

What is relevant and what is redundant?

M. ROSÁRIO OLIVEIRA

CEMAT AND DEPT. MATEMÁTICA, IST-ULISBOA

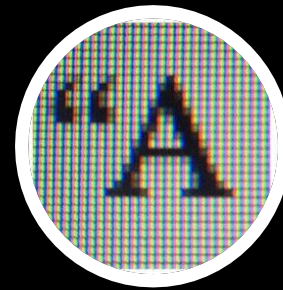


FEATURE SELECTION: Only the “*right*” features



FEATURE SELECTION: THE RIGHT DATA

More data is not necessarily
more information...



Feature Selection:

- Extract from the data **useful** and **valuable** knowledge for **real problem** solving



FEATURE SELECTION: THE RIGHT DATA

- Select a **small subset** of the original features
- Such that we remove **irrelevant** and **redundant** features
- In order to:
 - Reduce **computational complexity**
 - Improve model **accuracy**
 - Increase model **interpretability**

Feature Selection Methods

(classification)

Classifier dependent

Classifier independent

Wrapper

Embedded

Filter



WRAPPER METHODS

- **Idea:** search for feature subsets, using the **classifier accuracy** as the measure of **utility** for a candidate subset
- **Disadvantages:**
 - computational cost
 - selected features are classifier specific
- **Example:**
 - Stepwise regression



EMBEDDED METHODS

- **Idea:** Classifier estimations and feature selection are not separated and interact
- **Disadvantages:**
 - Selected features are classifier specific
 - Regularized_OF = OF + λ regularization_penalty
- **Example:**
 - Regularization methods



FILTER METHODS

- **Idea:** Classifier's estimation and feature selection are separated and depend on a specific measure of benefit
- Most popular ones: rely on **Mutual Information (MI)** and **Entropy**
- **Mutual Information:** measures **linear** and **non-linear** associations among features
- **Example:**
 - Forward feature selection methods based on MI

ENTROPY, MUTUAL INFORMATION



Entropy

Entropy

Motivated by problems in the field of telecommunications

A Mathematical Theory of Communication*

C. E. Shannon (1948)

- A measure of uncertainty
- One formula that changed the world...



Entropy Discrete rv

$$H(\mathbf{X}) = - \sum_{\mathbf{x} \in \mathcal{X}} P(\mathbf{X} = \mathbf{x}) \ln P(\mathbf{X} = \mathbf{x}).$$

- Does not depend on the values of X , only on its prob.
- $H(a)=0$
- $H(X) \geq 0$, Non-negative
- $H(X) = \ln(n)$, $X \sim \text{Unif}\{a_1, \dots, a_n\}$, is maximum



Differential Entropy Continuous rv

$$h(\mathbf{X}) = - \int_{\mathbf{x} \in \mathcal{X}} f_X(\mathbf{x}) \ln f_X(\mathbf{x}) d\mathbf{x}.$$

- Does not depend on the values of X , only on its pdf
- Can be negative
- $h(X) = \ln(a)$, $X \sim \text{Unif}(0, a)$,
 - $a=1$, $h(X)=0$
 - $a<1$, $h(X)<0$

Mutual Information

Mutual Information Discrete rv

$$MI(X, Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} P(X = x, Y = y) \ln \frac{P(X = x, Y = y)}{P(X = x)P(Y = y)}$$

Mutual Information Continuous rv

$$MI(X, Y) = \int_{y \in \mathcal{Y}} \int_{x \in \mathcal{X}} f_{X,Y}(x, y) \ln \frac{f_{X,Y}(x, y)}{f_X(x)f_Y(y)} dx dy$$

