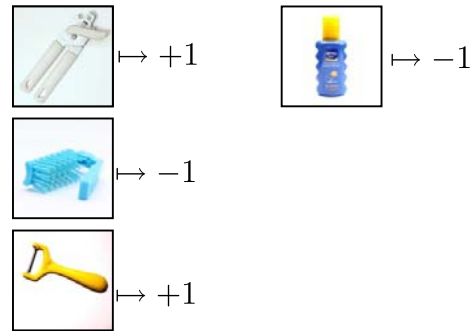


Kernel Methods for Structured In- & Outputs

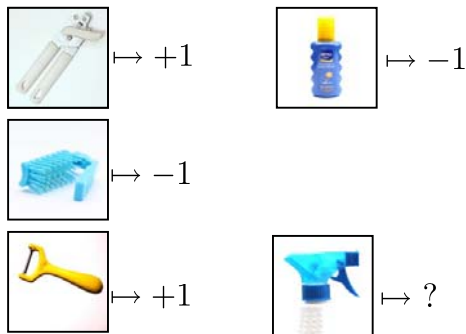
Thomas Gärtner

Supervised Classification Example



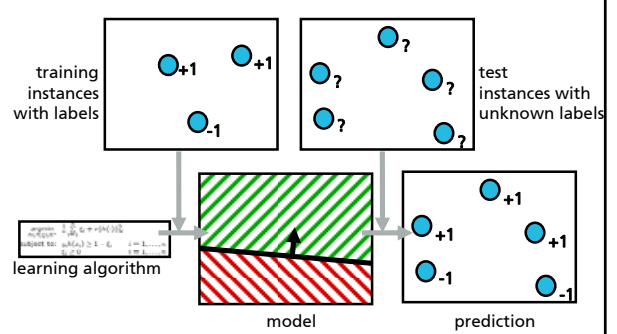
Seite 2

Supervised Classification Example



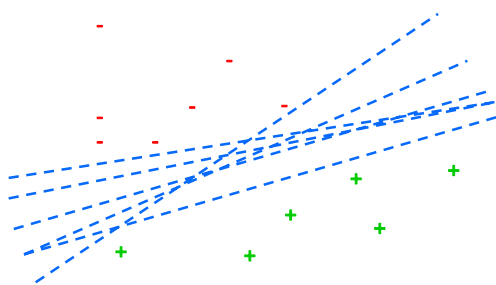
Seite 3

Supervised Classification



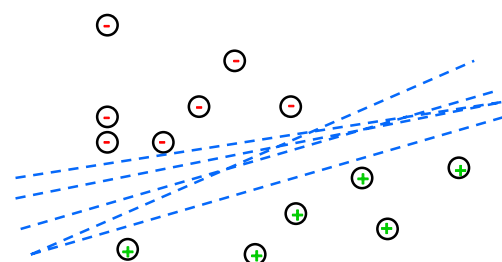
Seite 4

Classification by Linear Functions



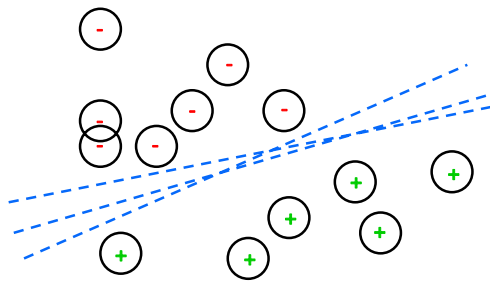
Seite 5

Motivation: Support Vector Machines



Seite 6

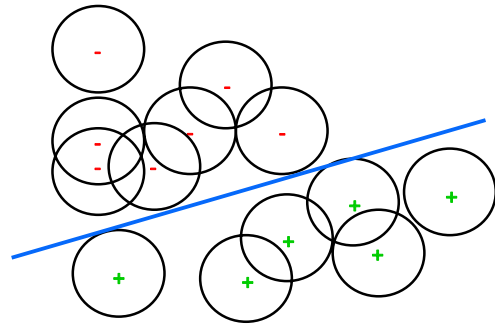
Motivation: Support Vector Machines



Seite 7

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

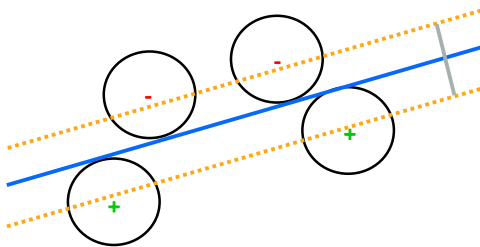
Motivation: Support Vector Machines



Seite 8

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

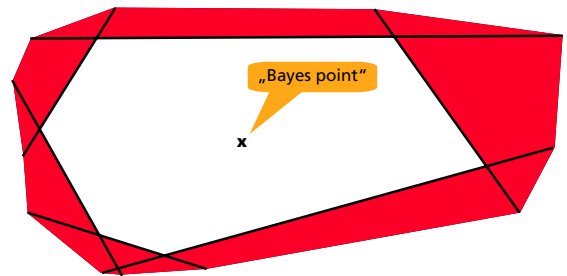
Motivation: Maximal Margin Hyperplane



Seite 9

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

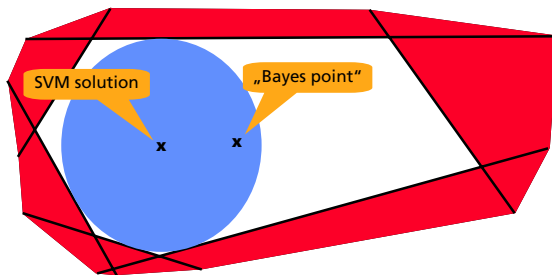
Version Space Geometry of SVMs



Seite 10

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

Version Space Geometry of SVMs



Seite 11

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

Regularised Risk Minimisation Principle

$$\hat{h} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) \quad \text{s.t.} \quad \Omega[f] \leq \lambda$$

