

# Computational Properties of Probabilistic Neural Networks

*Jiří Grim and Jan Hora*

**Institute of Information Theory and Automation  
Academy of Sciences of the Czech Republic, Prague**

**Department of Pattern Recognition**

<http://www.utia.cas.cz/RO>

**ICANN'10, Thessaloniki, September 15-18, 2010**

# Outline

- 1 Overfitting in Neural Networks
- 2 Probabilistic Neural Networks
  - Statistical Pattern Recognition Based on Mixtures
  - Structural Mixture Model
- 3 Numerical Experiment: Classification of Chess-Board Patterns
  - Randomly Generated Chess-Board Patterns
  - Marginal Probabilities of the Chess-Board Patterns
  - Recognition of Chess-board patterns
  - Recognition of Chess-board patterns
- 4 Concluding Remarks

# Overfitting in Neural Networks

## problem of overfitting

- small multidimensional training data sets ( $\approx$  insufficiently representative)
- $\Rightarrow$  "overfitting" of parameters to training data
- $\Rightarrow$  bad "generalizing property" caused by overfitting
- a general analysis of overfitting is difficult

## to reduce the risk of overfitting:

- dimensionality reduction and/or large data sets
- cross-validation techniques
- "under-computing": stopping rule for training
- optimal complexity of classifiers

## Probabilistic Neural Networks:

structural mixtures  $\Rightarrow$  reduced complexity  $\Rightarrow$  less prone to overfitting

# Statistical Pattern Recognition Based on Mixtures

$\mathbf{x} = (x_1, \dots, x_N) \in \mathcal{X}$ : N-dimensional data vectors

$\Omega = \{\omega_1, \omega_2, \dots, \omega_J\}$ : finite number of classes

$P(\mathbf{x}|\omega)p(\omega)$ ,  $\omega \in \Omega$ : conditional distributions of classes

**approximation of  $P(\mathbf{x}|\omega)$  by mixtures of product components:**

$$P(\mathbf{x}|\omega) = \sum_{m \in \mathcal{M}_\omega} f(m)F(\mathbf{x}|m), \quad P(\mathbf{x}) = \sum_{\omega \in \Omega} p(\omega)P(\mathbf{x}|\omega)$$

$$F(\mathbf{x}|m) = \prod_{n \in \mathcal{N}} \theta_{mn}^{x_n} (1 - \theta_{mn})^{1-x_n}, \quad 0 \leq \theta_{mn} \leq 1, \quad m \in \mathcal{M}_\omega, \quad \mathcal{M} = \sum_{\omega \in \Omega} \mathcal{M}_\omega$$

**decision making based on Bayes formula:**

$$p(\omega|\mathbf{x}) = \frac{P(\mathbf{x}|\omega)p(\omega)}{P(\mathbf{x})} = \sum_{m \in \mathcal{M}_\omega} q(m|\mathbf{x}), \quad q(m|\mathbf{x}) = \frac{w_m F(\mathbf{x}|m)}{\sum_{j \in \mathcal{M}} w_j F(\mathbf{x}|j)}$$

**probabilistic neuron  $\approx$  mixture component**

# Structural Mixture Model (Grim et al. 1986, 1999, 2002)

$$F(\mathbf{x}|m) = \prod_{n \in \mathcal{N}} f_n(x_n|m)^{\phi_{mn}} f_n(x_n|0)^{1-\phi_{mn}}, \quad \phi_{mn} \in \{0, 1\} \approx \text{structure}$$

$$f_n(x_n|m) = \theta_{mn}^{x_n} (1 - \theta_{mn})^{1-x_n}, \quad n \in \mathcal{N}, \quad \mathcal{N} = \{1, \dots, N\}$$

