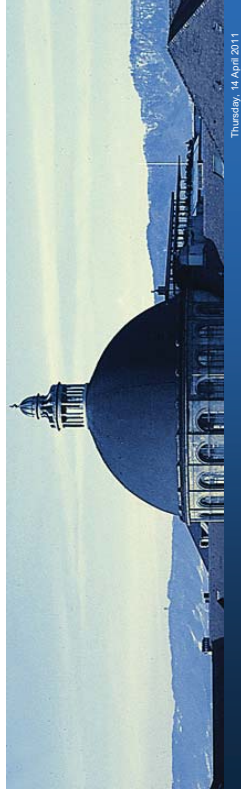


## Structure learning and the generalization capacity of algorithms

Joachim M. Buhmann

Computer Science Department, ETH Zurich



Thursday, 14 April 2011

Unsolved Problems in Computer Science:

- 1) What is useful information?
- 2) How can we design robust algorithms?

*myVision*: Information theory for pattern analysis

- Approximation capacity of algorithms
- Examples for cost function minimization
  - Cluster validation
  - Role based access control
  - (Robust SVD)
- Philosophical remarks

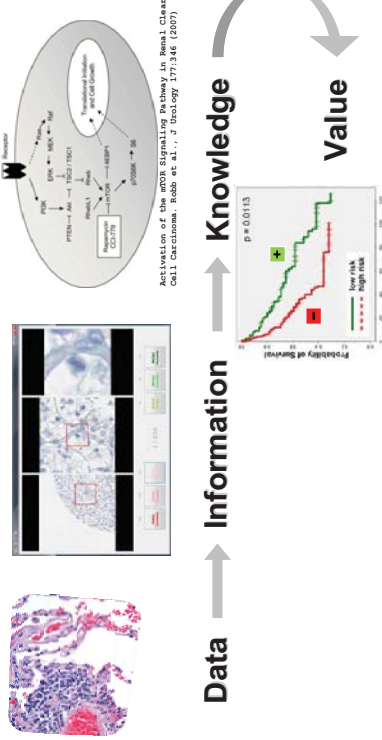
Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

4

## Modern information society and the Information Technology value chain



Thursday, April 14, 2011

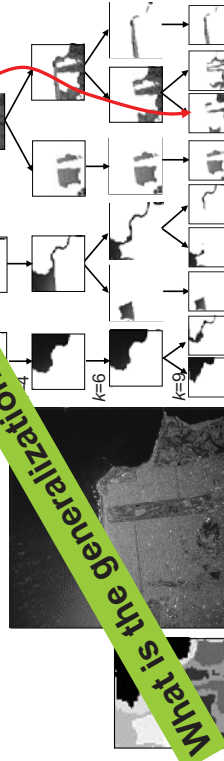
Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

3

## Pattern recognition and modeling

- Given are data  $X$  and hypotheses, i.e. interpretation of these data.
- An algorithm finds good hypotheses



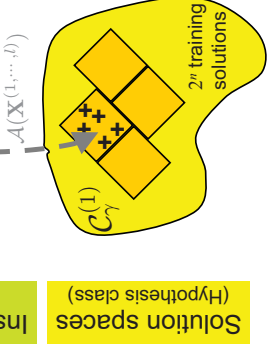
Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

5

## Generalization of algorithms



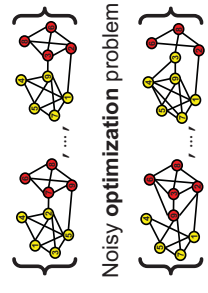
- Sample instance from a trial distribution (smoothed analysis)
- Map instances to solutions with algorithm  $\mathcal{A}(\mathbf{X}^{(1, \dots, l)})$
- Cover hypothesis class
- Test with  $\mathbf{X}^{(2)}$

## Information Theory & Pattern Recognition

- IT-Components
  - Pattern Recognition elements
- Code book  $\zeta$  {strings} = hypothesis class
 

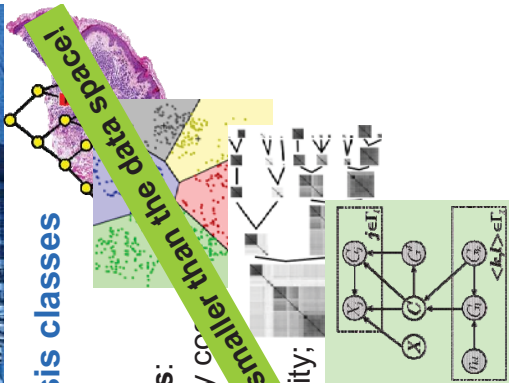
0000000	0100101	1000011	1100110
0001111	0101010	1001100	1101001
0010110	0110011	1010101	1110000
0011001	0111100	1011010	1111111

  - Noisy channel
    - 1100110  $\rightarrow$  0110110
  - Decoder: minimize Hamming distance
    - 0110110 - 1100110 | <
    - 0110110 - 1010101 |
- Noisy optimization problem
  - Decoding by **approximate optimization** of test instance



## Examples of hypothesis classes

- Classifiers
- Partitions or clusterings: compactness/connectivity
- Trees or dendrograms: partitions with ultrametricity; Tree depth
- Graph models: measured probability models; # nodes/edges?