- Measures linear and non-linear associations between X and Y

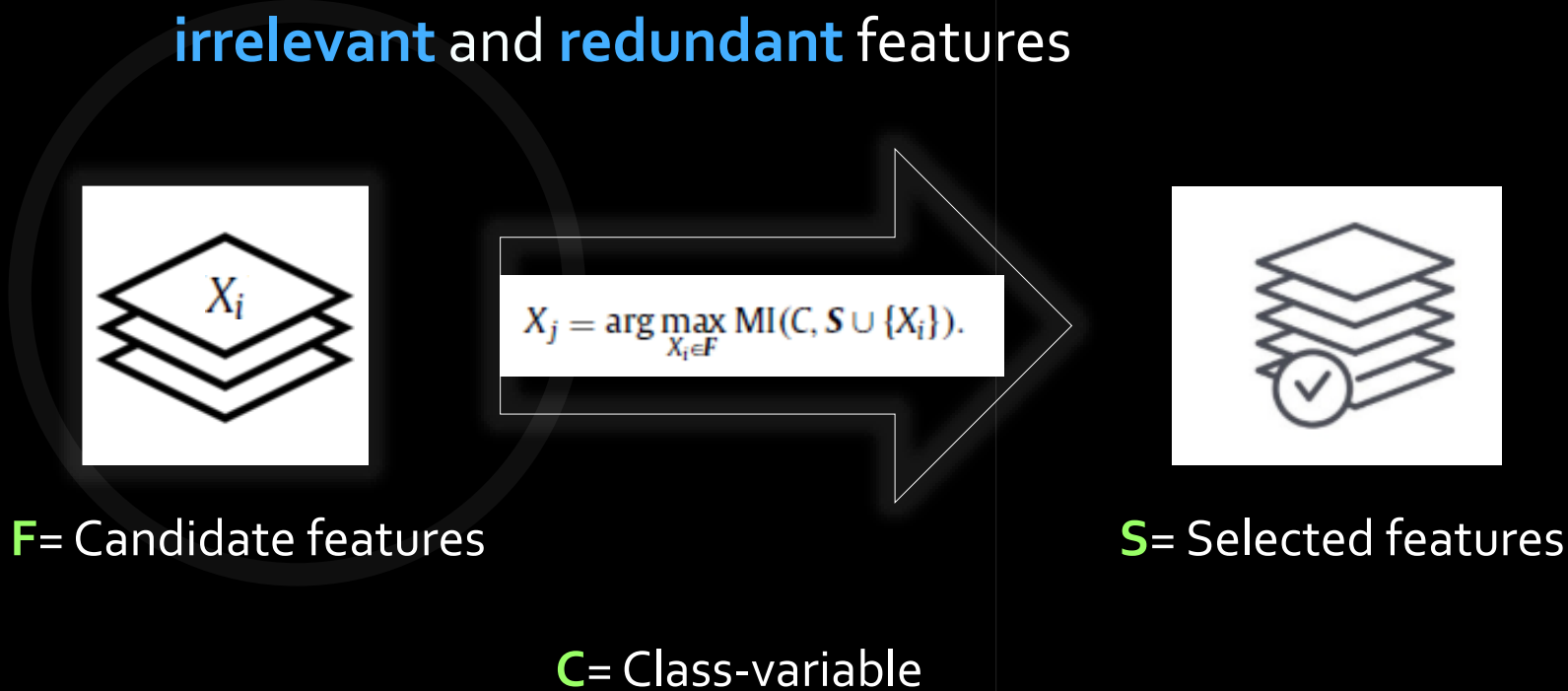
- $MI(X, Y) \geq 0$
- Symmetric
- $MI(X, Y) = 0$ iff $X \perp\!\!\!\perp Y$
- $MI(X, X) = H(X)$
- All properties hold, except
- $MI(X, X) = +\infty$

FORWARD FEATURE SELECTION

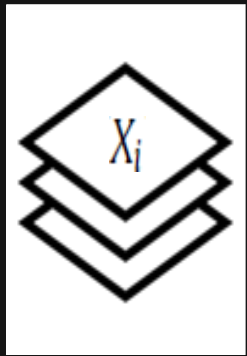


FORWARD FEATURE SELECTION

Goal: Select a **small subset** of the original features, excluding **irrelevant** and **redundant** features



FORWARD FEATURE SELECTION



$$X_j = \arg \max_{X_i \in S} MI(C, S \cup \{X_i\}).$$



Problem:

$$\begin{aligned} OF(X_i) &= MI(C, S) + MI(C, X_i | S) \\ &= MI(C, S) + MI(C, X_i) - TMI(C, X_i, S) \\ &= MI(C, S) + MI(C, X_i) - MI(X_i, S) + MI(X_i, S | C). \end{aligned}$$

Calculation / Estimation

OBJECTIVE FUNCTIONS: INTERPRETABILITY

$$X_j = \arg \max_{X_i \in F} \text{MI}(C, S \cup \{X_i\}).$$

max

$$\text{OF}(X_i) = \text{MI}(C, S) + \text{MI}(C, X_i | S)$$

$$\text{OF}(X_i) = H(C) - H(C | X_i, S) \quad \text{min}$$

$$\begin{aligned} \text{OF}(X_i) &= \text{MI}(C, S) + \text{MI}(C, X_i | S) \\ &= \text{MI}(C, S) + \text{MI}(C, X_i) - \text{TMI}(C, X_i, S) \\ &= \text{MI}(C, S) + \text{MI}(C, X_i) - \text{MI}(X_i, S) + \text{MI}(X_i, S | C). \end{aligned}$$

min

max

OBJECTIVE FUNCTIONS: INTERPRETABILITY

$$X_j = \arg \max_{X_i \in F} MI(C, S \cup \{X_i\}).$$

$$\begin{aligned} OF(X_i) &= MI(C, S) + MI(C, X_i|S) \\ &= MI(C, S) + MI(C, X_i) - TMI(C, X_i, S) \\ &= \cancel{MI(C, S)} + MI(C, X_i) - MI(X_i, S) - MI(X_i, S|C). \end{aligned}$$

