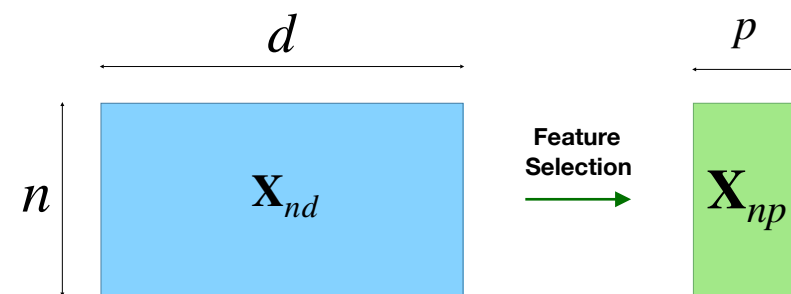


Feature Selection in tumor profiles

Given a gene signature tumor profiles can be classified by low and high survival rate. Tumors are characterized by more than $d = 20,000$ feature genes and cancer datasets use to have a low sample n size where $n \ll d$. To avoid the curse of dimensionality and to guide biomarker discovery feature selection methods are used and a subset of p genes are selected.

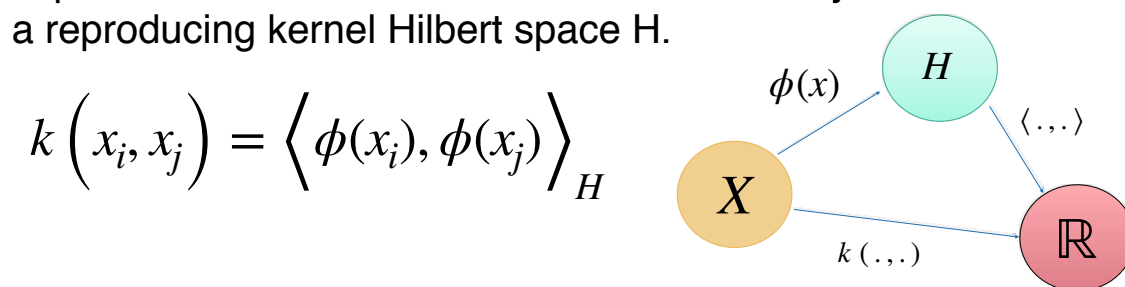


Instead of learning a selection function $f(x)$ from a supervised approach $y = f(x)$ which may be affected by overfitting we propose to learn a selection function from the labels and from the latent structure \mathbf{z} of the training data simultaneously via kernel methods. The objective is to regularize the selection and to improve generalization in classification tasks.

$$y = \hat{f}(x) \rightarrow (z, y) = \hat{f}(x)$$

Kernel Methods

In this work Kernel functions are used to mix latent structure with supervised labels. A Kernel measures similarity between vectors in a reproducing kernel Hilbert space H .



