

# O testování shody ROC křivek

J. Antoch, M. Betinec, L. Prchal a P. Sarda

UNIVERZITA KARLOVA V PRAZE  
UNIVERSITÉ PAUL SABATIER, TOULOUSE

4. února 2010

# REAL PROBLEM

# Real problem – collocation extraction

## Problem

Our colleagues from the Institute of Formal and Applied Linguistics (ÚFAL) are developing automatic method for two-word collocation extraction from a text corpus PDT 2.0 comprising more than  $2 \cdot 10^6$  annotated sentences.

## Examples of bigram collocations

- visí otazník – the question mark is hanging – **open question**
- mít pravdu – to have right – **to be right**

## Data set available

- **2557** have been annotated as **true collocations** ...  $\mathcal{C}_1$
- **9675** have been annotated as **normal bigrams** ...  $\mathcal{C}_0$

# Stochastic point of view

## Classification approach

Candidate's chance to be a true collocation is evaluated using a so called **association measure**  $Y \equiv Y(g, \vartheta)$ . ( $\Rightarrow$  **hyperparameter**)

These measures are supposed to separate  $\mathcal{C}_0$  and  $\mathcal{C}_1$  linearly

$$Y \geq \vartheta \implies g \in \mathcal{C}_1$$

$$Y < \vartheta \implies g \in \mathcal{C}_0$$

for a collocation candidate  $g$  and an arbitrary **threshold**  $\vartheta \in \mathbb{R}$ .

## Decision problem

The problem can be seen as a statistical decision  $g \in \mathcal{C}_0$  against  $g \in \mathcal{C}_1$  with a “critical value”  $\vartheta$

## General aim

To measure/display overall performance of  $Y$  for varying  $\vartheta$

# Analogy with statistical testing

## Choice of $\vartheta$

$$Y \geq \vartheta \implies g \in \mathcal{C}_1$$

$$Y < \vartheta \implies g \in \mathcal{C}_0$$

## Small values of $\vartheta$

- $\mathcal{C}_1$  is “preferred”
- the primary interest in “high power”
- **almost all true collocations** are labeled with a big amount of wrongly labeled normal words – **very fine translation**

## Large values of $\vartheta$

- the primary interest in “small level”
- only **the most evident true collocations** are labeled with a small amount of misclassified normal words – **rough translation**

# ROC REMINDER

# Two-sample classification

## Objects

$\mathcal{C}$  ... set of objects

- $\mathcal{C} = \mathcal{C}_0 \cup \mathcal{C}_1, \quad \mathcal{C}_0 \cap \mathcal{C}_1 = \emptyset$
- $\mathcal{C}_0$  ... class **without condition**
- $\mathcal{C}_1$  ... class **with condition**

## Condition

Existence of considered event

- illness
- bonita etc.

## Reality

For  $i = 1, \dots, n$

$$i\text{th object} \in \begin{cases} \mathcal{C}_0 \\ \mathcal{C}_1 \end{cases} \equiv G_i = \begin{cases} 0 \\ 1 \end{cases}$$

# Linear classifier

## Diagnostic variable $Y$

- (discriminant) score, marker etc.
- evaluation of object properties

## Threshold value $\vartheta$ and decision

$$\hat{G}(\vartheta) = \begin{cases} 0 & \text{if } Y \leq \vartheta \\ 1 & \text{if } Y > \vartheta, \end{cases} \quad \vartheta \in \mathbb{R}$$

## Hyperparameter

- given by classification method used up to the value of  
 $\Rightarrow$  **hyperparameter**  $\alpha \in \mathcal{A}$

$$Y : \mathcal{X} \times \mathcal{A} \longrightarrow \mathbb{R}$$

$$(\mathbf{x}, \alpha)^T \longmapsto y$$

- covers most typical methods as  
LDA, QDA, FLDA, LogReg, NNet, SVM, ...