Relevance



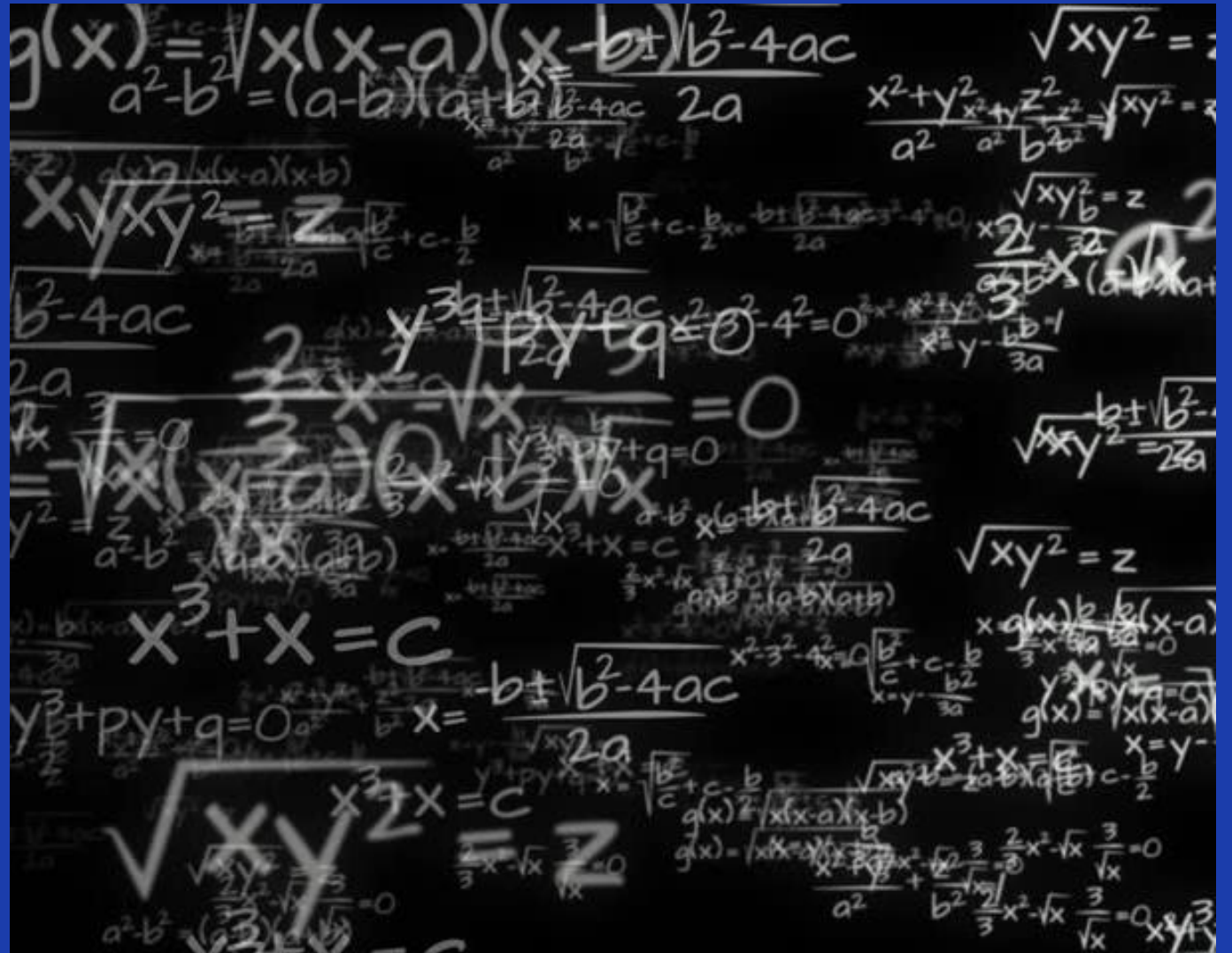
Inter-feature
redundance



Class-relevant
Redundancy



FEATURE SELECTION METHODS



FORWARD FEATURE SELECTION METHODS

3RD GROUP

$$OF(X_i) = MI(C, S) + MI(C, X_i) - MI(X_i, S) + MI(X_i, S|C).$$

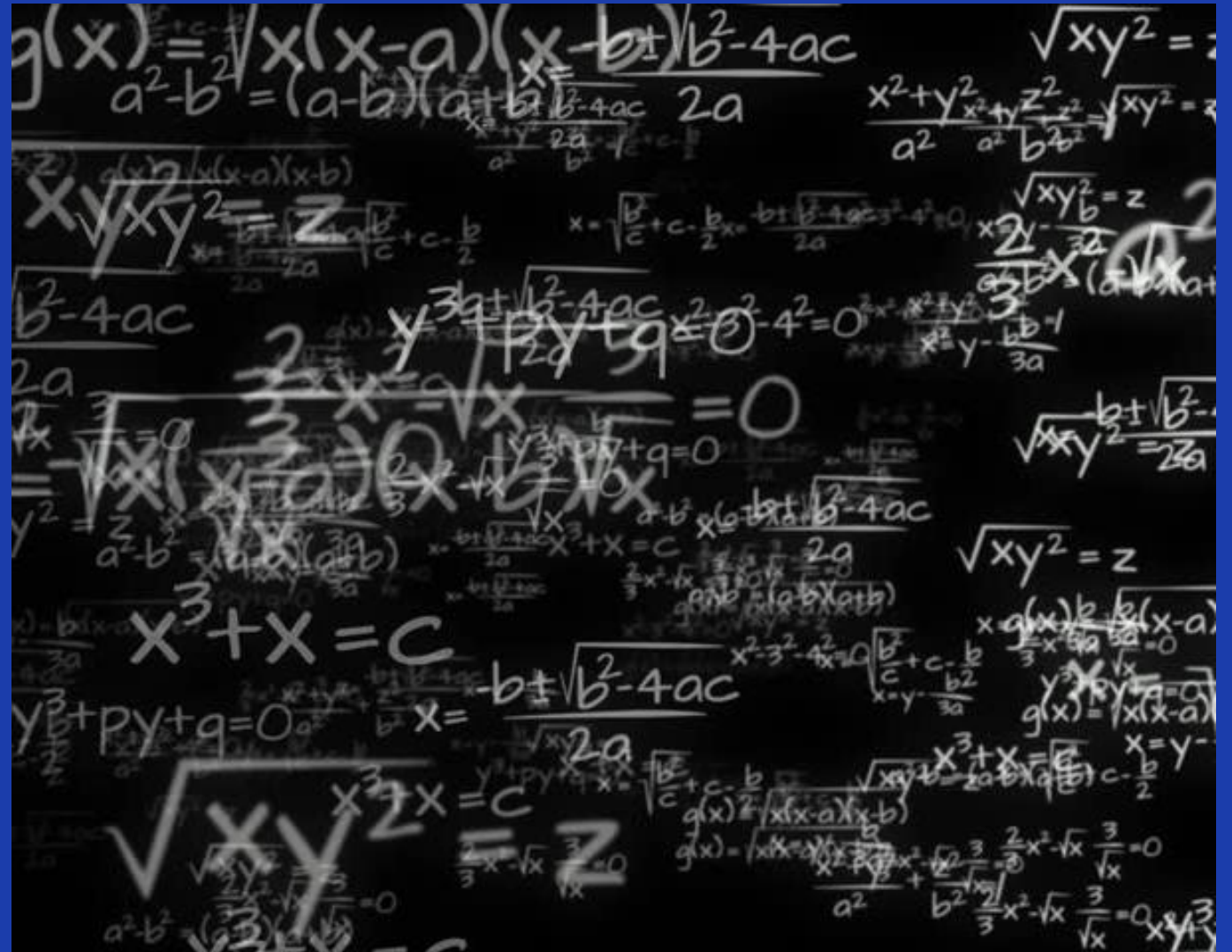
Third Group of Methods

Method	Objective function evaluated at X_i
CIFE	$MI(C, X_i) - \sum_{X_s \in S} (MI(X_i, X_s) - MI(X_i, X_s C))$
JMI	$MI(C, X_i) - \frac{1}{ S } \sum_{X_s \in S} (MI(X_i, X_s) - MI(X_i, X_s C))$
CMIM	$MI(C, X_i) - \max_{X_s \in S} \{MI(X_i, X_s) - MI(X_i, X_s C)\}$
