

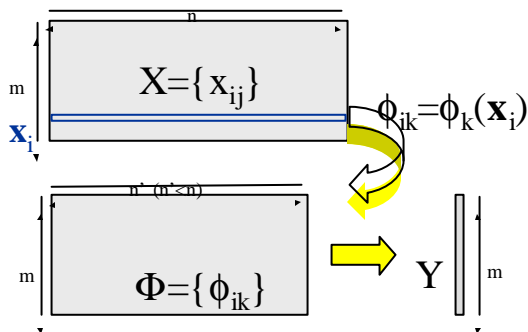
Lecture 11: Information Theoretic Methods

Isabelle Guyon
guyoni@inf.ethz.ch

Book Chapter 6 and
<http://www.jmlr.org/papers/volume3/torkkola03a/torkkola03a.pdf>

Mutual Information as “Information Gain”

Feature Transforms



Information, Entropy

- For an event happening with probability p , the quantity of information is $-\log p$.
- If the log is of base 2, the unit is a “bit”.
- The average quantity of information over all events is the entropy:

$$H = - \sum p \log p$$

Conditional Entropy

- Our case: Y =class labels, Φ feature representation.

- Entropy: $H(Y) = - \sum_y p(y) \log p(y)$

- Conditional entropy:

$$H(Y|\Phi) = - \int_{\mathbf{f}} p(\mathbf{f}) \sum_y p(y|\mathbf{f}) \log p(y|\mathbf{f}) d\mathbf{f}$$

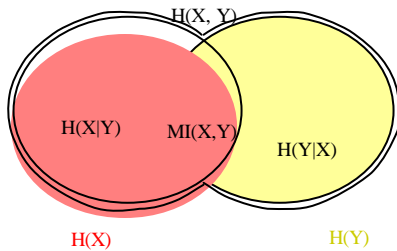
Mutual Information

- Information Gain: The amount of class uncertainty reduction by observing Φ is:

- $MI(Y, \Phi) = H(Y) - H(Y|\Phi)$

- $MI(Y, \Phi) = \sum_y \int_{\mathbf{f}} p(y, \mathbf{f}) \log \frac{p(y, \mathbf{f})}{p(y) p(\mathbf{f})} d\mathbf{f}$

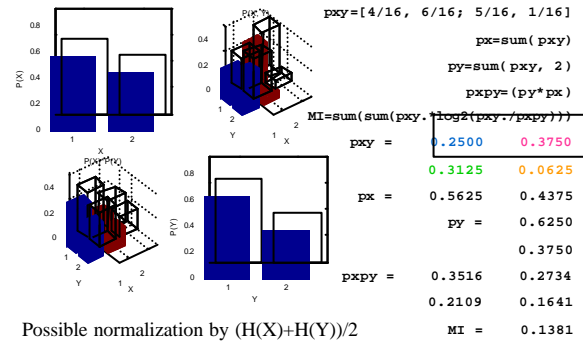
Mutual Information



$$H(X, Y) = H(X) + H(Y) - MI(X, Y)$$

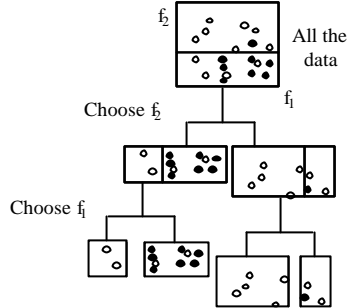
$$MI(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

MI Estimation



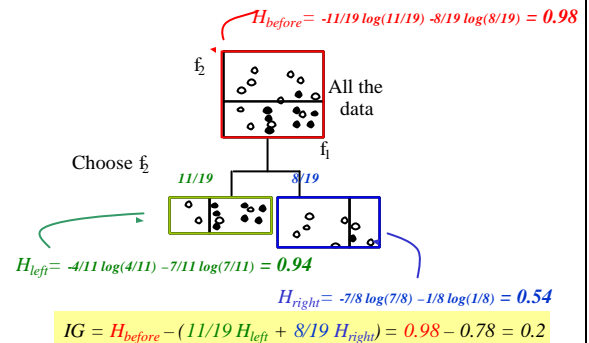
Tree Classifiers

CART (Breiman, 1984) or C4.5 (Quinlan, 1993)



At each step, choose the feature that "reduces entropy" most. Work towards "node purity".