$\phi_{mn} = 0 \Rightarrow f_n(x_n|m)$  is replaced by fixed “background”  $f_n(x_n|0)$

$$P(\mathbf{x}|\omega) = \sum_{m \in \mathcal{M}_\omega} F(\mathbf{x}|m) f(m) = F(\mathbf{x}|0) \sum_{m \in \mathcal{M}_\omega} G(\mathbf{x}|m, \phi_m) f(m)$$

$$F(\mathbf{x}|0) = \prod_{n \in \mathcal{N}} f_n(x_n|0), \quad G(\mathbf{x}|m, \phi_m) = \prod_{n \in \mathcal{N}} \left[ \frac{f_n(x_n|m)}{f_n(x_n|0)} \right]^{\phi_{mn}}$$

$G(\mathbf{x}|m, \phi_m) \approx$  defined on different subspaces

▶ OPTIMIZATION: EM algorithm

$$p(\omega|\mathbf{x}) = \frac{P(\mathbf{x}|\omega)p(\omega)}{P(\mathbf{x})} = \frac{\sum_{m \in \mathcal{M}_\omega} w_m G(\mathbf{x}|m, \phi_m)}{\sum_{j \in \mathcal{M}} w_j G(\mathbf{x}|j, \phi_j)}$$

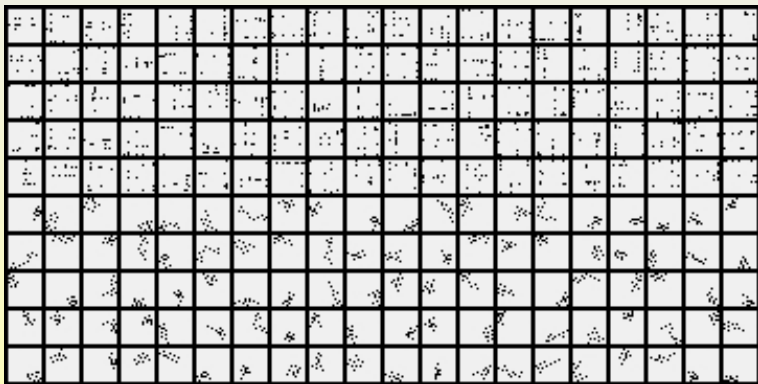
$\Rightarrow$  “background distribution”  $F(\mathbf{x}|0)$  cancels in the Bayes formula

# Chess-Board Patterns Made by Rook and Knight

**16x16 chess-board patterns:** 256-dimensional binary vectors

**generating patterns:** random moves until 10 different positions

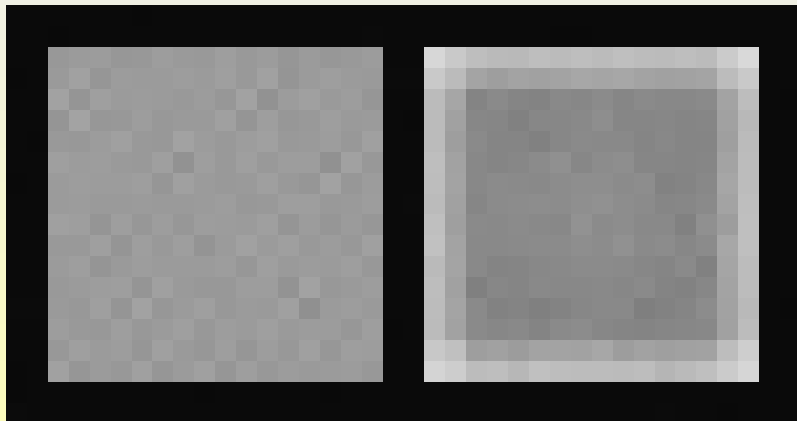
chess-piece position coded by  $x_n = 1$



**training and testing data:** 100000 binary vectors for each class

# Marginal Probabilities of the Chess-Board Patterns

class-means of training numerals (“mean images”)