JMIM	$MI(C, X_i) - \max_{X_s \in S} \{MI(X_i, X_s) - MI(X_i, X_s C) - MI(C, X_s)\}$

$$OF_{DMIM}(X_i) = MI(X_i, C) - \max_{X_s \in S} MI(X_i, X_s) + \max_{X_s \in S} MI(X_i, X_s|C).$$

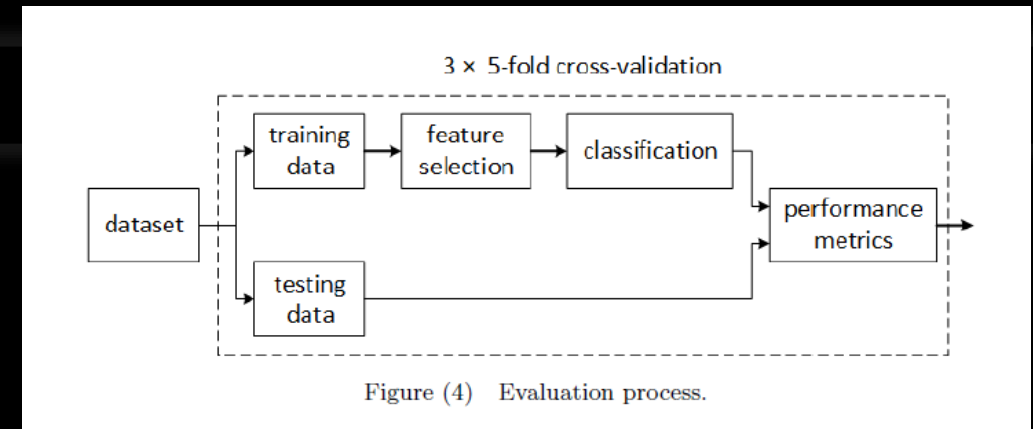
- **Class-relevant redundancy:** contribution of a candidate feature to the explanation of the class, when taken together with already selected features

THEORETICAL COMPARISON



NUMERICAL COMPARISON:

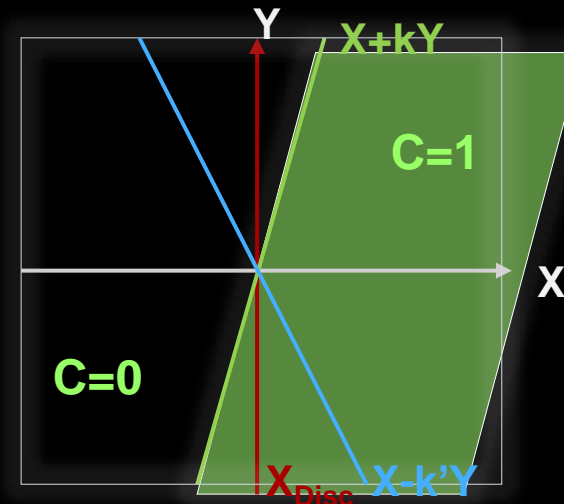
How comparisons are usually done:



Theoretical Setup:

Class-Variable: $C_k = \begin{cases} 0, & X + kY < 0 \\ 1, & X + kY \geq 0 \end{cases}$