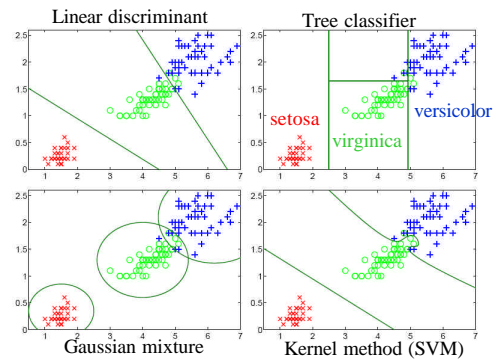
Information Gain



IG = MI

- $H_{before} = H(Y)$
- $H_{left} = H(Y|f_2=-1)$
- $H_{right} = H(Y|f_2=1)$
- $IG = H_{before} - (p(f_2=-1) H_{left} + p(f_2=1) H_{right})$
 $= H(Y) - (p(f_2=-1) H(Y|f_2=-1) + p(f_2=1) H(Y|f_2=1))$
 $= H(Y) - H(Y|f_2)$
 $= MI(Y, f_2)$

Iris Data (Fisher, 1936/Chapter 1)



MI for Feature Selection

- Random Forest (*Breiman, Cutler, 2002*)
- Feature ranking
- Forward selection
- Feature subset selection (Markov Blankets, scaling factors, shrinkage)

Forward Selection with MI

Fleuret, 2004. Practical only for binary features.

- Select a first feature $X_{2(1)}$ with maximum MI with the target.
- For each remaining feature X_i and each previously selected feature $X_{2(j)}$, compute the conditional mutual information:

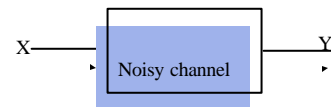


$$\text{CMI}(X_i, Y | X_{2(j)}) = \sum_{X_{2(j)}} P(X_{2(j)}) \text{MI}(X_i, Y | X_{2(j)})$$

- Select the feature with maximum CMI.

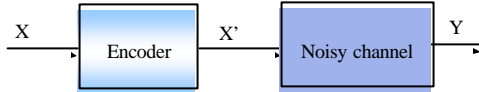
Transmission Channels

Transmission Channels



$$\begin{aligned} \text{Channel capacity} &= \max_{p(x)} \text{MI}(X, Y) \\ &= \text{(second Shannon theorem)} \\ &\quad \text{highest achievable information} \\ &\quad \text{transmission rate without error.} \end{aligned}$$

Coding



Code efficiency = length(code) / length(original signal)

Theorems (Shannon, 1948):

- Errorless transmission is achievable with sufficiently inefficient (well chosen) codes; the efficiency is bounded by the capacity.
- For a chosen error rate, one can find a shorter code.

MI and Error Rate

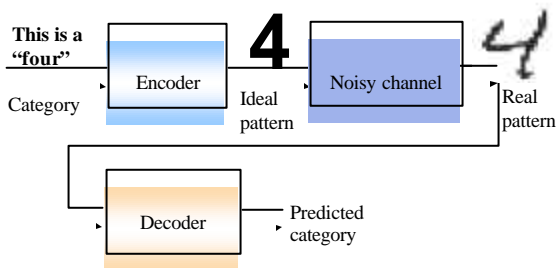
$$e_{\text{bayes}}(X) \geq 1 - \frac{I(Y, X) + \log 2}{\log(|Y|)}$$