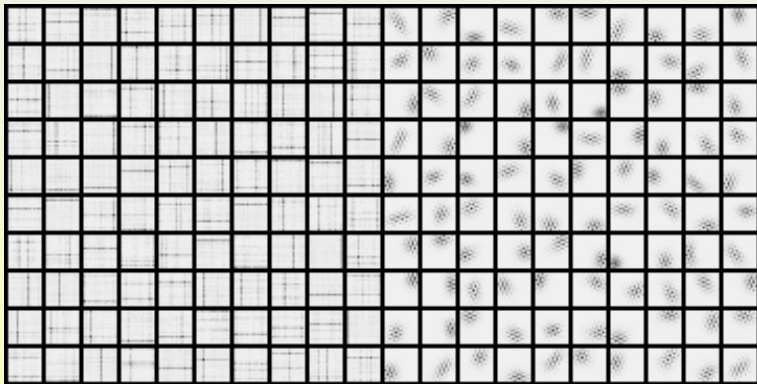


**left:** rook-made patterns

**right:** knight-made patterns

# Component Parameters of the Estimated Mixtures $P(\mathbf{x}|\omega)$

component parameters  $\theta_{mn}$  in chess-board arrangement



**left:** rook-made patterns

**right:** knight-made patterns



# Recognition of Chess-board patterns: full model

Recognition error of the full multivariate Bernoulli mixture model in %

M	1 000	200 000	10 000	200 000	100 000	200 000
2	<b>34.70</b>	41.56	<b>39.59</b>	40.38	<b>39.90</b>	40.02
4	<b>13.10</b>	15.83	<b>16.54</b>	16.65	<b>16.42</b>	16.48
10	<b>1.65</b>	7.72	<b>6.60</b>	7.00	<b>6.49</b>	6.60
20	<b>0.95</b>	9.21	<b>5.40</b>	5.90	<b>4.04</b>	4.34
40	<b>0.15</b>	8.76	<b>3.91</b>	4.90	<b>2.73</b>	2.89
100	<b>0.00</b>	9.35	<b>2.01</b>	4.54	<b>1.37</b>	1.90
200	<b>0.00</b>	11.02	<b>1.22</b>	5.57	<b>0.84</b>	1.68
400	<b>0.00</b>	15.40	<b>0.69</b>	8.35	<b>0.45</b>	1.92
1000	<b>0.00</b>	17.77	<b>0.20</b>	14.66	<b>0.14</b>	3.76

# Recognition of Chess-board patterns: subspace model

Recognition error of the structural Bernoulli mixture model in %

M	1 000	200 000	10 000	200 000	100 000	200 000
4	<b>13.80</b>	16.48	<b>25.53</b>	26.53	<b>16.91</b>	16.92
8	<b>6.45</b>	9.97	<b>10.96</b>	11.32	<b>7.28</b>	7.29
20	<b>5.70</b>	10.97	<b>5.14</b>	5.77	<b>4.70</b>	4.72
40	<b>8.65</b>	13.63	<b>4.20</b>	4.73	<b>3.29</b>	3.32
80	<b>4.20</b>	12.91	<b>6.12</b>	6.88	<b>1.91</b>	1.92
200	<b>0.25</b>	11.46	<b>3.36</b>	4.76	<b>1.83</b>	1.85
400	<b>0.00</b>	18.11	<b>3.54</b>	4.70	<b>3.10</b>	3.20
800	<b>0.00</b>	18.50	<b>3.88</b>	4.82	<b>5.42</b>	5.45
2000	<b>0.00</b>	18.75	<b>2.84</b>	6.39	<b>2.71</b>	2.73

# Concluding Remarks

## Overfitting of Structural Mixture Model

- probabilistic

# Structural Modification of EM Algorithm

**STRUCTURAL OPTIMIZATION:** can be included into EM algorithm

$$L = \frac{1}{|\mathcal{S}_\omega|} \sum_{\mathbf{x} \in \mathcal{S}_\omega} \log \left[ \sum_{m \in \mathcal{M}_\omega} F(\mathbf{x}|0) G(\mathbf{x}|m, \phi_m) w_m \right], \quad F(\mathbf{x}|0) = \prod_{n \in \mathcal{N}} f_n(x_n|0)$$