Candidate Features: $X, X - k'Y, X_{Disc}, Z$



2. THEORETICAL COMPARISON:

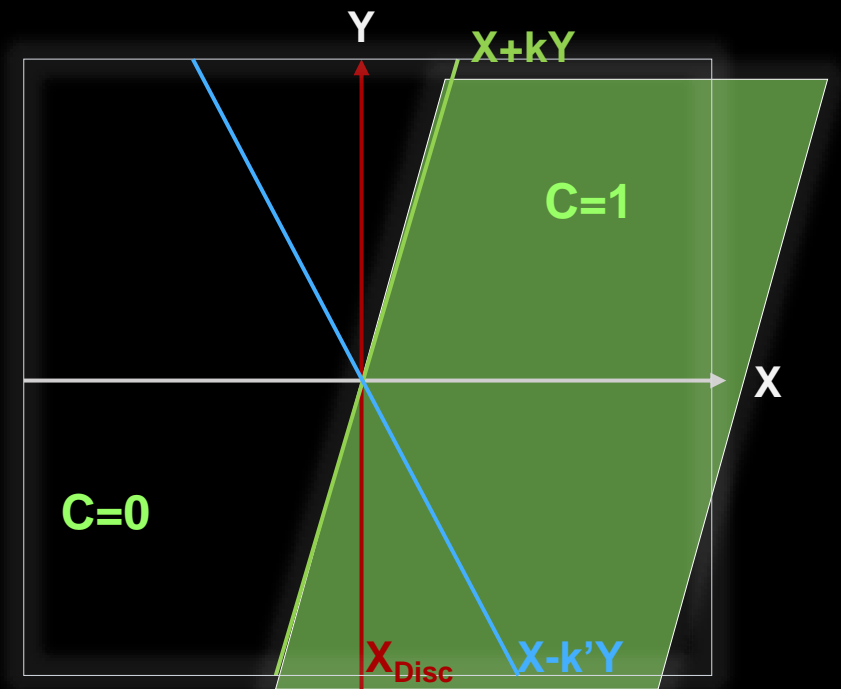
USING A DISTRIBUTIONAL SETTING

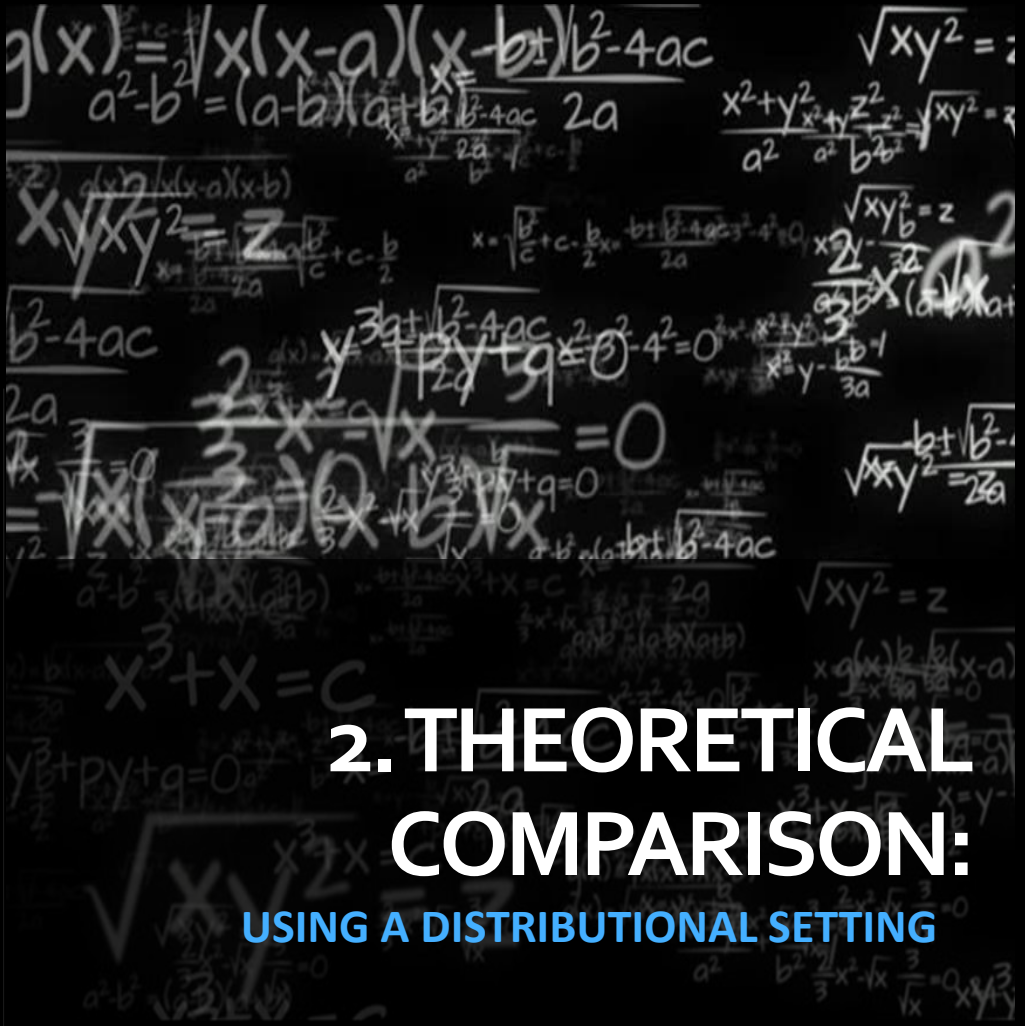
FEATURE (RELEVANCE) TYPES

- Class-Variable: $C_k = \begin{cases} 0, & X + kY < 0 \\ 1, & X + kY \geq 0 \end{cases}$
- Candidate Features: $X, X - k'Y, X_{Disc}, Z$

Features Categories:

- Irrelevant: Z
- Relevant: $X, X - k'Y, X_{Disc}$
 - Fully Relevant: $X - k'Y$
- Redundant:
 - If we chose X then X_{Disc} is redundant
 - If we chose X_{Disc} then X is redundant





2. THEORETICAL COMPARISON:

USING A DISTRIBUTIONAL SETTING

Features Order: OF were calculated theoretically assuming X, Y , and Z are indep. $N(0,1)$

