# Information theoretic feature selection in multi-label data through composite likelihood

Konstantinos Sechidis, <u>Nikolaos Nikolaou</u>, and Gavin Brown

School of Computer Science
University of Manchester

- Multi-label: Each datapoint can be associated to $> 1$ labels

- Multi-label: Each datapoint can be associated to $> 1$ labels
- Applications

- Multi-label: Each datapoint can be associated to $> 1$ labels
- Applications
  - ▶ Bioinformatics: 1 gene/protein, many functions

- Multi-label: Each datapoint can be associated to $> 1$ labels
- Applications
  - Bioinformatics: 1 gene/protein, many functions
  - Text Mining: 1 webpage/document, many categories

# Motivation: Multi-label Learning

- Multi-label: Each datapoint can be associated to $> 1$ labels
- Applications
  - Bioinformatics: 1 gene/protein, many functions
  - Text Mining: 1 webpage/document, many categories
  - Image Retrieval: 1 image, many semantic concepts



Male, Person,
Motorbike, Vehicle
Building

Female, Person,
Building

Male, Person

Rabbit, Animal
Car, Vehicle

- Multi-label: Each datapoint can be associated to $> 1$ labels
- Applications
  - ▶ Bioinformatics: 1 gene/protein, many functions
  - ▶ Text Mining: 1 webpage/document, many categories
  - ▶ Image Retrieval: 1 image, many semantic concepts



Male, Person,
Motorbike, Vehicle
Building

Female, Person,
Building

Male, Person

Rabbit, Animal
Car, Vehicle

- Common characteristic of these domains: Large number of features

# Feature Selection

- Feature Selection: Find minimal subset of features with maximal useful information
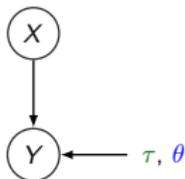
# Feature Selection

- Feature Selection: Find minimal subset of features with maximal useful information

- Filters: Functions that assign a "utility" score to each feature

# Feature Selection

- Feature Selection: Find minimal subset of features with maximal useful information

- Filters: Functions that assign a "utility" score to each feature

- In this work we discuss information-theoretic filters

- Feature Selection: Find minimal subset of features with maximal useful information

- Filters: Functions that assign a "utility" score to each feature

- In this work we discuss information-theoretic filters

- Filter Assumption: model and feature selection are independent

- Brown et al. (JMLR 2012) unified many heuristic information-theoretic filter criteria for feature selection

  Conditional Likelihood Maximization under model

# Feature Selection via Likelihood Maximization

- Brown et al. (JMLR 2012) unified many heuristic information-theoretic filter criteria for feature selection

  Conditional Likelihood Maximization under model



- Negative log-likelihood asymptotically decomposes into 3 terms:
  $$\lim_{N\to\infty} -\ell = \text{model term} + \text{feature selection term} + \text{Bayes error}$$

  feature selection is mutual info $I(X_\theta; Y)$

Feature space independence assumptions:
Full:

| Features | Label |
|----------|-------|

$$X_1$$

$$X_2$$

$$X_3 \qquad Y$$

$$\vdots$$

$$X_d$$

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

| Features | Label |
|----------|-------|

$$X_1$$

$$X_2$$

$$X_3$$

$$\vdots$$

$$X_d$$

$$Y$$

Feature space independence assumptions:
Full:



Features    Label

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$    $Y$

$J_{MIM}(X_k) = I(X_k; Y)$

Feature space independence assumptions:

Full:

| Features | Label |
|---|---|

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y$

$J_{MIM}(X_k) = I(X_k; Y)$

Partial
(i.e. pairwise dependencies):

| Features | Label |
|---|---|

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y$

# Single-label Feature Selection Criteria

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}(X_k) = I(X_k; Y)$$

# Single-label Feature Selection Criteria

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$J_{MIM}(X_k) = I(X_k; Y)$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$J_{MIM}(X_k) = I(X_k; Y)$