**EM Algorithm:** ( $m \in \mathcal{M}_\omega, n \in \mathcal{N}, \mathbf{x} \in \mathcal{S}_\omega$ )

$$q(m|\mathbf{x}) = q(m|\mathbf{x}, \omega) = \frac{G(\mathbf{x}|m, \phi_m) w_m}{\sum_{j \in \mathcal{M}_\omega} G(\mathbf{x}|j, \phi_j) w_j},$$

$$w'_m = \frac{1}{|\mathcal{S}_\omega|} \sum_{\mathbf{x} \in \mathcal{S}_\omega} q(m|\mathbf{x}), \quad \theta'_{mn} = \frac{1}{\sum_{\mathbf{x} \in \mathcal{S}_\omega} q(m|\mathbf{x})} \sum_{\mathbf{x} \in \mathcal{S}_\omega} x_n q(m|\mathbf{x})$$






**structural criterion:** (Kullback-Leibler  $\mathbb{I}$ -divergence)

$$\gamma'_{mn} = w'_m \left[ \theta'_{mn} \log \frac{\theta'_{mn}}{\theta_{0n}} + (1 - \theta'_{mn}) \log \frac{(1 - \theta'_{mn})}{(1 - \theta_{0n})} \right]$$






**structural parameter optimization:**  $\phi'_{mn} = 1$  for the  $r$  highest values  $\gamma'_{mn}$

**Remark.** The “structural” EM algorithm converges monotonically.






# References 1/4

-  Dempster, A.P., Laird, N.M., & Rubin, D.B., (1977): Maximum likelihood from incomplete data via the EM algorithm. *J. of the Royal Stat. Soc.*, **B 39**, pp. 1-38.
-  Dietterich, T.: Overfitting and undercomputing in machine learning. *ACM Computing Surveys*, 3 **27** (1995) 326-327
-  Grim J.: Multivariate statistical pattern recognition with nonreduced dimensionality. *Kybernetika*, **22** (1986) 142-157.
-  Grim J.: Information approach to structural optimization of probabilistic neural networks. In proceedings of: *4th System Science European Congress*, Ferrer, L. et al. (Eds.), (pp. 527-540), Valencia: Soc. Espanola de Sistemas Generales, (1999)
-  Grim, J.: Neuromorphic features of probabilistic neural networks. *Kybernetika.*, 5 **43** (2007) 697-712

# References 2/4

-  Grim, J., Pudil, P., Somol, P.: Recognition of handwritten numerals by structural probabilistic neural networks. In: *Proc. Second ICSC Symposium on Neural Computation*. (Bothe, H., Rojas, R. eds.). ICSC, Wetaskiwin (2000) 528-534
-  Grim, J., Just, P., Pudil, P.: Strictly modular probabilistic neural networks for pattern recognition. *Neural Network World*, **13** (2003) 599-615
-  Grim J., Hora, J.: "Iterative principles of recognition in probabilistic neural networks." *Neural Networks, Special Issue, 6* **21** (2008) 838-846
-  McLachlan, G.J., Peel, D.: *Finite Mixture Models*, John Wiley and Sons, New York, Toronto (2000)
-  Vajda, I., (1992): *Theory of Statistical Inference and Information*. Boston: Kluwer.

# References 3/4

-  Sarle, W.S.: Stopped training and other remedies for overfitting. In: Proceedings of the 27th Symposium on the Interface. (1995) Available via <ftp://ftp.sas.com/pub/neural/inter95.ps.Z>.
-  Schaffer C.: Overfitting avoidance as Bias. *Machine Learning*, 2 **10** (1993) 153–178.
-  Schlesinger, M.I.: Relation between learning and self-learning in pattern recognition. (in Russian), *Kibernetika*, (Kiev), No. 2 (1968) 81-88.
-  Specht, D.F.: Probabilistic neural networks for classification, mapping or associative memory. In: *Proc. IEEE Int. Conf. on Neural Networks*, 1 (1988) 525–532
-  Yinyin L., Starzyk J.A., Zhen Zhu: Optimized Approximation Algorithm in Neural Networks Without Overfitting. *IEEE Tran. Neural Networks*, **19** (2008) 983–995