# Evaluation of classifier

## Representation of classifier

- diagnostic variable  $Y$
- fixed choice of hyperparameter  $\alpha$  ... training

## Description of behavior

- for one  $\vartheta \Rightarrow$  one **classification** ...  $\Rightarrow$  **traditional criteria**  
(Acc, Err, risk)
- for all  $\vartheta \in \mathbb{R} \Rightarrow$  all possible classifications ... **ROC curve**

# ROC curve definition

## RANGE OF VALUES

$$\begin{aligned} \mathbf{r} : \mathbb{R} &\longrightarrow [0, 1] \times [0, 1] \\ \vartheta &\longmapsto [\text{FPR}(\vartheta), \text{TPR}(\vartheta)] = [1 - F_0(\vartheta), 1 - F_1(\vartheta)] \end{aligned}$$

True Positive Rate ..... sensitivity, recall, hit rate

$$\text{TPR}(\vartheta) = \mathbf{P} \left( \widehat{G}(Y, \vartheta) = 1 \mid G = 1 \right) = 1 - F_1(\vartheta)$$

False Positive Rate ..... nonspecificity, fallout, alarm rate

$$\text{FPR}(\vartheta) = \mathbf{P} \left( \widehat{G}(Y, \vartheta) = 1 \mid G = 0 \right) = 1 - F_0(\vartheta)$$

$$F_0(y) = \mathbf{P}(Y \leq y \mid G = 0)$$

$$\text{and } F_1(y) = \mathbf{P}(Y \leq y \mid G = 1)$$

# ROC curve definition

Theoretical ROC curve is the range of

$$\begin{aligned}\varrho(\cdot; F_0, F_1) : \mathbb{R} &\rightarrow [0, 1] \times [0, 1] \\ \vartheta &\mapsto [1 - F_0(\vartheta), 1 - F_1(\vartheta)].\end{aligned}$$

where

$$F_0(\vartheta) = \mathbf{P}(Y \leq \vartheta \mid \mathcal{C}_0) \equiv \mathbf{P}(Y_0 \leq \vartheta),$$

$$F_1(\vartheta) = \mathbf{P}(Y \leq \vartheta \mid \mathcal{C}_1) \equiv \mathbf{P}(Y_1 \leq \vartheta)$$

It is a curve in  $[0, 1] \times [0, 1]$  square consisting of  $1 - F_1(\vartheta)$  on the vertical axis plotted against  $1 - F_0(\vartheta)$  on the horizontal axis  $\forall t \in \mathbb{R}$

$$\text{ROC}_Y = \left\{ \mathbf{r} \in [0, 1]^2 : \exists \vartheta \in \mathbb{R} \quad \varrho(\vartheta; F_0, F_1) = \mathbf{r} \right\}$$

# Collocation extraction revisited

## Aim

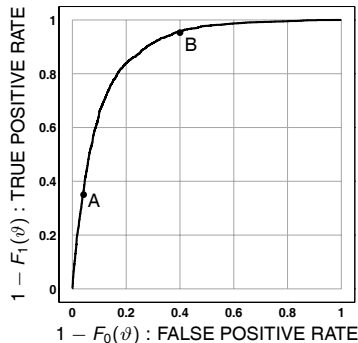
To measure/display overall performance of  $Y$  for varying  $\vartheta$ .

