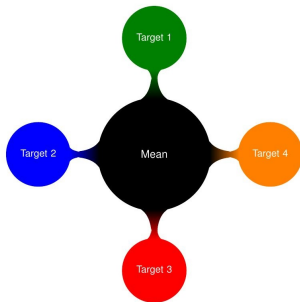


Advanced Machine Learning

Multi-Target Prediction: Methods Part 2



Learning goals

- Know how to leveraging and constructing the target similarity in multi-target learning

KRONECKER KERNEL RIDGE REGRESSION

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

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

where ϕ is feature mapping for features and ψ is feature mapping for target (features) and \otimes is Kronecker product. **TODO: Define \mathbf{t}**

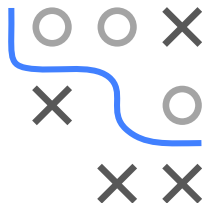
- This yields Kronecker product pairwise kernel in the dual:

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

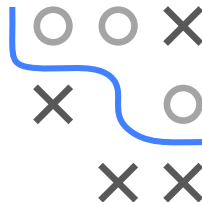
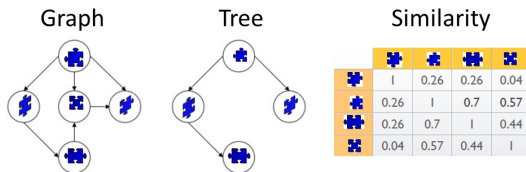
where k is kernel for feature map ϕ , g kernel for feature map ψ and $\alpha_{(\mathbf{x}', \mathbf{t}')}$ are dual parameters determined by:

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

- Commonly used in zero-shot learning.



EXPLOITING RELATIONS IN REGULARIZATION TERMS



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

$$\min_{\Theta} \|Y - \Phi\Theta\|_F^2 + \lambda \sum_{m=1}^I \sum_{m' \in \mathcal{N}(m)} \|\theta_m - \theta_{m'}\|^2,$$

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

- The graph or tree is given as prior information.
- Can be extended to a weighted version aware of the similarities (or correlations).

PROBABILISTIC CLASSIFIER CHAINS

- Estimate the joint conditional distribution $\mathbb{P}(\mathbf{y} \mid \mathbf{x})$.
- For optimizing the subset 0/1 loss:

$$L_{0/1}(\mathbf{y}, \hat{\mathbf{y}}) = \mathbb{1}_{[\mathbf{y} \neq \hat{\mathbf{y}}]}$$

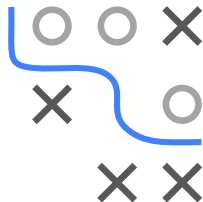
- Repeatedly apply the *product rule* of probability:

$$\mathbb{P}(\mathbf{y} \mid \mathbf{x}) = \prod_{j=m}^l \mathbb{P}(y_m \mid \mathbf{x}, y_1, \dots, y_{m-1}).$$

- Learning relies on constructing probabilistic classifiers for estimating

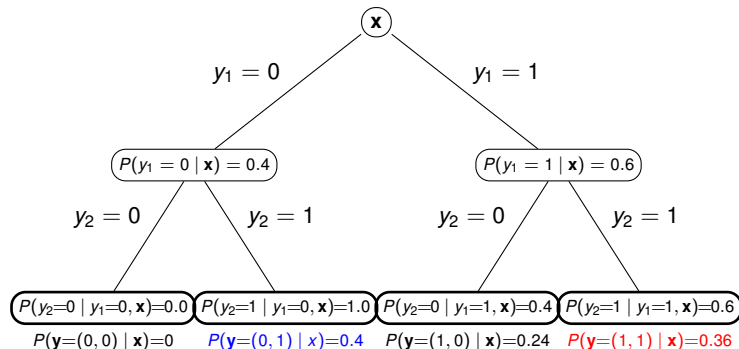
$$\mathbb{P}(y_m \mid \mathbf{x}, y_1, \dots, y_{m-1}),$$

independently for each $m = 1, \dots, l$.



PROBABILISTIC CLASSIFIER CHAINS

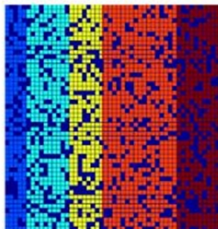
- Inference relies on exploiting a probability tree:



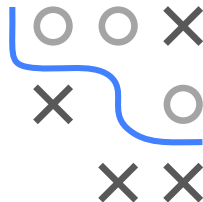
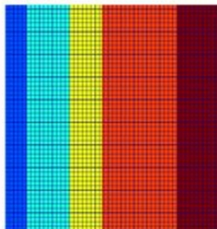
- For subset 0/1 loss one needs to find $f^*(\mathbf{x}) = \arg \max_{\mathbf{y}} \mathbb{P}(\mathbf{y} | \mathbf{x})$.
- Greedy and approximate search techniques with guarantees exist.
- Other losses: compute the prediction on a sample from $\mathbb{P}(\mathbf{y} | \mathbf{x})$.

LOW-RANK APPROXIMATION

High rank matrix



Low rank matrix



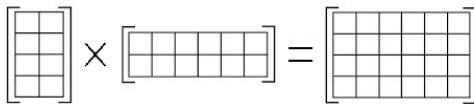
- Low rank materializes the idea that some structure is shared across different targets.
- Typically perform a low-rank approximation of the parameter matrix:

$$\min_{\Theta} \|Y - \Phi\Theta\|_F^2 + \lambda \text{rank}(\Theta)$$

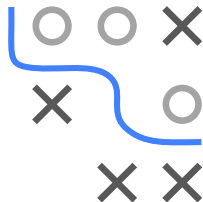
LOW-RANK APPROXIMATION

- Θ : 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 Θ



OVERVIEW OF METHODS

- Popular for multi-output regression, multi-task learning and multi-label classification.
- Linear as well as nonlinear methods.
- Algorithms:
 - Principal component analysis, Canonical correlation analysis, Partial least squares.
 - Singular value decomposition, Alternating structure optimization.
 - Compressed sensing, Output codes, Landmark labels, Bloom filters, Auto-encoders.