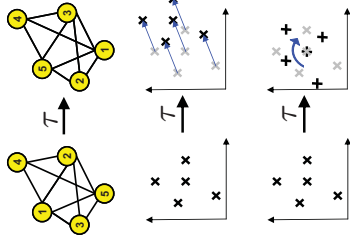


## Symmetries of the Learning Problem

- Assume (!) that the algorithm  $\mathcal{A}$  is equivariant under the transformations
 
$$\mathcal{T} = \{\tau : \tau \circ \mathcal{A}(\mathbf{X}) = \mathcal{A}(\tau \mathbf{X} \circ \mathbf{X})\}$$
- Idea: define a code by transforming a given problem  $\Rightarrow$  codebook  $\mathbb{T} = \{\tau_i \in \mathcal{T} : 1 \leq i \leq 2^{n\rho}\}$
- What is an appropriate set of transformations? That depends on the algorithm and the hypothesis class!

## How to generate code problems?

1. Combinatorial optimization problems: **permutation** of e.g. vertices in graphs
2. Localization problems: **shifts** of data
3. Orientation problems (PCA, SVD): **rotations & projection**
4. Sparse linear regression: **Scaling & permutation**



## Approximate Optimization

- Define weights for hypotheses with low costs

$$w_\beta(c, \mathbf{X}) = \exp(-\beta R(c, \mathbf{X})) \quad c \in \mathcal{C}(\mathbf{X})$$

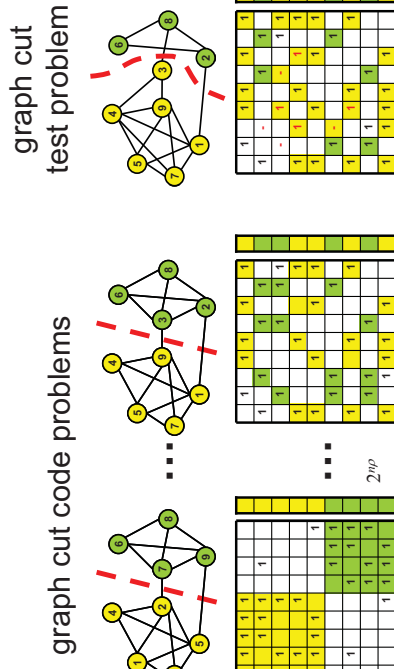
- Special case: **Approximation set**

$$w_\beta(c, \mathbf{X}) = \begin{cases} 1, & R(c, \mathbf{X}) \leq R(c^\perp, \mathbf{X}) + \gamma; \\ 0, & \text{otherwise.} \end{cases}$$

- Total weight of low cost hypotheses

$$W_\beta(\mathbf{X}) = \{w_\beta(c, \mathbf{X}) : c \in \mathcal{C}(\mathbf{X})\}$$