## Rough translation

**A** ... most evident collocations are labeled; level 5 %, power 40 %

## Fine translation

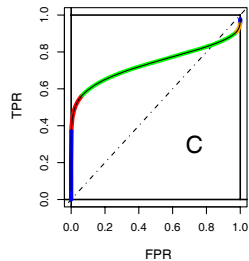
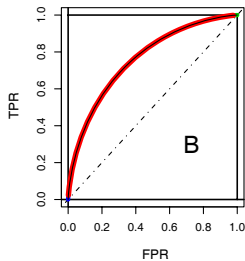
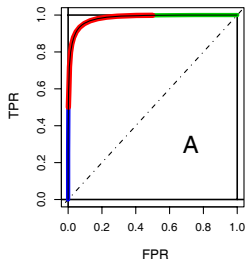
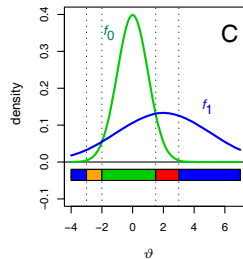
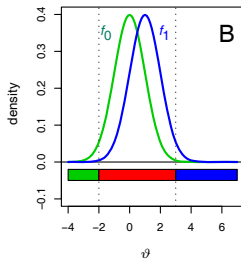
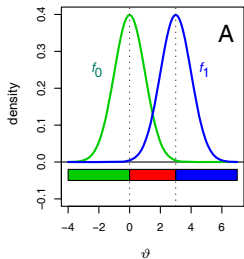
**B** ... almost all collocations are labeled; level 40 %, power 95 %



## Remind

$F_0(\vartheta) = P(Y \leq \vartheta | C_0) \equiv P(Y_0 \leq \vartheta)$  &  $F_1(\vartheta) = P(Y \leq \vartheta | C_1) \equiv P(Y_1 \leq \vartheta)$

# Examples of ROC curves



# Alternative definition

## ROC curve – TPR as function of FPR

$$\text{ROC}(\xi) = \text{TPR}(\xi) = 1 - F_1(F_0^{-1}(1 - \xi)),$$

where  $\xi := \text{FPR} \in [0, 1]$

## Properties

- Assumptions
  - $F_0$  and  $F_1$  are absolutely continuous (may be weakened)
  - supports  $f_0$  and  $f_1$  are identical
- used in parametric models

## Remind

$$\begin{aligned}\text{TPR} &= 1 - F_1(F_0^{-1}(1 - \text{FPR})) = 1 - F_1(F_0^{-1}(1 - (1 - F_0(\vartheta)))) \\ &= 1 - F_1(\vartheta) \equiv \text{TPR}(\vartheta)\end{aligned}$$

# NONPARAMETRIC APPROACH

# ROC curves and confusion matrix

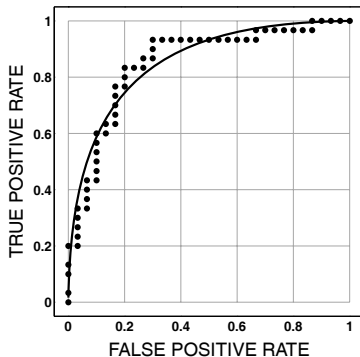
Result of classification for fixed $\vartheta = \vartheta_0$		
Reality	$\mathbf{g} \in \mathcal{C}_0$ (negatives)	$\mathbf{g} \in \mathcal{C}_1$ (positives)
$\mathbf{g} \in \mathcal{C}_0$ (negatives)	true negatives $TN(\vartheta_0)$	false positives $FP(\vartheta_0)$
$\mathbf{g} \in \mathcal{C}_1$ (positives)	false negatives $FN(\vartheta_0)$	true positives $TP(\vartheta_0)$

False positive rate

$$FPR \equiv \frac{FP(\vartheta_0)}{TN(\vartheta_0) + FP(\vartheta_0)}$$

True positive rate

$$TPR \equiv \frac{TP(\vartheta_0)}{TP(\vartheta_0) + FN(\vartheta_0)}$$





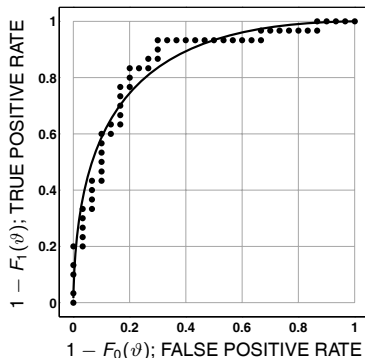
# ROC curves and confusion matrix (cont.)