error on training data (drawn iid)

capacity of hypothesis space

$$\hat{h} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) + \nu \Omega[f]$$

Seite 12

© Fraunhofer Institut für intelligente Analyse- und Informationssysteme IAS

Regularised Risk Minimisation Principle

$$\hat{h} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) \quad \text{s.t.} \quad \Omega[f] \leq \lambda$$

error on training data (drawn iid)

capacity of hypothesis space

$$\hat{h} = \operatorname{argmin}_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(x_i)) + \nu \Omega[f]$$

kernel methods have commonly F: RKHS with pd kernel k ; $\Omega[\cdot] = \|\cdot\|^2$; V convex

-> resulting optimisation problem is convex

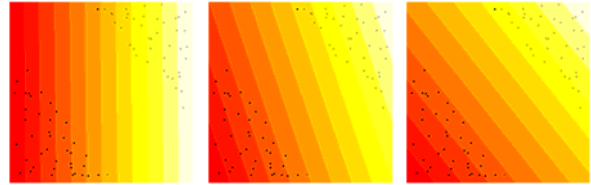
-> minimiser lies in the span of $\{k(x_i, \cdot)\}_i$

Seite 13

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Stability with respect to Outliers



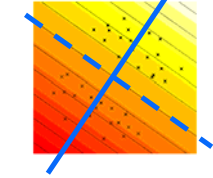
Seite 14

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

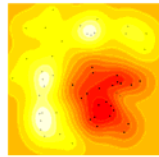


Kernel Methods

in hypothesis space: $f(\cdot) = \sum c_i k(x_i, \cdot)$



in input space:



before, we were looking for a hyperplane in input space, now, we are looking for a hyperplane in a function space (a RKHS with kernel k) where each x is represented by the function $k(x, \cdot)$