Fano, 1961

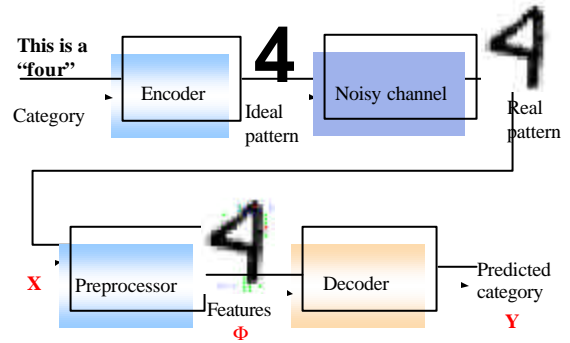
$$e_{\text{bayes}}(X) \leq \frac{1}{2}H(Y|X) = \frac{1}{2}(H(Y) - I(Y, X))$$

Hellman and Raviv (1970) Feder and Merhav (1994)

Noisy Channel in PR



Noisy Channel in PR



Information Bottleneck

Tishby et al., 1999



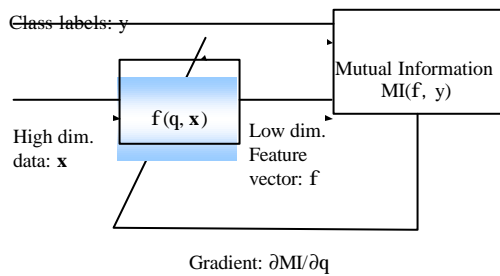
maximize $MI(X, Y | \Phi) = MI(\Phi, Y)$

but

minimize $MI(X, \Phi)$

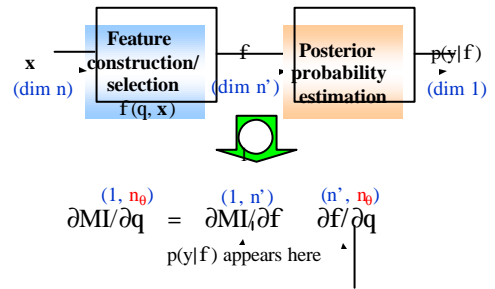
Information Bottleneck Methods

Gradient Descent



Torkkola, 2001

Computational Tricks (1/4)



For example $f(\mathbf{W}, \mathbf{x}) = \mathbf{W} \mathbf{x}^T$ (linear model) $\Rightarrow \partial f / \partial \mathbf{w} = \mathbf{x}$

Computational Tricks (2/4)

$$MI(Y, \Phi) = \sum_y \int_{\Phi} p(y, f) \log \frac{p(y, f)}{p(y) p(f)} df$$



$$MI(Y, \Phi) = \sum_y \int_{\Phi} (p(y, f) - p(y) p(f))^2 df$$

Justification: Renyi entropy $H(Y) = -\log \sum_y p(y)^2$

Computational Tricks (3/4)

$$G(\mathbf{y}, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \mathbf{y}^T \Sigma^{-1} \mathbf{y}\right)$$

$$\int_{\mathbf{y}} G(\mathbf{y} - \mathbf{a}_i, \Sigma_1) G(\mathbf{y} - \mathbf{a}_j, \Sigma_2) d\mathbf{y} = G(\mathbf{a}_i - \mathbf{a}_j, \Sigma_1 + \Sigma_2)$$



Use densities that are mixtures of Gaussians.

Examples of Suitable Densities

Gaussian mixture:

$$p(f|y) = G(f - f_y, \Sigma_y)$$

$$p(f) = \sum_y p(y) p(f|y) = \sum_y (m_y/m) G(f - f_y, \Sigma_y)$$

one per class

Parzen windows:

$$p(f|y) = (1/m_y) \sum_{c \in y} G(f - f_c, \sigma)$$

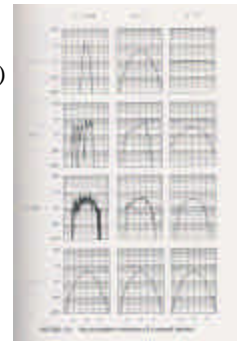
$$p(f|y) = (1/m) \sum_i G(f - f_i, \sigma)$$

the same for all
one per example

Parzen Windows (Rosenblatt 1956, Parzen 1962)