False positive rate

$$\text{FPR} = \frac{\text{FP}(\vartheta_0)}{\text{TN}(\vartheta_0) + \text{FP}(\vartheta_0)}$$

True positive rate

$$\text{TPR} = \frac{\text{TP}(\vartheta_0)}{\text{TP}(\vartheta_0) + \text{FN}(\vartheta_0)}$$



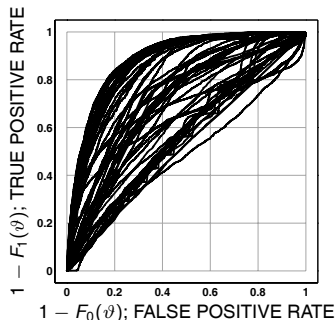
$$\frac{\text{FP}(\vartheta)}{\text{TN}(\vartheta) + \text{FP}(\vartheta)} = \frac{1}{n_0} \text{FP}(\vartheta) = \frac{1}{n_0} \sum_{i=1}^{n_0} \mathbb{I}(Y_{0i} > t) = 1 - \hat{F}_0(\vartheta), \quad \forall \vartheta \in \mathbb{R},$$

$$\frac{\text{TP}(\vartheta)}{\text{TP}(\vartheta) + \text{FN}(\vartheta)} = \frac{1}{n_1} \text{TP}(\vartheta) = \frac{1}{n_1} \sum_{i=1}^{n_1} \mathbb{I}(Y_{1i} > t) = 1 - \hat{F}_1(\vartheta), \quad \forall \vartheta \in \mathbb{R}.$$

# ÚFAL ROC curves

## Linguistic collocation measures

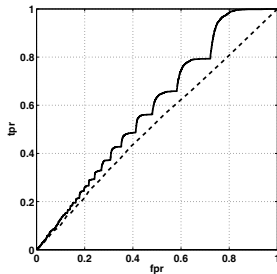
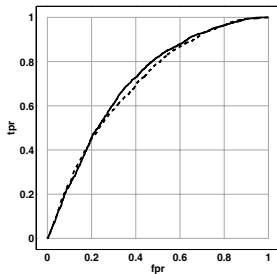
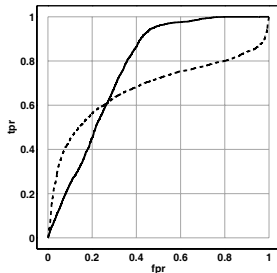
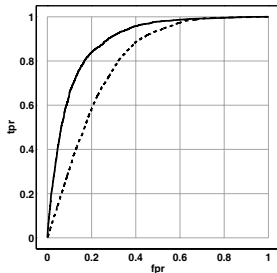
Linguist in ÚFAL use **86** different association measures for collocation extraction.



## Tasks for statisticians

- How to **compare** these measures?
- How to **detect groups** (clusters) of collocation measures that behave analogously?

# Typical linguistic ROC curves



# Basic situation – reminder

## Setup

- **Diagnostic variable**  $Y$  with conditional cdfs  $F_0(\vartheta)$  and  $F_1(\vartheta)$ , i.e.  $F_k(\vartheta) = P(Y \leq \vartheta | \mathcal{C}_k) \equiv P(Y_k \leq \vartheta)$ ,  $k = 0, 1$
- $Y_0$  and  $Y_1$  follow **continuous distributions** with densities  $f_0(\vartheta)$  and  $f_1(\vartheta)$  such that  $f_0(\vartheta) > 0$ ,  $f_1(\vartheta) > 0$  on the same interval  $\mathcal{I}_Y \subseteq \mathbb{R}$ .
- $Y_0$  and  $Y_1$  are independent.