Performance Measure:
Bayes Risk and Bayes Classifier
=
 $\min\{\text{Total Probability of Misclassification}\}$

THEORETICAL COMPARISON

A	MI(C _k , A)
X	$\frac{1}{2} \ln(2\pi e) - \frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X C_k=j}(u) \ln f_{X C_k=j}(u) du$
X - k'Y	$\frac{1}{2} \ln(2\pi e(1+k'^2)) - \frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X-k'Y C_k=j}(u) \ln f_{X-k'Y C_k=j}(u) du$
Z	0
X _{disc}	$2 \ln(2) + \frac{\arctan k}{\pi} \ln\left(\frac{\arctan k}{\pi}\right) + (1 - \frac{\arctan k}{\pi}) \ln\left(\frac{1 - \arctan k}{2\pi}\right)$

- ^a $X|C_k = j \sim \text{SN}(0, 1, \frac{(-1)^{j+1}}{k})$, $j = 0, 1$.
^b $X - k'Y|C_k = j \sim \text{SN}(0, \sqrt{1+k'^2}, (-1)^{j+1} \frac{(1-kk')}{(k+k')})$, $j = 0, 1$.

Table 8
MI between pairs of input features.

A	B	MI(·, ·)
X	X - k'Y	$\frac{1}{2} \ln(1 + \frac{1}{k'^2})$
X	X _{disc}	$\ln(2)$
X - k'Y	X _{disc}	$\frac{1}{2} \ln(2\pi e) - \frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X C_k=j}(u) \ln f_{X C_k=j}(u) du$
Z	B	0, $B \in \{X, X - k'Y, X_{\text{disc}}\}$

- ^a $X|C_k = j \sim \text{SN}(0, 1, \frac{(-1)^{j+1}}{k})$, $j = 0, 1$.

A	B	MI(·, · C _k)
X	X - k'Y	$\frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X C_k=j}(u) \ln f_{X C_k=j}(u) du - \frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X-k'Y C_k=j}(u) \ln f_{X-k'Y C_k=j}(u) du$
X	X _{disc}	$-\frac{\arctan k}{\pi} \ln\left(\frac{\arctan k}{\pi}\right) - (1 - \frac{\arctan k}{\pi}) \ln\left(\frac{1 - \arctan k}{2\pi}\right)$
X - k'Y	X _{disc}	$\frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X-k'Y C_k=j}(u) \ln f_{X-k'Y C_k=j}(u) du$
Z	B	0, $B \in \{X, X - k'Y, X_{\text{disc}}\}$

- ^a $X|C_k = j \sim \text{SN}(0, 1, \frac{(-1)^{j+1}}{k})$, $j = 0, 1$.
^b $X - k'Y|C_k = j \sim \text{SN}(0, \sqrt{1+k'^2}, (-1)^{j+1} \frac{(1-kk')}{(k+k')})$.
^c $h(X - k'Y|X_{\text{disc}}, C_k)$ in (A.5).

$$= \frac{\sqrt{2}}{\sqrt{1+k'^2}} \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}\sqrt{1+k'^2}(k+k')} \times \Phi\left(-\frac{(1-kk')z}{(k-k')\sqrt{1+k'^2}}\right) dz$$

- ^a $X|C_k = j \sim \text{SN}(0, 1, \frac{(-1)^{j+1}}{k})$, $j = 0, 1$.

The desired probability density function is $f_{X-k'Y|X_{\text{disc}}=1, C_k=0}(v) = \begin{cases} \frac{2\pi}{\pi - \arctan k} \int_{-\frac{v}{k'}}^{\frac{v}{k}} \phi(z) \phi(v+k'z) dz, & v \geq 0 \\ 0, & v < 0 \end{cases}$

We finally consider $u = 0$ and $j = 0$. For $v \geq 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 0, C_k = 0) = \int_{-\infty}^0 \int_{-\infty}^{\infty} \phi(w) \phi(z) dw dz - \int_0^{\infty} \int_{-kz}^0 \phi(w) \phi(z) dz dw - \int_{-\infty}^{-\frac{v}{k'}} \int_{v+kz}^0 \phi(w) \phi(z) dz dw - \frac{1}{2} \int_{-\infty}^0 \phi(z) [\Phi(-kz) - \frac{1}{2}] dz - \int_{-\infty}^{-\frac{v}{k'}} [\frac{1}{2} - \Phi(v+k'z)] \phi(z) dz - \frac{3}{4} - \frac{1}{2} F_{\text{SN}(0,1,-k)}(0) - \frac{1}{2} \Phi(-\frac{v}{k'}) + \int_{-\infty}^{-\frac{v}{k'}} \Phi(v+k'z) \phi(z) dz$

As for the case $v < 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 0, C_k = 0) = \int_{-\infty}^{\frac{v}{k'}} \int_{-\infty}^{v+kz} \phi(w) \phi(z) dw dz + \int_{-\frac{v}{k'}}^{\infty} \int_{-\infty}^{-kz} \phi(w) \phi(z) dz dw - \int_{-\infty}^{\frac{v}{k'}} \Phi(v+k'z) \phi(z) dz + \int_{-\frac{v}{k'}}^{\infty} \Phi(-kz) \phi(z) dz - \int_{-\infty}^{\frac{v}{k'}} \Phi(v+k'z) \phi(z) dz + \frac{1}{2} - \frac{1}{2} F_{\text{SN}(0,1,-k)}(-\frac{v}{k'})$

