELATIONS IN REGULARIZATION

- In MTP with target features, we often use kernel methods.
- Consider the following pairwise model representation in the primal:



where  $\phi$  is feature mapping for features and  $\psi$  is feature mapping for target (features) and  $\otimes$  is Kronecker product.

- This yields Kronecker product pairwise kernel in the dual:
- Graph-based regularization for graph-type relations in targets:

$$f(\mathbf{x}, \mathbf{t}) = \min_{\Theta} \|Y - \Phi\Theta\|_F^2 + \lambda \sum_{m=1} \sum_{m' \in \mathcal{N}(m)} \|\theta_{m'} - \theta_m\|^2 = \sum_{(\mathbf{x}, \mathbf{t}) \in \mathcal{D}} \alpha_{(\mathbf{x}, \mathbf{t})} \Gamma((\mathbf{x}, \mathbf{t}), (\mathbf{x}', \mathbf{t}')),$$

where  $k$  is kernel for feature map  $\phi$ ,  $g$  kernel for feature map  $\psi$  where  $\mathcal{N}(j)$  is the set of targets related to target  $j$  and  $\alpha_{(\mathbf{x}, \mathbf{t})}$  are dual parameters determined by:

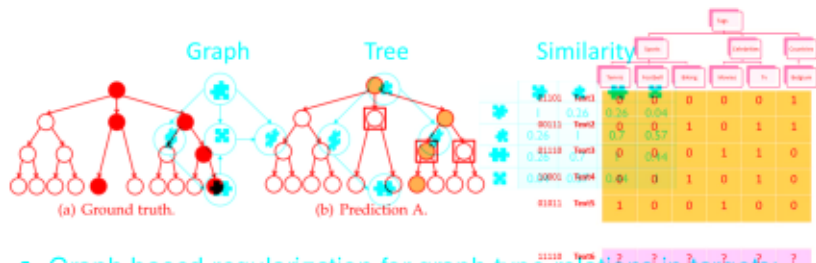
- The graph or tree is given as prior information.
- Can be extended to a weighted version aware of the similarities
- Commonly used in zero-shot learning

Gopal and Yang, Recursive regularization for large-scale classification with hierarchical and graphical dependencies, KDD 2013.

Stock et al., A comparative study of pairwise learning methods based on kernel ridge regression, Neural Computation 2018.



# HIERARCHICAL MULTI-LABEL CLASSIFICATION



- Graph-based regularization for graph-type relations in targets:

- Hierarchies can also be used to define specific loss functions, such as the Hierarchy-loss:

$$\min_{\Theta} \|Y - \Phi\Theta\|_F + \lambda \sum_j \sum_{m' \in \mathcal{N}(j)} \|\theta_m - \theta_{m'}\|^2,$$

$$L_{Hier}(\mathbf{y}, f) = \sum_{m, m' \in \mathcal{N}(j)} c_m \mathbb{I}_{[anc(y_m) = anc(\hat{y}_{m'})]},$$

where  $\mathcal{N}(j)$  is the set of targets related to target  $j$ .

- The graph or tree is given as prior information.
- This is rather common in multi-label classification problems.
- Can be extended to a weighted version aware of the similarities

Bi and Kwok, Bayes-optimal hierarchical multi-label classification, IEEE Transactions on Knowledge and Data Engineering, 2014.

Gopal and Yang, Recursive regularization for large-scale classification with hierarchical and graphical dependencies, KDD 2013.

# PROBABILISTIC CLASSIFIER CHAINIFICATION

- Estimate the joint conditional distribution  $P(\mathbf{y} | \mathbf{x})$ .

- For optimizing the subset 0/1 loss:



	Year1	Year2	Year3	Year4	Year5	Year6
01100 Year1	0	0	0	0	0	1
00111 Year2	0	0	1	0	1	1
10100 Year3	0	0	0	1	1	0
10000 Year4	0	0	1	0	1	0
01010 Year5	1	0	0	1	0	0
11110 Year6	?	?	?	?	?	?

- Repeatedly apply the *product rule* of probability:

- Hierarchies can be used to define *specific* loss functions, such as the Hierarchy-loss:

$$L_{\text{Hier}}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{m=1}^l c_m \mathbb{1}_{[\text{anc}(y_m) \neq \text{anc}(\hat{y}_m)]},$$

$$P(\mathbf{y}_m | \mathbf{x}, y_1, \dots, y_{m-1}),$$

- This is rather common in multi-label classification problems. independently for each  $m = 1, \dots, l$ .

Bi and Kwok, Bayes-optimal hierarchical multi-label classification, IEEE Transactions on Knowledge and Data Engineering, 2014.