$$p(x) = (1/m) \sum_{i=1:m} G(f - f_i, \sigma_m)$$

$$\sigma_m = \sigma / \sqrt{m}$$



From Duda and Hart, 1974

Computational Tricks (4/4)

$$MI(Y, \Phi) = \sum_y \int_{\phi} (p(y, f) - p(y) p(f))^2 df$$



$$MI(Y, \Phi) = \sum_y \int_{\phi} p(y, f)^2 df \quad V_{in}$$

$$+ \sum_y \int_{\phi} p(y)^2 p(f)^2 df \quad V_{all}$$

$$- 2 \sum_y \int_{\phi} p(y, f) p(y) p(f) df \quad V_{bw}$$

Plugging In

$$MI(Y, \Phi) = V_{in} + V_{all} - V_{bw}$$

$$V_{in} = \text{function} (G(f_{ci} - f_{cj}, 2\sigma))$$

$$V_{all} = \text{function} (G(f_i - f_j, 2\sigma))$$

$$V_{bw} = \text{function} (G(f_{ci} - f_j, 2\sigma))$$

Information Forces

- Max MI by hill climbing (gradient descent):

$$q(t+1) = q(t) + \eta \frac{\partial MI}{\partial q}$$

$$= q(t) + \eta \frac{\partial MI}{\partial f} \frac{\partial f}{\partial q}$$

$$= q(t) + \eta [\partial V_{in}/\partial f + \partial V_{all}/\partial f - \partial V_{bw}/\partial f] \partial f / \partial q$$
- Force = - gradient (Energy) Energy \leftrightarrow -MI
- $\partial V_{in}/\partial f$ Cohesion within classes
- $\partial V_{all}/\partial f$ Overall cohesion
- $-\partial V_{bw}/\partial f$ Repulsion to other classes

Feature Construction

<http://members.cox.net/torkkola/mmi>, Torkkola, 2001

Example 1: Three classes in 3-dimensional space. Artificial data.

1.a $f = \text{linear}$, $p(y|f) = \text{Parzen windows}$

1.b $f = \text{RBF}$, $p(y|f) = \text{Parzen windows}$

Example 2: Three classes in 12-dimensional space. This is the [oil-pipeline data](#) from [NCRG](#) at Aston University.

2.a $f = \text{linear}$, $p(y|f) = \text{Parzen windows}$

2.b $f = \text{linear}$, $p(y|f) = \text{Gaussian mixture}$

Example 3: Six classes in 36-dimensional space. This is the Landsat satellite image database from [Statlog-project](#), $f = \text{linear}$, $p(y|f) = \text{Parzen windows}$

Application to Feature Selection

$\Phi_i = \sigma_i x_i$ scaling factors, $\sigma_i \geq 0$

$S^* = \operatorname{argmax}_S \operatorname{MI}(\Phi, Y)$

subject to: $\sum_i \sigma_i^2 = 1$ (shear)

$S^* = \operatorname{argmin}_S \operatorname{MI}(X, \Phi) - \beta \operatorname{MI}(\Phi, Y)$

Conclusion

- MI characterizes the **dependency between random** variables:
 - Feature ranking
 - Forward selection with CMI
- MI is also an **Information Gain**
 - Tree classifiers
 - Random Forests
- MI characterizes the **channel capacity** and the **rate of distortion**
 - Information bottleneck methods
 - $\min_q \operatorname{MI}(X, \Phi) - \beta \operatorname{MI}(\Phi, Y)$

Exercise Class

Madelon
How to review a paper?
Exercises on IT methods

Madelon (continued)

```
my_classifier=svc({'coef0=1', 'degree=0', 'gamma=1',  
'shrinkage=1'});  
my_model=chain({'probe(relief',{'p_num=2000',  
'pval_max=0'})', standardize, my_classifier})  
'f_max=20'
```

Earn 1 point with a result better than the baseline method ($\text{testBER} < 7.33$ or $[\text{testBER} = 7.33 \text{ and } \text{feat_num} < 20$ (4%)).