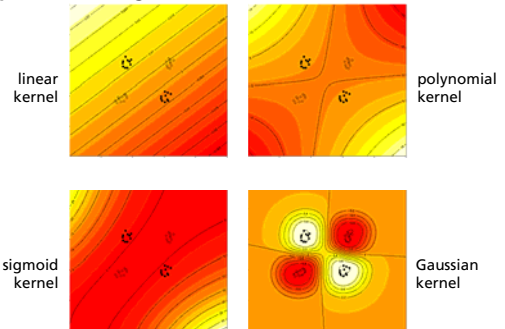
requiring an RKHS means basically everything works as before

Seite 15

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Hypothesis Change with Kernels

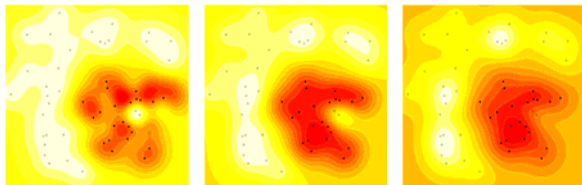


Seite 16

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Stability with respect to Outliers

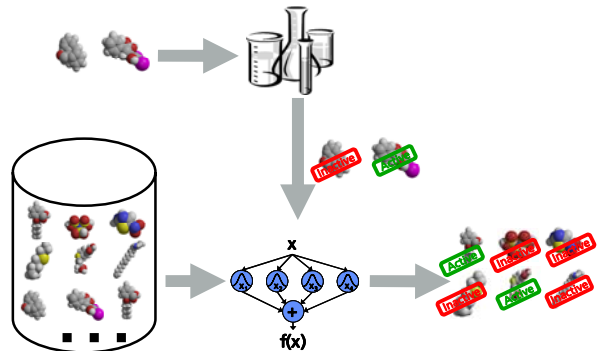


Seite 17

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Prediction of Biochemical Properties



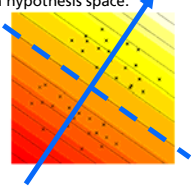
Seite 18

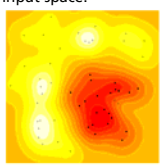
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

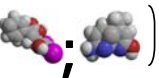


Approach: Kernel Methods for Graphs

$f(\cdot) = \sum c_i k(x_i, \cdot)$

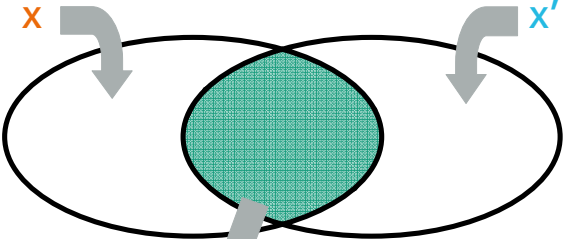
in hypothesis space: 

in input space: 

we need valid kernel functions for molecules: $k(\text{molecule}_1; \text{molecule}_2)$ 

Seite 19
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

The "Intersection Kernel"-Principle (1)

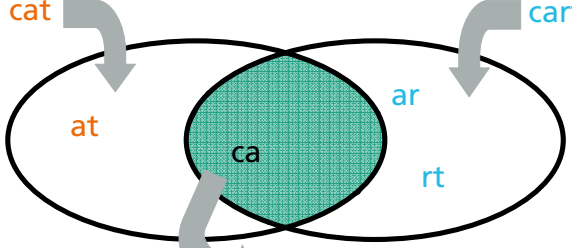


can be feasible even if x is only given intensionally and even if $|x|$ is exponential in its representation!!!

$k(x, x') = \mu(x \cap x')$

Seite 21
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

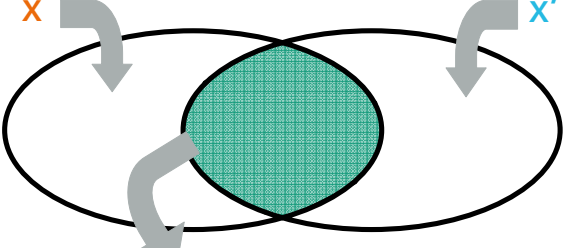
n-gram Kernels (n = 2)