We again need to take the derivative of the two expressions with respect to v to obtain the corresponding conditional density functions. In the case of $v \geq 0$, $\frac{d}{dv} [\frac{1}{4} + \frac{1}{2} F_{\text{SN}(0,1,k)}(0) - \frac{1}{2} \Phi(-\frac{v}{k'}) + \int_{-\infty}^{-\frac{v}{k'}} \Phi(v+k'z) \phi(z) dz]$

A	B	MI(·, ·)
X	X - k'Y	$\frac{1}{2} \ln(1 + \frac{1}{k'^2})$
X	X _{disc}	$\ln(2)$
X - k'Y	X _{disc}	$\frac{1}{2} \ln(2\pi e) - \frac{1}{2} \sum_{j=0}^1 \int_{\mathbb{R}} f_{X C_k=j}(u) \ln f_{X C_k=j}(u) du$
Z	B	0, $B \in \{X, X - k'Y, X_{\text{disc}}\}$

- ^a $X|C_k = j \sim \text{SN}(0, 1, \frac{(-1)^{j+1}}{k})$, $j = 0, 1$.

$= \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz$, while, for $v < 0$, $\frac{d}{dv} [\int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz - \frac{1}{2} \Phi(\frac{v}{k'})] = \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz + \frac{1}{2} \phi(\frac{v}{k'}) \frac{1}{k'} - \frac{1}{2} \phi(\frac{v}{k'}) \frac{1}{k'}$

Note that the following important result [51, Ch. 3] was required in order to obtain both final expressions above: $\frac{d}{dx} \int_{a(x)}^{b(x)} g(x, y) dy = \int_{a(x)}^{b(x)} \frac{dg(x, y)}{dx} dy + g(x, b(x)) b'(x) - g(x, a(x)) a'(x)$. (A.4)

This result was applied to $\frac{d}{dv} \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz$, in the expression for $v \geq 0$, and also to $\frac{d}{dv} \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz$, concerning the case $v < 0$. The desired probability density function is $f_{X-k'Y|X_{\text{disc}}=1, C_k=1}(v) = \begin{cases} \frac{2\pi}{\pi - \arctan k} \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz, & v \geq 0 \\ \frac{2\pi}{\pi - \arctan k} \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz, & v < 0 \end{cases}$

We now consider $u = 0$ and $j = 1$. We have to split again in two cases. For $v \geq 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 0, C_k = 1) = \int_0^{\infty} \int_{-\frac{v}{k'}}^0 \phi(w) \phi(z) dz dw + \int_0^{\infty} \phi(z) [\frac{1}{2} - \Phi(-kz)] dz - \frac{1}{2} F_{\text{SN}(0,1,-k)}(0) - \frac{1}{4}$

For $v < 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 0, C_k = 1) = \int_{-\frac{v}{k'}}^0 \int_{-\infty}^{v+kz} \phi(w) \phi(z) dz dw + \int_{-\frac{v}{k'}}^0 \int_{-\infty}^{-kz} \phi(w) \phi(z) dz dw - \int_{-\frac{v}{k'}}^0 \phi(z) [\Phi(v+k'z) - \Phi(-kz)] dz + \int_{-\frac{v}{k'}}^0 \phi(z) [\frac{1}{2} - \Phi(-kz)] dz dw - \int_{-\frac{v}{k'}}^0 \phi(z) \Phi(v+k'z) dz - \frac{1}{2} [F_{\text{SN}(0,1,-k)}(-\frac{v}{k'})]$

functions. In the case of $v \geq 0$, it is simply 0 since there is no dependency on v . As for $v < 0$, $\frac{d}{dv} [\int_{-\frac{v}{k'}}^{\infty} \phi(z) \Phi(v+k'z) dz + \frac{1}{2} F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k}) + \frac{1}{2} \Phi(\frac{v}{k'})] = \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz + \frac{1}{2} f_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k}) \frac{1}{k'+k} - \frac{1}{2} \phi(\frac{v}{k'}) \frac{1}{k'} - \frac{1}{2} f_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k}) \frac{1}{k'+k} + \frac{1}{2} \phi(\frac{v}{k'}) \frac{1}{k'}$

Once again, (A.4) was applied, in this case to $\int_{-\frac{v}{k'}}^{\infty} \phi(z) \Phi(v+k'z) dz$. The desired probability density function is $f_{X-k'Y|X_{\text{disc}}=0, C_k=1}(v) = \begin{cases} 0, & v \geq 0 \\ \frac{2\pi}{\pi - \arctan k} \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz, & v < 0 \end{cases}$

We now consider $u = 1$ and $j = 0$. For $v \geq 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 1, C_k = 0) = \int_{-\frac{v}{k'}}^{\infty} \int_0^{v+kz} \phi(w) \phi(z) dz dw + \int_{-\frac{v}{k'}}^{\infty} \int_0^0 \phi(w) \phi(z) dz dw = \int_{-\frac{v}{k'}}^{\infty} \phi(z) [\Phi(v+k'z) - \frac{1}{2}] dz + \int_{-\frac{v}{k'}}^{\infty} \phi(z) [\Phi(-kz) - \frac{1}{2}] dz dw - \int_{-\frac{v}{k'}}^{\infty} \phi(z) \Phi(v+k'z) dz - \frac{1}{2} [\Phi(-\frac{v}{k'+k}) - \Phi(-\frac{v}{k'})] + \frac{1}{2} [F_{\text{SN}(0,1,-k)}(0) - F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k})] - \frac{1}{2} \Phi(\frac{v}{k'+k}) - \int_{-\frac{v}{k'}}^{\infty} \phi(z) \Phi(v+k'z) dz - \frac{1}{2} \Phi(\frac{v}{k'+k}) + \frac{1}{2} F_{\text{SN}(0,1,-k)}(0) - \frac{1}{2} F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k})$