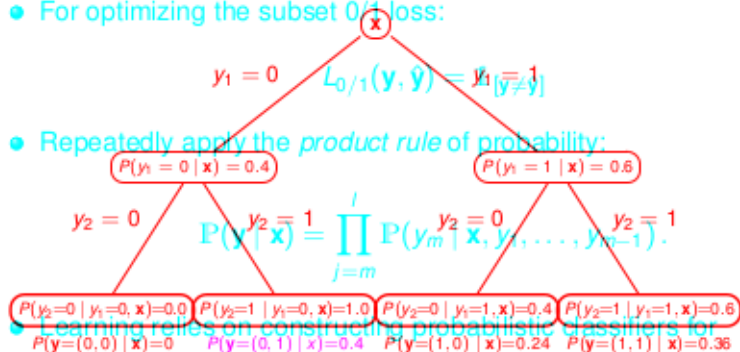


# PROBABILISTIC CLASSIFIER CHAINS

- Inference relies on exploiting a probability tree  $\mathbb{P}(\mathbf{y} | \mathbf{x})$ .

- For optimizing the subset 0/1 loss:

- Repeatedly apply the *product rule of probability*:



- Learning relies on constructing probabilistic classifiers for

- For subset 0/1 loss one needs to find  $h(\mathbf{x}) = \arg \max_{\mathbf{y}} \mathbb{P}(\mathbf{y} | \mathbf{x})$ .
- Greedy and approximate search techniques with guarantees exist.
- Other losses: compute the prediction on a sample from  $\mathbb{P}(\mathbf{y} | \mathbf{x})$ .

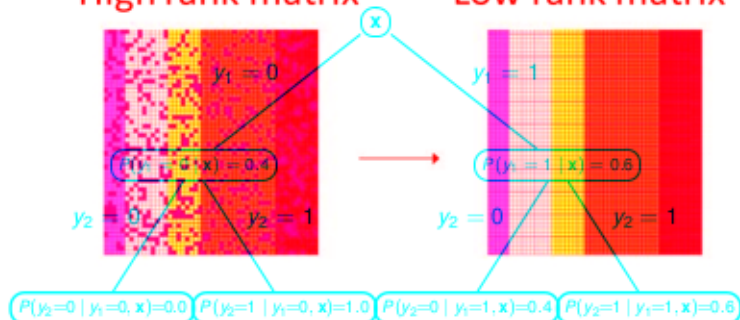


# LOW-RANK APPROXIMATION CHAINS

- Inference relies on exploiting a probability tree:

High rank matrix

Low rank matrix



- Low rank = some structure is shared across targets

- Typically perform low-rank approx of param matrix:

- For subset 0/1 loss one needs to find  $h(\mathbf{x}) = \arg \max_{\mathbf{y}} \mathbb{P}(\mathbf{y} | \mathbf{x})$ .

- Greedy and approximate search techniques with guarantees exist.

- Other losses: compute the prediction on a sample from  $\mathbb{P}(\mathbf{y} | \mathbf{x})$ .

Chen et al., A convex formulation for learning shared structures from multiple tasks, ICML 2009.

Dembczynski et al., An analysis of chaining in multi-label classification, ECAI 2012.



# LOW-RANK APPROXIMATION

- $\Theta$ : parameter matrix of dimensionality  $p \times l$

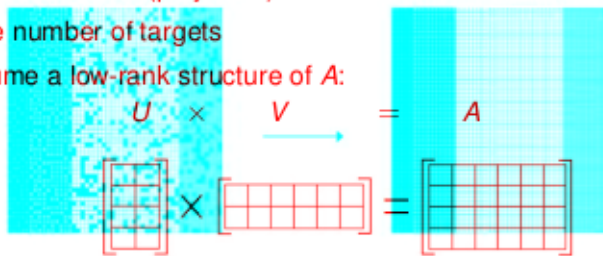
High rank matrix

Low rank matrix

- $p$ : the number of (projected) features

- $l$ : the number of targets

- Assume a low-rank structure of  $A$ :



- We can write  $\Theta = UV$  and  $\Theta \mathbf{x} = UV\mathbf{x}$
- Low rank = some structure is shared across targets
- $V$  is a  $p \times \hat{l}$  matrix
- Typically perform low-rank approx of param matrix:
- $U$  is an  $\hat{l} \times l$  matrix

- $\hat{l}$  is the rank of  $\Theta$   $\min_{\Theta} \|Y - \Phi\Theta\|_F^2 + \lambda \text{rank}(\Theta)$



# LOW-RANK APPROXIMATION

- $\Theta$ : parameter matrix of dimensionality  $p \times l$
- $p$ : the number of (projected) features
- $l$ : the number of targets
- Assume a low-rank structure of  $A$ :

$$U \times V = A$$



- We can write  $\Theta = UV$  and  $\Theta \mathbf{x} = UV\mathbf{x}$
- $V$  is a  $p \times \hat{l}$  matrix
- $U$  is an  $\hat{l} \times l$  matrix
- $\hat{l}$  is the rank of  $\Theta$