## Theoretical ROC curve is the range of

$$\begin{aligned}\varrho(\cdot; F_0, F_1) : \mathbb{R} &\rightarrow [0, 1] \times [0, 1] \\ \vartheta &\mapsto [1 - F_0(\vartheta), 1 - F_1(\vartheta)]\end{aligned}$$

$$\text{ROC}_Y = \{ \mathbf{r} \in [0, 1]^2 : \exists \vartheta \in \mathbb{R} \quad \varrho(\vartheta; F_0, F_1) = \mathbf{r} \}$$

# Equivalence test – setting for two ROC curves

## Two ROC curves

- $\text{ROC}_Y = \{ \mathbf{r} \in [0, 1]^2 : \exists \vartheta \in \mathbb{R} \quad \varrho(\vartheta; F_0, F_1) = \mathbf{r} \}$
- $\text{ROC}_Z = \{ \mathbf{r} \in [0, 1]^2 : \exists \vartheta \in \mathbb{R} \quad \varrho(\vartheta; G_0, G_1) = \mathbf{r} \}$

## Observed data

- $n_0$  objects from  $\mathcal{C}_0$  and  $n_1$  objects from  $\mathcal{C}_1$
- For each object two (different) measures  $Y$  and  $Z$  are evaluated
- It yields samples  $Y_{01}, \dots, Y_{0n_0}$  distributed according to  $F_0(\vartheta)$ , and  $Y_{11}, \dots, Y_{1n_1}$  distributed according to  $F_1(\vartheta)$
- Analogously,  $Z_{01}, \dots, Z_{0n_0}$  – each follows  $G_0(\vartheta)$ , and  $Z_{11}, \dots, Z_{1n_1} \sim G_1(\vartheta)$

# Equivalence test – idea

## Equivalence of two ROC curves

for us means that for any  $\mathbf{r}_Y \in \text{ROC}_Y$  exists  $\mathbf{r}_Z \in \text{ROC}_Z$  such that

$$\mathbf{r}_Y = \mathbf{r}_Z$$

i.e.

$$\text{ROC}_Y \equiv \text{ROC}_Z \iff$$

$$\forall t_Y \in \mathcal{I}_Y \exists t_Z \in \mathcal{I}_Z : F_0(t_Y) = G_0(t_Z) \text{ and } F_1(t_Y) = G_1(t_Z)$$

## Transformation function

Define  $\tau_0, \tau_1 : \mathcal{I}_Y \rightarrow \mathcal{I}_Z$  such that

$$\tau_0(\vartheta) = G_0^{-1}(F_0(\vartheta)) \quad \text{and} \quad \tau_1(\vartheta) = G_1^{-1}(F_1(\vartheta)) \quad \forall \vartheta \in \mathcal{I}_Y.$$

# Equivalence test – formal definition

## Hypothesis

ROC curves are **equivalent** if and only if  $\tau_0(\vartheta) \equiv \tau_1(\vartheta)$ , i.e.

$$H : \forall \vartheta \in \mathcal{I}_Y \quad \tau_0(\vartheta) = \tau_1(\vartheta),$$

## Alternative

We aim to test H against the **alternative**

$$A : \exists \tilde{\mathcal{I}}_Y \subseteq \mathcal{I}_Y, \tilde{\mathcal{I}}_Y \neq \emptyset \quad \tau_0(\vartheta) \neq \tau_1(\vartheta) \quad \forall \vartheta \in \tilde{\mathcal{I}}_Y$$

$$\tau_0(\vartheta) = G_0^{-1}(F_0(\vartheta)) \quad \text{and} \quad \tau_1(\vartheta) = G_1^{-1}(F_1(\vartheta)) \quad \forall \vartheta \in \mathcal{I}_Y.$$

# Equivalence test – test statistic

## Test statistic

$$T = n \int_{\mathcal{I}_Y^*} (\hat{\tau}_0(\vartheta) - \hat{\tau}_1(\vartheta))^2 d\vartheta,$$

$$\hat{\tau}_0(\vartheta) = \hat{G}_0^{-1}(\hat{F}_0(\vartheta)), \quad \hat{\tau}_1(\vartheta) = \hat{G}_1^{-1}(\hat{F}_1(\vartheta)), \quad \forall \vartheta \in \mathcal{I}_Y,$$