## Code problem generation for Graph Cut



## Coding with Graph Cut approximation sets

define a set of code problems

problem generator PG

$$R(\cdot, \mathbf{X}^{(1)})$$

$$R(\cdot, \mathbf{X}^{(1)})$$

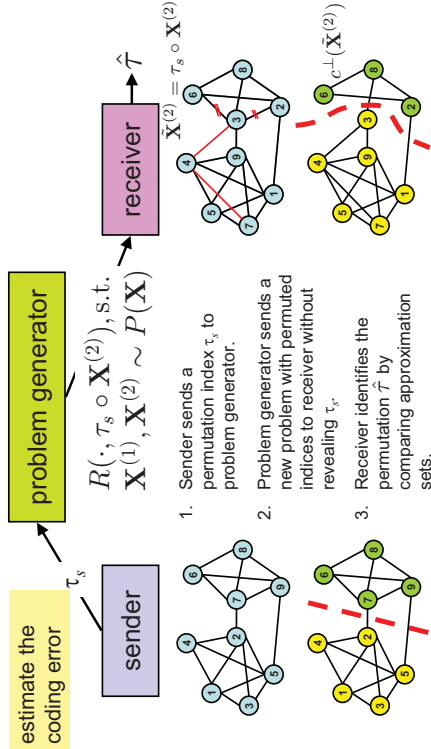
receiver

$$\{\tau_1, \dots, \tau_{2^{np}}\}$$

$$\tau \circ \mathbf{X}^{(1)}$$



## Communication by approximation sets



## Condition of vanishing total error

- $\lim_{n \rightarrow \infty} P(\hat{\tau} \neq \tau_s | \tau_s) = 0$  yields
    - Rate is bounded by mutual information
- $$\rho \log 2 < \frac{1}{n} \log \frac{|\{\tau \neq s\}| Z_{\beta}^{(1)} Z_{\beta}^{(2)}}{Z_{\beta}^{(1)} Z_{\beta}^{(2)}}$$
- $$= \frac{1}{n} \left( \log \frac{|\{\tau \neq s\}|}{Z_{\beta}^{(1)}} + \log \frac{|C^{(2)}|}{Z_{\beta}^{(2)}} - \log \frac{|C^{(2)}|}{Z_{\beta, s}^{(1 \& 2)}} \right)$$
- $$\equiv \mathcal{I}_{\gamma}(\tau_s, \hat{\tau})$$
- Lower bound: generalize Fano's inequality to ASC (work in progress)

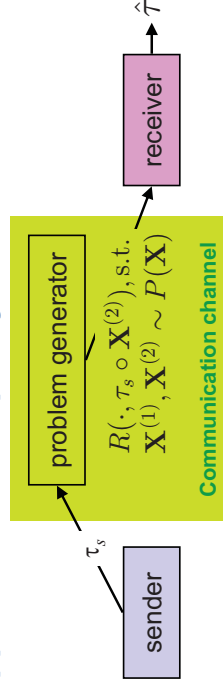
## Communication Process

- Receiver compares sets of hypothesis weights  $W_{\beta}(\mathbf{X}^{(1)})$  of training instance (code problem) with approximate clusterings  $W_{\beta}(\mathbf{X}^{(2)})$  of the test data.
- Define a mapping  $\psi : \mathbf{X}^{(1)} \rightarrow \mathbf{X}^{(2)}$
- Decoding by maximizing weight overlap

$$\hat{\tau} \in \arg \max_{\tau \in \Pi} \sum_{c \in C(\mathbf{X}^{(1)})} w_{\beta}(c, \tau(\mathbf{X}^{(1)})) w_{\beta}(c, \tilde{\mathbf{X}}^{(2)})$$

with  $\tilde{\mathbf{X}}^{(2)} := \psi \circ \tau_s \circ \psi^{-1}(\mathbf{X}^{(2)})$

## Model Selection by Maximization of Approximation Capacity

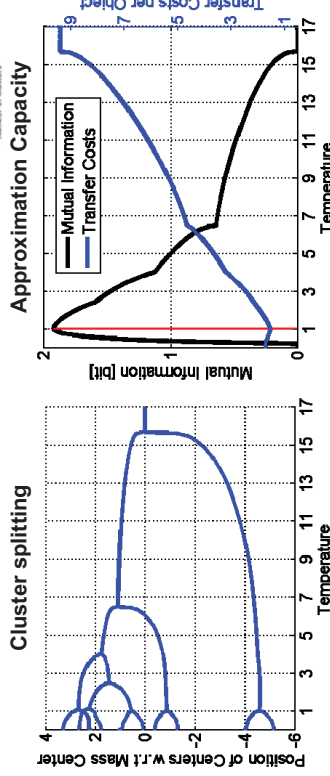


- Optimize the communication channel w.r.t. approximation quality  $\gamma(\beta)$ , topology and metric of solution space, cost function  $R(\cdot, \cdot)$ , transfer function  $\psi$

## 2d Mixture Model Estimation

### Experimental Setting:

5 Gaussians,  $n=10000$ ,  $d=2$ ,  $k^{\max}=10$



Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

21

## Conclusion

- **Quantization:** Noise quantizes mathematical structures (hypothesis classes)  $\Rightarrow$  symbols
  - These symbols can be used for **coding!**
  - Optimal error free coding scheme determines **approximation capacity** of a model class.
- $\Rightarrow$  Bounds for robust optimization.
- $\Rightarrow$  **Quantization** of hypothesis class measures **structure specific information** in data.

Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

26

## Future Work

- **Generalization:** replace approximation sets based on cost functions by smoothed outputs of **algorithms** ("smoothed generalization")
- **Model reduction** in dynamical systems: quantize sets of ODEs or PDEs (systems biology)
- Relate **statistical complexity**, i.e. the approximation capacity, to algorithmic or **computational complexity**.

Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

27

## Philosophical speculations

- We experience a **paradigm shift from model driven reasoning to algorithm dominated reasoning** (Bernard Chazelle "The Algorithm: Idiom of Modern Science")
- $\Rightarrow$  **model validation** more essential than modeling since modeling can be algorithmically formulated as exploration of model space.
- *Ceterum censeo:* The coupling of **statistical complexity** and **algorithmic complexity** should be reconsidered in the light of **statistical learning theory** and information theory.

Thursday, April 14, 2011

Joachim M. Buhmann

UPPR 2011 Workshop, Heidelberg

28