# Nominal Association Measures and Feature Selection for Categorical Data

#### Xiaogang (Steven) Wang

Dept. of Mathematics and Statistics, York University

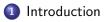
Joint Work with Wenxue Huang and Yong Shi

Dec.12, 2011

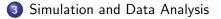
Xiaogang (Steven) Wang (York University)Nominal Association Measures and Feature

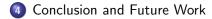
Dec.12, 2011 1 / 20

#### Outline



2 General Nominal Association Measures





Xiaogang (Steven) Wang (York University)Nominal Association Measures and Feature

### Nominal Associations

- We are interested in association measures for nominal data.
- The response variables have more than two categories.
- Goodman and Kruskal have argued that many nominal association measures are based on  $\chi^2$  statistic for testing independence.

#### Goodman and Kruskal's au

Goodman and Kruskal (1954) proposed the following measure:

$$\tau^{Y|X} = \frac{\sum_{i=1}^{n_Y} \sum_{j=1}^{n_X} p(Y=i; X=j)^2 / p(X=j) - \sum_{i=1}^{n_Y} p(Y=i)^2}{1 - \sum_{i=1}^{n_Y} p(Y=i)^2}, \quad (1)$$

where  $n_X$  and  $n_Y$  represent the number of classes/categories for X and Y respectively.

### Goodman and Kruskal's $\tau$

The principle can be stated as

$$r(Y|X) = [V(Y) - V(Y|X)]/V(Y)$$

where V symbolizes certain measure of uncertainty.

Goodman and Kruskal's  $\tau$  can be derived by using Gini concentration index defined as

$$V_G(X) := \sum_{i=1}^{n_X} p(X=i)(1-p(X=i)).$$
(3)

Xiaogang (Steven) Wang (York University)Nominal Association Measures and Feature

(2)

#### Association Vector and General Association Degree

The association vector

$$\Theta^{Y|X} := (\theta^{(Y=1)|X}, \theta^{(Y=2)|X}, \dots, \theta^{(Y=n_Y)|X})$$

$$\tag{4}$$

is given by

$$\theta^{(Y=s)|X} := \frac{E[p(Y=s|X)^2] - p(Y=s)^2}{p(Y=s)(1-p(Y=s))}, \qquad s = 1, 2, \dots, n_Y.$$
(5)

The global association degree is defined as

$$\tau_{\alpha}^{Y|X} = \sum_{s=1}^{n_Y} \alpha_s \, \theta^{Y=s|X}.$$
 (6)

with  $\sum_{s} \alpha_{s} = 1$  and  $\alpha_{s} \ge 0$  for all  $s = 1, 2, \dots, n_{Y}$ ,

Dec.12, 2011 6 / 20

#### **Theoretical Properties**

**THEOREM 1.** Assume  $\alpha$  is a regular weight vector.

$$\begin{array}{l} \bullet \quad 0 \leq \theta^{(Y=s)|X} \leq 1 \text{ and } 0 \leq \tau_{\alpha}^{Y|X} \leq 1; \\ \bullet \quad \tau_{\alpha}^{Y|X} = 0 \iff Y \text{ and } X \text{ are independent.} \end{array}$$

- $\ \, \bullet \ \, \tau_{\alpha}^{Y|X} = 1 \iff Y \text{ is completely determined by } X;$
- If the weight vector is assigned as in the following:

$$\alpha^{P} = \frac{1}{V_{G}(Y)} (p(Y=1) - p(Y=1)^{2}, \dots, p(Y=n_{Y}) - p(Y=n_{Y})^{2});$$
(7)

where 
$$V_G(Y) = \sum_{s} p(Y = s)(1 - p(Y = s))$$
, then

$$\tau^{Y|X} = \tau^{Y|X}_{\alpha^P}.$$
(8)

Dec.12, 2011 7 / 20

#### Association Matrix

The association matrix is given by

$$\gamma(Y|X) := (\gamma^{st}(Y|X)), \tag{9}$$

where

$$\gamma^{st}(Y|X) := \frac{E[p(Y=s|X) \ p(Y=t|X)]}{p(Y=s)},$$

where  $s, t = 1, 2, \dots, n_Y$ .

We have the following properties:

γ(Y|X) is a row stochastic matrix;
 θ<sup>(Y=s)|X</sup> is the normalization of γ<sup>ss</sup>(Y|X);

→ □ → → □ → □ → ○ ○ ○

# Hierarchical Equivalence Structure

- *E-1* equivalent, if  $\tau^{X_1|X_2} = \tau^{X_2|X_1} = \tau^{Y|X_1} = 1$ ;
- 2 *E-2* equivalent, if  $\tau^{Y|X_1} = 1 = \tau^{Y|X_2}$ ;
- **3** *E-3* equivalent, if  $\gamma(Y|X_1) = \gamma(Y|X_2)$ ;
- E-4 equivalent, if  $\Theta^{Y|X_1} = \Theta^{Y|X_2}$ ;
- **③** E-5 equivalent with respect to a weight vector  $\alpha$ , if  $\tau_{\alpha}^{Y|X_1} = \tau_{\alpha}^{Y|X_2}$ .

**THEOREM 2.** If  $X_1$  and  $X_2$  are E-i equivalent (with respect to Y), then they are E-(i+1) equivalent (with respect to Y), for i = 1, 2, 3, 4.

・「「・・」・ 「」・ 「」・

### Structural Base

#### Definition

A subset

$$\{V_{i_1}, V_{i_2}, \ldots, V_{i_k}\} \subseteq V$$

is called an  $\alpha$ -structural base for S for a given regular weight vector  $\alpha$  if  $\alpha$ B1. for each  $V_i \in V$ ,  $\tau_{\alpha}^{V_i|(V_{i_1}, V_{i_2}, ..., V_{i_k})} = 1$ ;  $\alpha$ B2. for any  $V \in \{V_{i_1}, \ldots, V_{i_k}\}$ ,  $\tau_{\alpha}^{V|(\{V_{i_1}, \ldots, V_{i_k}\} \setminus \{V\})} < 1$ .