Earn a total of 2 points with a result $\text{testBER} < 6.79\%$ or $[\text{testBER} = 6.79\% \text{ and } \text{feat_num} < 17$ (3.4%)].

Homework

- Write a review of a paper in the “Feature Extraction” book (the paper you presented or will present or any paper/chapter you want).
- Questions:
 - What is the goal of a review?
 - What is the goal of publishing?

Review form

- A. Eliminating criteria** (a zero means rejection)
1. Scope
 2. Novelty
 3. Usefulness
 4. Sanity
- B. Other fundamental criteria** (a low grade implies major revisions)
1. Quantity
 2. Reproducibility
 3. Demonstration
 4. Comparison

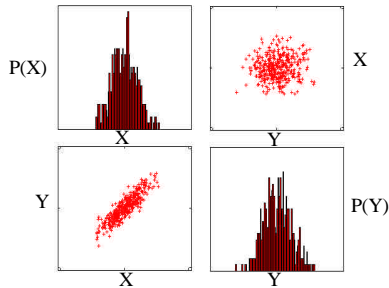
Review Form (continued)

- C. Forgivable problems of contents** (minor revisions may be required)
1. Completeness
 2. Take-aways
 3. Bibliography
 4. Outlook
- D. Bonus reproducibility criteria** (give additional credit)
1. Data availability
 2. Code availability

Review Form (end)

- E. Message delivery** (major/minor revisions may be required)
1. Readability
 2. Notations
 3. Figures
 4. Formalism
 5. Density
 6. Language

MI and Pearson Correlation



$$MI(X, Y) = -\frac{1}{2} \log(1-R^2)$$

Exercise: prove $MI = -\frac{1}{2} \log(1-R^2)$

Prove:

- Entropy Gaussian source: $H(X) = \frac{1}{2} \log(2\pi e) \sigma^2$
- Solution of one-variable least-square regression $y=ax$: $a = \sigma_{xy} / \sigma_x^2$. We will call a the sol. of $x=a'y$.
- Pearson correlation coefficient: $R^2 = aa'$
- Residual error $\sigma^2 = E(y-ax)^2 = (1-R^2)\sigma_y^2$.
- $Y = \text{function}(X) + N$, N is random noise
 - $MI(X, Y) = H(Y) - H(N)$
 - $MI(X, Y) \geq H(X) - H(N)$
- For a Gaussian channel with Gaussian source: $MI(X, Y) = -\frac{1}{2} \log \sigma^2 / \sigma_y^2$. Conclude the proof.

Markov Blankets

- Notations: All features X , Features selected Φ , complement Φ' , $X = [\Phi, \Phi']$.
- Definition: A **Markov Blanket** is a subset of features Φ such that:

$$P(X, Y | \Phi) = P(X | \Phi) P(Y | \Phi)$$

We write $X \perp Y | \Phi$ (X is independent of Y given Φ)

Exercise (Markov Blankets)

- If Φ is a MB, prove $P(Y | \Phi) = P(Y | X)$ for all assignments of values to Φ' .
- We define $DMI(\Phi) = E_{p(X)} D_{KL}(P(Y | X) || P(Y | \Phi))$
 Prove: $DMI(\Phi) = E_{p(X, Y)} P(Y | X) / P(Y | \Phi)$
- We define $A = \log[P(X, Y) / (P(X)P(Y))]$ and $B = \log[P(\Phi, Y) / (P(\Phi)P(Y))]$, prove:
 - $A - B = P(Y | X) / P(Y | \Phi)$
 - $\sum_{\{\Phi, Y\}} P(\Phi, Y) B = \sum_{\{X, Y\}} P(X, Y) A$
 - $DMI(\Phi) = MI(X, Y) - MI(\Phi, Y)$
- Define an optimization problem to find a minimum MB. How can you relate it to the information bottleneck approach?

Complement of exercise

- What are the similarities and differences between the “bottleneck” IT methods and the Bayesian/regularized risk methods?
- Can you suggest hybrid approaches?
- Can you suggest forward selection and backward elimination IT methods using the scaling factor approach?