$k(\text{cat}, \text{cart}) = 1$

Seite 22
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

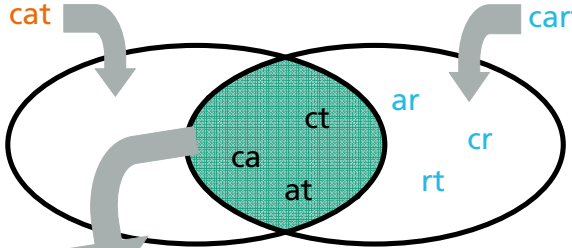
The "Intersection Kernel"-Principle (2)



$k(x, x') = \sum_{p \in x \cap x'} f(|p|_x) \cdot f(|p|_{x'})$

Seite 23
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

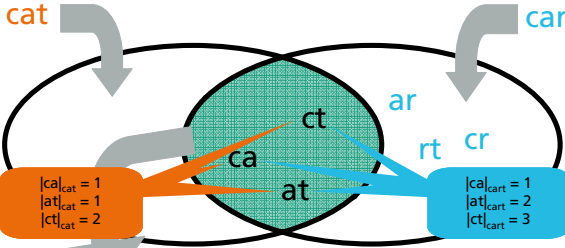
Subsequence Kernels (n = 2)



$k(\text{cat}, \text{cart}) = \lambda^2 + \lambda^3 + \lambda^5$

Seite 24
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

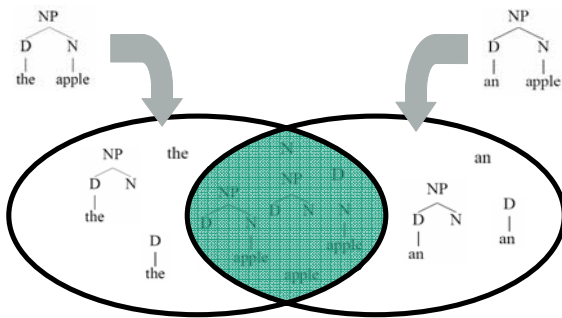
Subsequence Kernels (n = 2)



$k(\text{cat}, \text{cart}) = \lambda^2 + \lambda^3 + \lambda^5$

Seite 25
© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAS

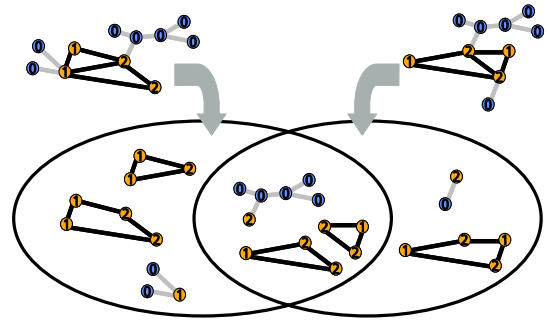
(Parse-) Tree Kernel



Seite 27

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

Graph Kernels

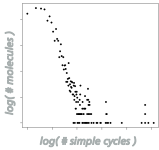


Seite 28

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

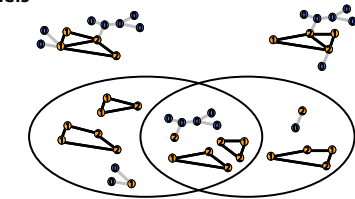
Cyclic Pattern Kernels

$k(\text{graph}_1, \text{graph}_2)$



computing CPK is NP-hard

but efficient learning is still possible as most compounds contain only a few cycles



'complete' as well as all-subgraphs graph kernels are hard polynomial algorithm for all-walks kernels [colt'03]

(enumeration of) cyclic patterns [kdd'04]

walk kernel for RRL [mlj'06]

[ilp'03] "best algorithmic paper"

Seite 29

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

Graph Kernels --- Negative Results

Computing the "all patterns" graph kernel is hard for instance for patterns that are "all paths", "all cycles", or "all connected subgraphs" with subgraph-isomorphism as the embedding operator.