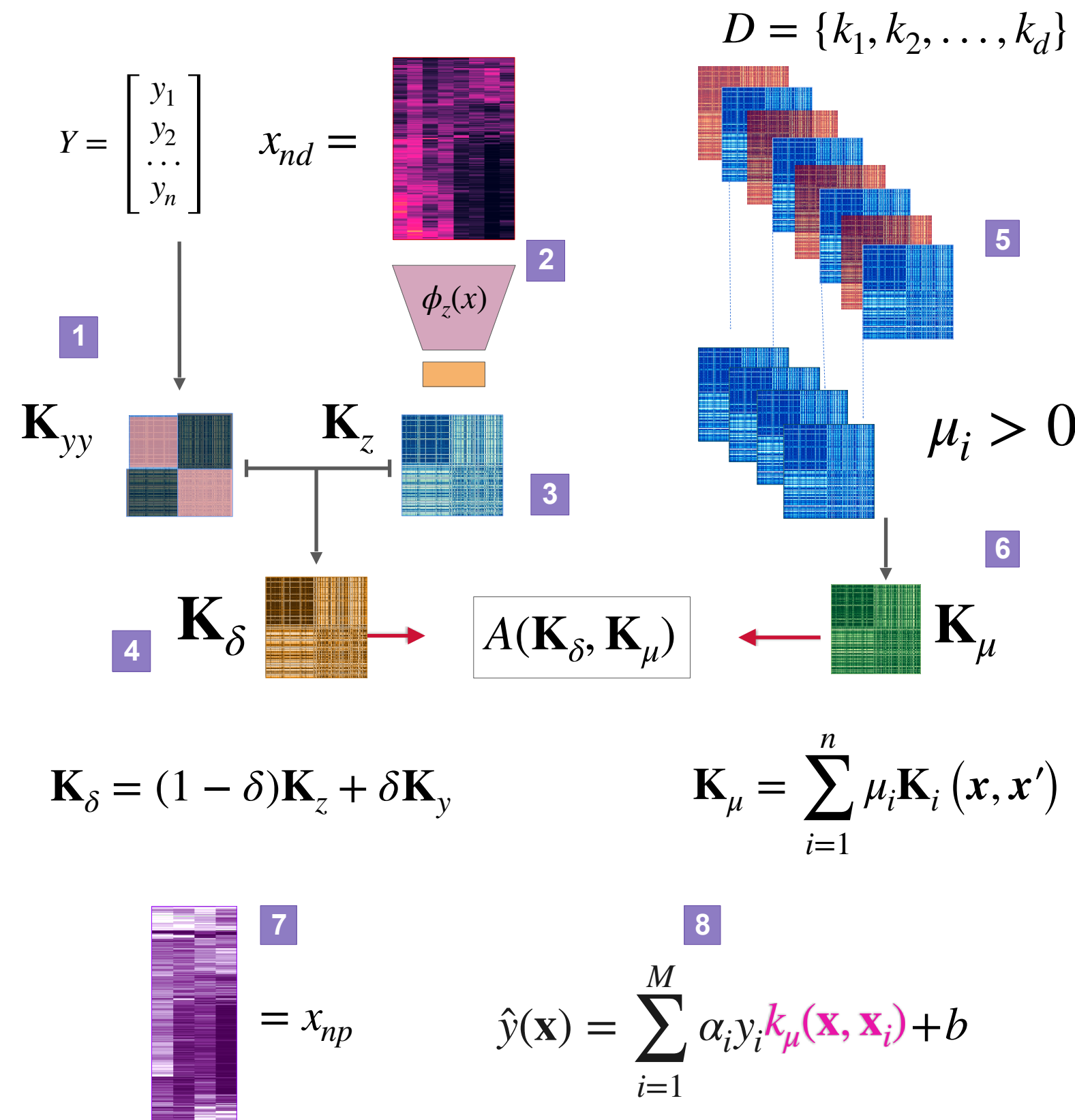
To learn from \mathbf{z} and \mathbf{y} kernel alignment is used. This score measures the similarity between a target kernel and a sample kernel. The higher the score the more similar both kernels.

$$A(K_1, K_2) = \frac{\langle K_1, K_2 \rangle_F}{\sqrt{\langle K_1, K_1 \rangle_F \langle K_2, K_2 \rangle_F}}$$

To select features multiple kernel learning is used. A linear combination of feature-wise kernels is done to obtain a new custom kernel that maximizes the alignment with the target kernel.

$$k_\mu(x, x') = \sum_{i=1}^n \mu_i k_i(x, x'), \mu_i \geq 0$$

Proposed method: KLR-FS (arxiv.org/abs/2004.04866)



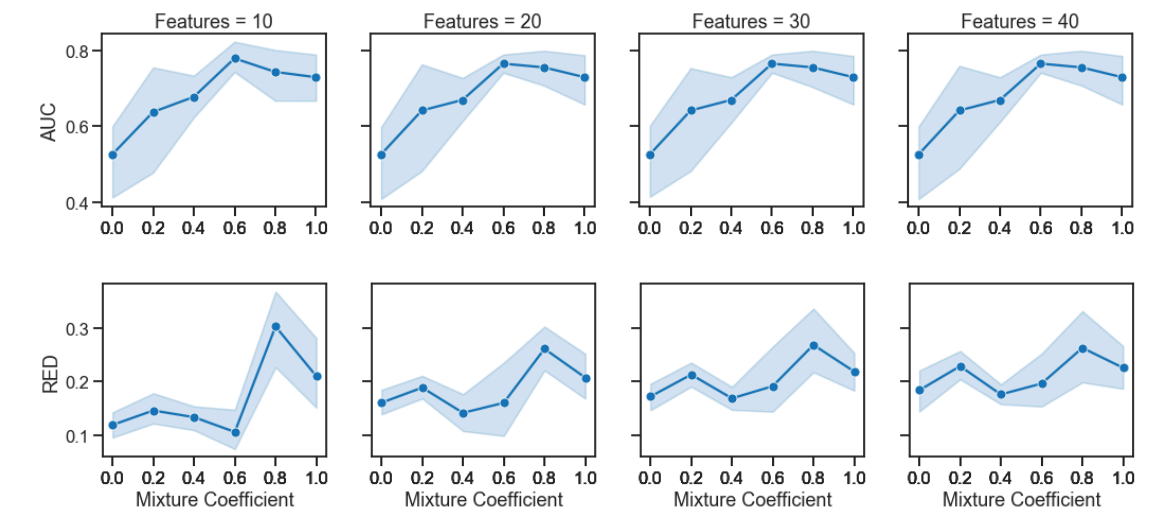
$$K_\delta = (1 - \delta)K_z + \delta K_y$$