As for the case $v < 0$, $P(X - k'Y \leq v, X_{\text{disc}} = 1, C_k = 0) = 0$. We now need to take the derivative with respect to v of the expression obtained for $v \geq 0$ (the derivative of the one for $v < 0$ is 0) to obtain the corresponding conditional density functions. We have $\frac{d}{dv} [\int_{-\frac{v}{k'}}^{\infty} \phi(z) \Phi(v+k'z) dz - \frac{1}{2} \Phi(-\frac{v}{k'+k}) + \frac{1}{2} F_{\text{SN}(0,1,-k)}(0) - \frac{1}{2} F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k})] = \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz - \frac{1}{2} \phi(-\frac{v}{k'+k}) + \frac{1}{2} f_{\text{SN}(0,1,-k)}(0) - \frac{1}{2} f_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k})$

$= \int_{-\frac{v}{k'}}^{\infty} \phi(z) \phi(v+k'z) dz - \frac{1}{2} \phi(-\frac{v}{k'+k})$

$P(X - k'Y \leq v, X_{\text{disc}} = 1, C_k = 0) = \int_{-\frac{v}{k'+k}}^0 \int_{-kz}^{v+kz} \phi(w) \phi(z) dw dz = \int_{-\frac{v}{k'+k}}^0 \phi(z) [\Phi(v+k'z) - \frac{1}{2}] dz + \int_0^{\infty} \phi(z) [\Phi(v+k'z) - \frac{1}{2}] dz - \int_{-\frac{v}{k'+k}}^0 \phi(z) \Phi(v+k'z) dz - F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k}) + \int_{-\frac{v}{k'+k}}^{\infty} \frac{2\pi}{\pi - \arctan k} \zeta(z, v) dz$

$= \int_{-\frac{v}{k'+k}}^0 \int_{-kz}^{v+kz} \phi(w) \phi(z) dw dz = \int_{-\frac{v}{k'+k}}^0 \phi(z) [\Phi(v+k'z) - \frac{1}{2}] dz + \int_0^{\infty} \phi(z) [\Phi(v+k'z) - \frac{1}{2}] dz - \int_{-\frac{v}{k'+k}}^0 \phi(z) \Phi(v+k'z) dz - F_{\text{SN}(0,1,-k)}(-\frac{v}{k'+k}) + \int_{-\frac{v}{k'+k}}^{\infty} \frac{2\pi}{\pi - \arctan k} \zeta(z, v) dz$

$$g(x) = \frac{1}{2} \sqrt{x(x-a)(x-b) \pm \sqrt{b^2 - 4ac}}$$

$$a^2 - b^2 = (a-b)(a+b)$$

$$\sqrt{xy^2} = z$$

$$x^2 + y^2 + z^2 = \sqrt{xy^2} = z$$

$$x = \frac{b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$y = \frac{3a \pm \sqrt{b^2 - 4ac}}{2a}$$

$$x^2 - 3^2 - 4^2 = 0$$

$$x^2 - 3^2 = 4^2$$

$$x = \pm \sqrt{3^2 + 4^2} = \pm 5$$

$$x^3 + x = c$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$\sqrt{xy^2} = z$$

$$x^2 - 3^2 - 4^2 = 0$$

$$x = \pm \sqrt{3^2 + 4^2} = \pm 5$$

CONCLUSIONS

- Theoretical framework for the comparison of feature selection methods.
- Derivation of upper and lower bounds for the target objective functions.
- Distributional setting to highlight deficiencies of feature selection methods.
- Identification of feature selection methods to be avoided and preferred.



MIM, MIFS, mRMR, maxMIFS:
Ignore complementary

JMI, CMIM or even
better: DMIM

Main References



Neurocomputing 513 (2022) 215–232

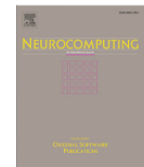


ELSEVIER

Contents lists available at [ScienceDirect](#)

Neurocomputing

journal homepage: www.elsevier.com/locate/neucom



Feature selection using Decomposed Mutual Information Maximization



Francisco Macedo^a, Rui Valadas^{b,c,*}, Eunice Carrasquinha^{a,b}, M. Rosário Oliveira^a, António Pacheco^a

Theoretical foundations of forward feature selection methods based on mutual information

Francisco Macedo^{a,b}, M. Rosário Oliveira^{a,*}, António Pacheco^a, Rui Valadas^c

^aCEMAT and Department of Mathematics, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, Lisboa 1049-001, Portugal

^bEPF Lausanne, SB-MATHICSE-ANCHP, Station 8, Lausanne CH-1015, Switzerland

^cIT and Department of Electrical and Computer Engineering, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, Lisboa 1049-001, Portugal

Theoretical evaluation of feature selection methods based on mutual information

Cláudia Pascoal^a, M. Rosário Oliveira^{a,*}, António Pacheco^a, Rui Valadas^b

^aCEMAT and Dep. Mathematics, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1049-001 Lisboa, Portugal

^bIT and Dep. Electrical and Computer Engineering, Instituto Superior Técnico, Universidade de Lisboa, Av. Rovisco Pais, 1049-001 Lisboa, Portugal