$\hat{F}_k(\vartheta)$  and  $\hat{G}_k(\vartheta)$ ,  $k = 0, 1$ , denote the empirical distribution functions,

$$\hat{G}_k^{-1}(u) = \inf \{t : \hat{G}_k(\vartheta) > u\}, \quad k = 0, 1,$$

and closed interval  $\mathcal{I}_Y^* \subseteq \mathcal{I}_Y$  is chosen such that

$$0 < g_0(\tau_0(\vartheta)) < \infty, \quad 0 < g_1(\tau_1(\vartheta)) < \infty, \quad \forall \vartheta \in \mathcal{I}_Y^*.$$



# Equivalence test – test statistic

## Theorem

Under the null hypothesis, the test statistic  $T$  converges weakly to the infinite weighted sums of independent  $\chi_1^2$  variables  $\eta_1^2, \eta_2^2, \dots$

$$T \xrightarrow[n \rightarrow \infty]{w} T^B = \sum_{j=1}^{\infty} \lambda_j \eta_j^2,$$

where  $\{\lambda_j\}$  represent the eigenvalues of the covariance operator of the zero mean gaussian process  $B(t)$  with the covariance structure

$$\text{cov}(B(s), B(t)) = c_0 \frac{F_0(s)(1 - F_0(t))}{g_0(\tau_0(s))g_0(\tau_0(t))} + c_1 \frac{F_1(s)(1 - F_1(t))}{g_1(\tau_1(s))g_1(\tau_1(t))}, \quad s \leq t,$$

$c_0, c_1$  are positive constants.

# Equivalence test – critical values

To obtain (asymptotic) critical values we need

- 1 to **estimate** eigenvalues  $\{\lambda_j\}$
- 2 to **evaluate** the distribution function of a weighted sum of  $\chi_1^2$  variables

To estimate the covariance structure and  $\lambda_j$ 's

$$\widehat{\text{cov}}(B(s), B(t)) = c_0 \frac{\hat{F}_0(s)(1 - \hat{F}_0(t))}{\tilde{g}_0(\hat{\tau}_0(s))\tilde{g}_0(\hat{\tau}_0(t))} + c_1 \frac{\hat{F}_1(s)(1 - \hat{F}_1(t))}{\tilde{g}_1(\hat{\tau}_1(s))\tilde{g}_1(\hat{\tau}_1(t))},$$

for  $s, t \in \{t_1, \dots, t_p\} \subset \mathcal{I}_Y$ , with  $\tilde{g}_k, k = 0, 1$ , being density kernel estimators. The spectral decomposition of the matrix

$$\left( \widehat{\text{cov}}(B(t_i), B(t_j)) \right)_{i,j=1}^p$$

# Equivalence test – Monte Carlo critical values

## Trimming and Monte Carlo

Suppose that eigenvalues  $\lambda_1, \dots, \lambda_J$  are estimated. It allows an approximation of  $T^B$  by its first  $J$  estimated components

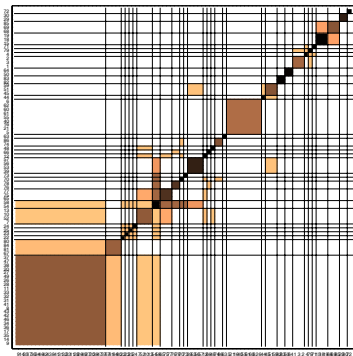
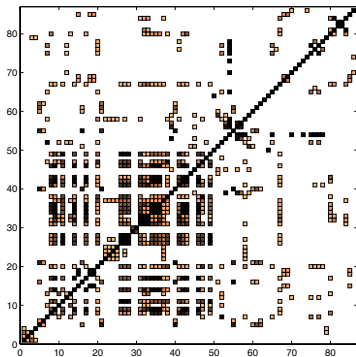
$$T^B \approx \sum_{j=1}^J \hat{\lambda}_j \eta_j^2 = S^J.$$

As distribution of  $S^J$  is not explicitly known, we perform **Monte Carlo simulations** in order to obtain the corresponding quantiles.