Computing any graph kernel for which

$$k(G, \cdot) = k(G', \cdot) \Leftrightarrow G \simeq G'$$

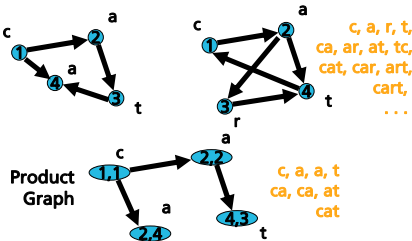
is at least as hard as deciding graph isomorphism

Seite 30

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

All-Walks Kernels

c, a, a, t
ca, ca, at, ta
cat
cata



$$k_{\times}(G_1, G_2) = 1^T \left[\lim_{n \rightarrow \infty} \sum_{i=0}^n \lambda_i E_{\times}^i \right] 1_{\text{eg}} = 1^T (I - \gamma E_{\times})^{-1} 1$$

Seite 31

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

Limits of matrix power series

eigenvalue decomposition

$$E = TDT^{-1} \Rightarrow \lambda_i(E^n) = (\lambda_i(E))^n$$

exponential series

$$\sum_{n=0}^{\infty} \frac{\beta^n}{n!} E^n \Rightarrow \lambda_i \left(\sum_{n=0}^{\infty} \frac{\beta^n}{n!} E^n \right) = e^{\beta \lambda_i(E)}$$

geometric series

$$\sum_{n=0}^{\infty} \gamma^n E^n \Rightarrow \lambda_i \left(\sum_{n=0}^{\infty} \gamma^n E^n \right) = \frac{1}{1 - \gamma \lambda_i(E)}$$

Seite 32

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS

Kernels between Vertices in (Hyper-) Graphs

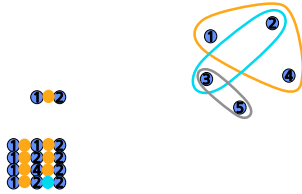
$$v_1^\top v_2 = 0$$

$$v_1^\top B B^\top v_2 = 1$$

$$v_1^\top B B^\top B B^\top v_2 = 4$$

$$\left[\sum_{n=0}^{\infty} \lambda_n (B B^\top)^n \right]_{ij}$$

weighted sum over walks from vertex i to j in the hypergraph.

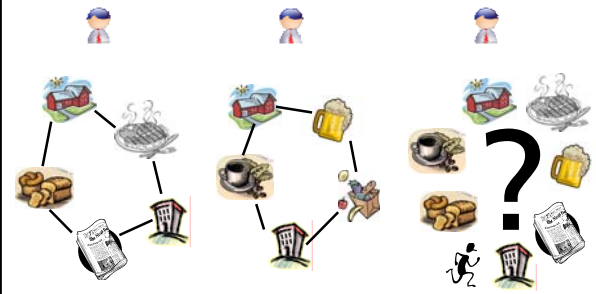


Seite 33

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Travelling John Q Public Problem



Seite 34

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Structured Output Prediction

input & output spaces $\mathcal{X}, \mathcal{Y}(\Sigma)$

training data $D = \{(x_i, Y_i)\}_{i \in [n]} \subseteq \mathcal{X} \times 2^{\mathcal{Y}(\Sigma)}$ (often $|Y_i| = 1$)

Ansatz

- joint scoring function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$
- prediction / decoding** $g(x) = \operatorname{argmax}_{y \in \mathcal{Y}} f(x, y)$

Seite 35

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Structured Output Prediction

input & output spaces $\mathcal{X}, \mathcal{Y}(\Sigma)$

training data $D = \{(x_i, Y_i)\}_{i \in [n]} \subseteq \mathcal{X} \times 2^{\mathcal{Y}(\Sigma)}$ (often $|Y_i| = 1$)

Ansatz

- joint scoring function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$
- prediction / decoding** $g(x) = \operatorname{argmax}_{y \in \mathcal{Y}} f(x, y)$

typically, the output space is exponential in $|D| + |\Sigma|$

Seite 36

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Finding the Scoring Function (simplified)

$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \Omega[f]$$

subject to $f(x_i, y) > f(x_i, z) \quad (\forall i, \forall z \in \mathcal{Y} \setminus Y_i, y \in Y_i)$

has exponentially many constraints !!