**THEOREM 3.** Let S be a data set with independent variables  $X_1, \ldots, X_n$  and a dependent variable Y, and  $\alpha$  a weight vector. Then there exists an  $\alpha$ -association base.

# ALGORITHM

- Ocompute all pairwise associations and find the best one.
- **2** Compute the augmented associations by adding one more variables.
- Sind the best combination from Step 2.
- Repeat Step 2 -3 until a stopping rule is reached.

The computational complexity for computing each  $\tau_{\alpha}$  is approximately O(n m). Bootstrapped confidence interval might be needed when the sample size is moderate.

### SIMULATION

The joint distribution of (Y, X1, X2) is given by the following:

$(X_1, X_2)$	P(X1; X2)	$P(Y=0 X_1,X_2)$	$P(Y=1 X_1;X_2)$	$P(Y=2 X_1;X_2)$
(0, 0)	9/16	95%	5%	0%
(0, 1)	3/16	30%	70%	0%
(1, 0)	3/16	50%	50%	0%
(1, 1)	1/16	0%	5%	95%

We also generate redundant variables X3, X4 and X5.  $P(X_3 = 1|X_1 = 1) = P(X_4 = 1|X_2 = 1) = 0.90.$   $X_5 = I(X_2 = 1) * I(X_3 = 1) * Z$ , where P(Z = 1) = 0.8.Note that  $P_Y = (0.6875, 0.2531, 0.0594).$ 

# Weighting Schemes

For comparison purposes, we consider three weighting schemes:

$$\alpha_{k}^{GK} = \frac{p(Y=k)(1-p(Y=k))}{\sum_{j=1}^{n_{Y}} p(Y=j)(1-p(Y=j))};$$
(10)  

$$\alpha_{k}^{EW} = \frac{1}{n_{Y}};$$
(11)  

$$\alpha_{k}^{IPW} = \frac{1}{p(Y=k)} / [\sum_{j=1}^{n_{Y}} \frac{1}{p(Y=j)}];$$
(12)

assuming that p(Y = k) > 0 for any k.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

# SIMULATION RESULTS

	$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
$\tau(GK)$	0.2382	0.1010	0.2060	0.0878	0.1511
$\tau$ (EW)	0.2221	0.1206	0.1923	0.1050	0.2943
$\tau$ (IPW)	0.2382 0.2221 0.1900	0.1597	0.1648	0.1393	0.5806

	$X_1 + X_4$	$X_2 + X_3 + X_5$	$X_1 + X_4 + X_5$	$X_1 + X_2$	ALL
$\tau$ (GK)	0.4627	0.4669	0.4823	0.5018	0.5018
$\tau$ (EW)	0.5570	0.5731	0.5884	0.6078	0.6078
$\tau$ (IPW)	0.7372	0.7940	0.8004	0.8198	0.8199

э

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# DATA ANALYSIS

Data Set: n=24000; m=120.

$\tau$ (GK)	0.0825	0.0760	0.0719	0.0680	0.0540
Variables	V34	V2	V48	V8	V12
$\tau$ (EW)	0.0911	0.0835	0.0809	0.0787	0.0638
Variables	V34	V2	V8	V48	V12
$\tau$ (IPW)	0.1075	0.0888	0.0773	0.0682	0.0605
Variables	V8	V34	V12	V111	V2

Table: Measure of associations using different weighting schemes for significant variables and some of their combinations.

A B + A B +

# DATA ANALYSIS

Data Set: n=24000; m=120.

$\tau(GK)$	0.1195	0.1136	0.1076	0.1021	0.1013
V34	+V2	+V48	+V8	+V7	+V5
$\tau$ (EW)	0.1333	0.1268	0.1268	0.1128	0.1107
V34	+V2	+V8	+V48	+V7	+V5
$\tau$ (IPW)	0.1945	0.1779	0.1598	0.1589	0.1569
V8	+V1	+V22	+V34	+V2	+V48

Table: Top five measure of associations based on the combination of two variable. The best single variables are three weighting schemes are V34, V34, and V2, V8 respectively.

A B F A B F

# DATA ANALYSIS

Data Set: n=24000; m=120.

$\tau$ (GK)	0.1454	0.1397	0.1393	0.1380	0.1368
V34+V2	+V8	+V102	+V5	+V12	+V111
$\tau$ (EW)	0.1709	0.1583	0.1563	0.1554	0.1554
V34+V2	+V8	+V12	+V111	+V102	+V5
$\tau$ (IPW)	0.2576	0.2527	0.2379	0.2333	0.2321
V8+V1	+V2	+V48	+V7	+V34	+V102

Table: Top five measure of associations based on the combinations of triple variables by using V34+V2 and another explanatory variable using different weighting schemes.

A B + A B +



- We introduce a association vector and matrix to measure nominal associations.
- These measures provide both local and global evaluations.
- We also discover the hierarchical equivalence structure.
- The existence of a structural basis is proved.

#### Future Works

- Extension to ordinal categorical data.
- Bayesian approach when the sample size is small.
- Application to risk management and biological data.

- Huang, W., Shi, Y and Wang, X. (2011). Nominal Association Measures and Feature Selection for Categorical Data. To be submitted.
- Goodman, L.A. (1996). JASA 91, 408-428.
- Goodman, L.A. and Kruskal, W. H. (1954). JASA, 49, 732-764.
- Lloyd, C. J. (1999). Statistical analysis of categorical data.