# Proximity matrix for linguistic association measures

## Application on linguistic association measures

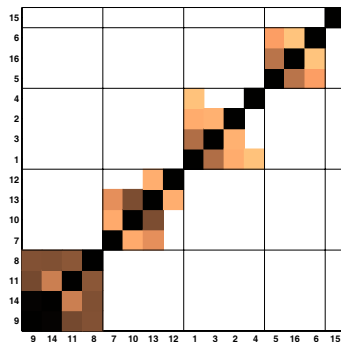
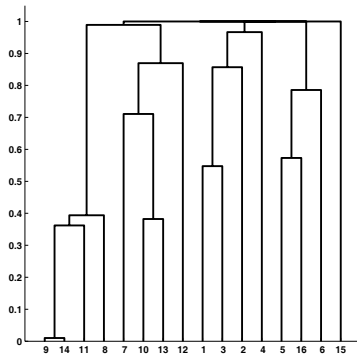
We applied our test on all pairs of 86 collocation association measures and used  $1 - p\text{-value}$  as the proximity distance between two ROC curves.



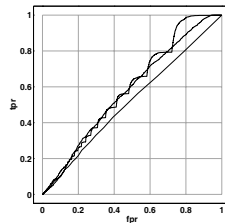
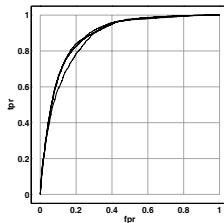
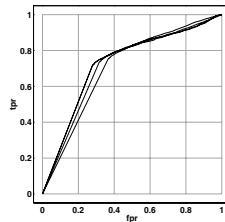
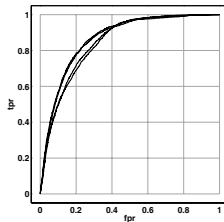
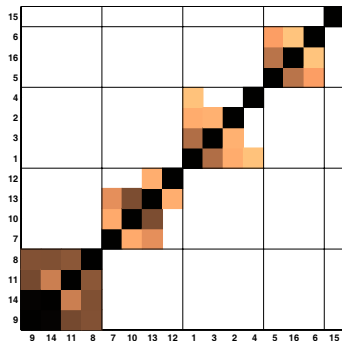
# Collocation extraction – selected results

## Proximity matrix and dendrogram

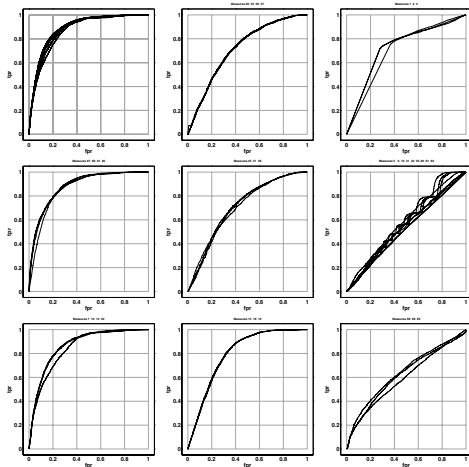
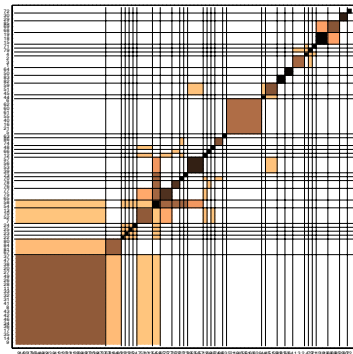
Rearranged (permuted) proximity matrix and the corresponding dendrogram, both providing insight ideas on natural similarity clusters of the observed ROC curves.



# Equivalent classes of selected collocation rules



# Equivalent classes of collocation rules



# Linguistic problem – summary

- ROC curves are useful to display overall performance of a binary classifier
- ROC curve has a theoretical definition
- Statistical theory helps to understand properties of ROC curves and derive new analytical methods
- Our test is essentially based on a proper definition of a ROC curve. However, even straightforward ideas lead to quite complicated theoretical tasks
- Our test can be used to cluster ROC curves
- Clusters may serve to construct a superclassifier more efficient than individual measures



# DĚKUJI