Seite 37

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Structured Perceptron

make as few as possible mispredictions on a sequence $(x_1, y_1), (x_2, y_2), \dots$

```

initialise  $f_0 \leftarrow 0$ 
for  $x_i$  predict  $\operatorname{argmax}_z f_{i-1}(x, z)$ 
if misprediction
    then let  $f_i(\dots) \leftarrow f_{i-1}(\dots) + k((x_i, y_i), (\dots))$ 
    else let  $f_i(\dots) \leftarrow f_{i-1}(\dots)$ 
    
```

Seite 38

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAI5



Structured Perceptron

make as few as possible mispredictions on a sequence $(x_1, y_1), (x_2, y_2), \dots$

```

initialise  $f_0 \leftarrow 0$ 
for  $x_i$  predict  $\operatorname{argmax}_z f_{i-1}(x, z)$ 
if misprediction
  then let  $f_i(\dots) \leftarrow f_{i-1}(\dots) + k((x_i, y_i), (\dots))$ 
  else let  $f_i(\dots) \leftarrow f_{i-1}(\dots)$ 

```

Novikoff's Theorem applies and allows to bound the number of mispredictions as $O(R^2 / \delta^2)$

Seite 39

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



SVM-like Structured Output with Decoding-Oracle

$$\begin{aligned} \hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \quad & \nu \|f\|^2 + \|\xi_i\|_1 \\ \text{subject to} \quad & f(x_i, y) > \Delta(y, z) + f(x_i, z) - \xi_i \quad (\forall i, \forall z \in \mathcal{Y} \setminus \{y_i\}) \end{aligned}$$

can be solved efficiently by cutting-plane techniques

$$\begin{aligned} \hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \quad & \|w\|_1 + \|\xi_i\|_1 \\ \text{subject to} \quad & \langle w, \phi(x_i, y_i) \rangle > \Delta(y, z) + \langle w, \phi(x_i, z) \rangle - \xi_i \quad (\forall i, \forall z \in \mathcal{Y} \setminus \{y_i\}) \end{aligned}$$

can be solved efficiently by the ellipsoid method

Seite 41

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Finding Violated Constraints

decoding (polynomial time?)

- given f, x find $z \in \operatorname{argmax}_{y \in \mathcal{Y}} f(x, y)$

separation (polynomial time?)

- given f, x, y find any $z : f(x, y) < f(x, z)$ or proof that none exists

optimality (polynomial time? / in NP?)

- given f, x, y decide $\nexists z \in \mathcal{Y} : f(x, y) < f(x, z)$

"decoding is in P" is the strongest assumption
while "optimality is in NP" is the weakest

Seite 42

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



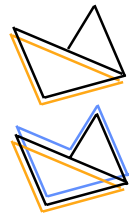
Optimality vs Non-Optimality

optimality (in NP?) is there **no** longer cycle

- given f, x, y decide $\nexists z \in \mathcal{Y} : f(x, y) < f(x, z)$
- what is a short certificate of optimality?

sub-optimality (in NP!) is there **any** longer cycle

- given f, x, y decide $\exists z \in \mathcal{Y} : f(x, y) < f(x, z)$
- certificate of non-optimality is short



we are interested mostly in problems where
sub-optimality is NP-complete

Seite 43

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



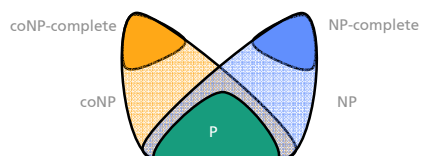
Optimality vs Non-Optimality

optimality (coNP-complete and hence not in NP)

- given f, x, y decide $\nexists z \in \mathcal{Y} : f(x, y) < f(x, z)$

sub-optimality (NP-complete)