$$K_\mu = \sum_{i=1}^n \mu_i K_i(x, x')$$

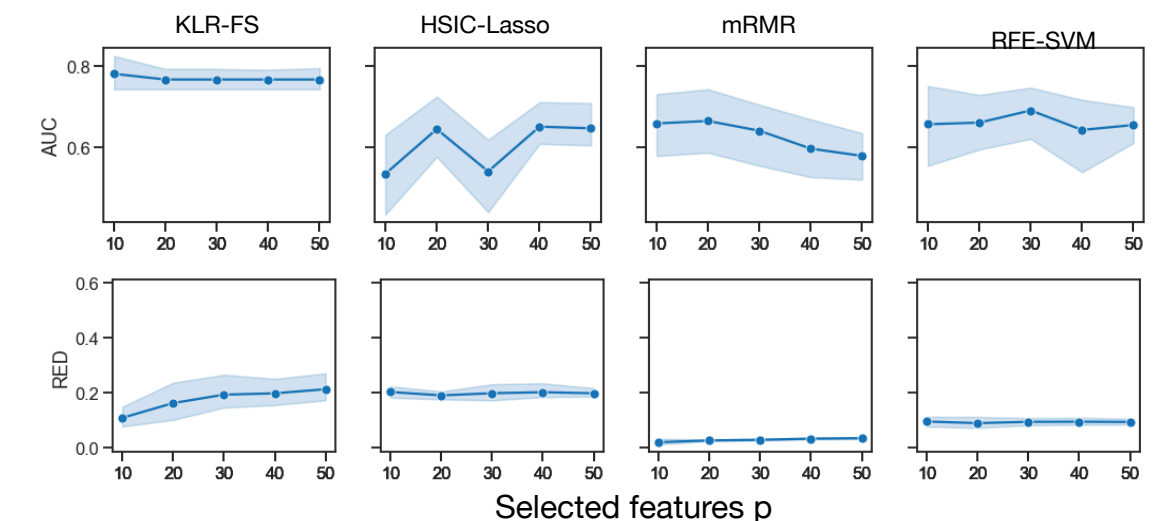
- (1) A supervised K_y kernel is built with tumor labels.
- (2) A latent space from training data is obtained via a non-linear dimensionality reduction method (kPCA).
- (3) An unsupervised kernel K_z is built on the latent space.
- (4) Both kernels are mixed via a mixture coefficient Δ to obtain a target kernel that keeps supervised and unsupervised information.
- (5) A set of feature-wise kernels are built and linearly combined via multiple kernel learning (MKL).
- (6) The target kernel used in the MKL step is the delta kernel, then only the feature-wise kernels that approximate to the target kernel are kept and the rest discarded.
- (7) With the feature-weights higher than zero features are selected.
- (8) The kernel obtained from MKL is used in Support Vector Classification.

Results

The method is evaluated with the Breast Cancer patient survival data from the BRCA-US repository at the International Cancer Genome Consortium. Area Under the ROC Curve is used to evaluate tumor classification. The Redundancy score (RED) is used to evaluate how redundant are the selected features, the lower the RED the better the selection. The train-test split is done five times randomly and selection and classification results are averaged.



KLR-FS is evaluated for different values of selected features and different mixtures (delta parameter) of supervised and unsupervised target kernels. At $\delta = 0.6$ there is a peak of AUC and simultaneously a low redundancy, suggesting that learning partially from latent structure improves the feature selection and the classification on test samples.



KLR-FS with $\delta=0.6$ is compared with HSIC-Lasso, mRMR and RFE-SVM in support vector classification. The highest AUC for different values of p is achieved by the proposed method. Despite the fact that the RED of the selected features by KLR-FS is not the smallest, it is at similar levels as the benchmark methods.

Conclusions

Learning from both latent structure and tumor labels improves the classification performance and serves as regularization.