- given f, x, y decide $\exists z \in \mathcal{Y} : f(x, y) < f(x, z)$



Seite 44

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Structured Output Ranking

$$h^* = \operatorname{argmin}_{f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}} \quad \text{consider regularised AUC of all ranked outputs}$$

- representer theorem $h^* \in \operatorname{span}$ exponentially many coefficients

Seite 45

© Fraunhofer Institut für Intelligente Analyse- und Informationssysteme IAIIS



Structured Output Ranking

$h^* = \operatorname{argmin}_{f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}}$ consider **regularised AUC** of all ranked outputs

- representer theorem $h^* \in \operatorname{span}$ **exponentially** many coefficients

assume a **tensor product of Hilbert spaces** $h \in \mathcal{H} = \mathcal{H}_X \otimes \mathcal{H}_Y$
assume a small output basis (d)

- **factorised representer theorem**

$h^* \in \operatorname{span}$ **polynomially** many coefficients

minimise a quadratic upper bound on the AUC

Seite 46

© Fraunhofer Institut für Intelligente
Analyse- und Informationssysteme IAIIS



Structured Output Ranking

the following terms occur in the objective

$$\sum_{y \in \mathcal{Y}} f(x_i, y) = f_\alpha^i{}^\top \sum_{z \in \mathcal{Y}} \phi(z) = f_\alpha^i{}^\top \Phi$$

$$\sum_{y \in \mathcal{Y}} f^2(x_i, y) = f_\alpha^i{}^\top \sum_{z \in \mathcal{Y}} \phi(z) \phi^\top(z) f_\alpha^i = f_\alpha^i{}^\top C f_\alpha^i$$

Seite 47

© Fraunhofer Institut für Intelligente
Analyse- und Informationssysteme IAIIS



Structured Output Ranking

the following terms occur in the objective

$$\sum_{y \in \mathcal{Y}} f(x_i, y) = f_\alpha^i{}^\top \sum_{z \in \mathcal{Y}(\Sigma)} \phi(z) = f_\alpha^i{}^\top \Phi$$

$$\sum_{y \in \mathcal{Y}} f^2(x_i, y) = f_\alpha^i{}^\top \sum_{z \in \mathcal{Y}(\Sigma)} \phi(z) \phi^\top(z) f_\alpha^i = f_\alpha^i{}^\top C f_\alpha^i$$

often these **exponential sums** can be computed in **polynomial time**

Seite 48

© Fraunhofer Institut für Intelligente
Analyse- und Informationssysteme IAIIS

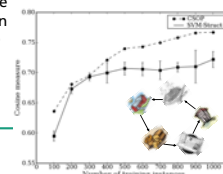
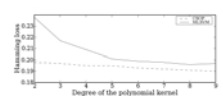
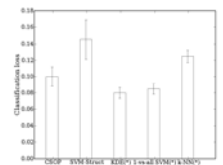


Summary

most approaches assume (at least) that **optimality** is in NP. However, it is often **co-NP complete**.

we proposed a **ranking** based approach that results in an **unconstrained convex optimisation problem**.

The new approach complements the existing one and is based on an orthogonal assumption. It can be adapted to new output spaces by coding one new matrix.

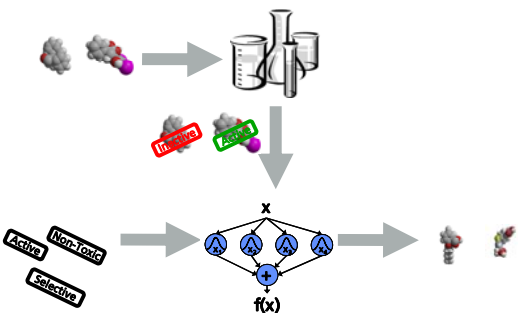


Seite 49

© Fraunhofer Institut für Intelligente
Analyse- und Informationssysteme IAIIS



Constructive Machine Learning for De Novo Drug Design



Seite 50

© Fraunhofer Institut für Intelligente
Analyse- und Informationssysteme IAIIS