# Single-label Feature Selection Criteria

Feature space independence assumptions:

Full:

Partial
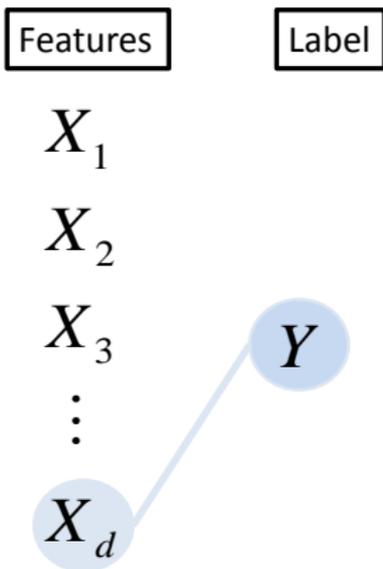(i.e. pairwise dependencies):



$$J_{MIM}(X_k) = I(X_k; Y)$$

$$J_{JMI}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_{\theta_j} X_k; Y)$$

# Single-label Feature Selection Criteria

Feature space independence assumptions:
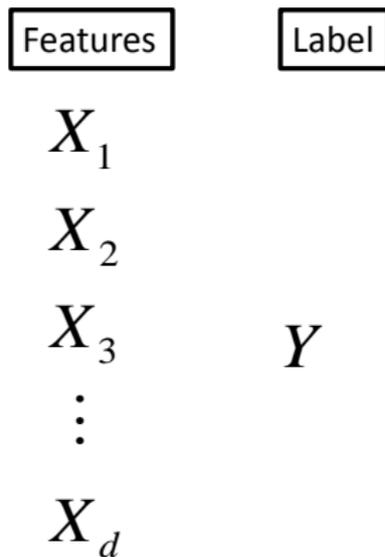
Full:

| Features | | Label |

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y$

$$J_{MIM}(X_k) = I(X_k; Y)$$

Partial (i.e. pairwise dependencies):

| Features | | Label |

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y$

$$J_{JMI}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_{\theta_j} X_k; Y)$$
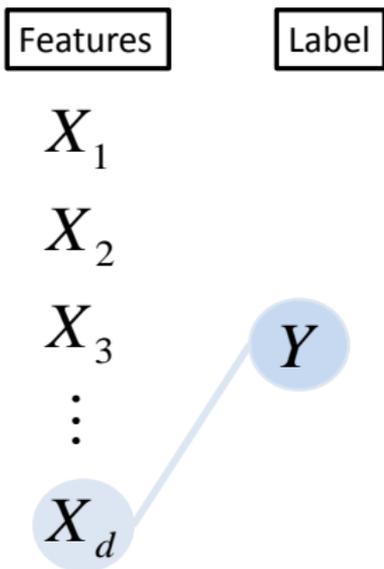
None:

| Features | | Label |

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y$

# Single-label Feature Selection Criteria

Feature space independence assumptions:



Full:

Partial (i.e. pairwise dependencies):

None:

$$J_{MIM}(X_k) = I(X_k; Y)$$

$$J_{JMI}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_{\theta_j} X_k; Y)$$

# Single-label Feature Selection Criteria

Feature space independence assumptions:
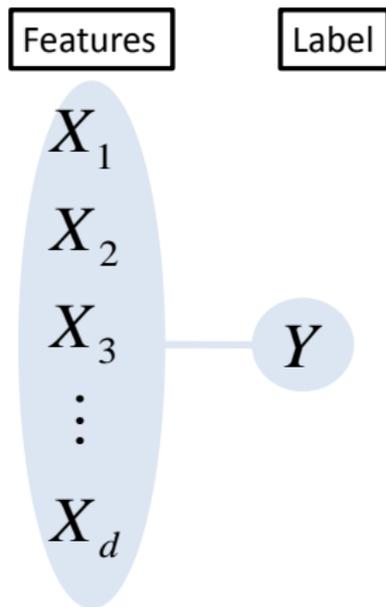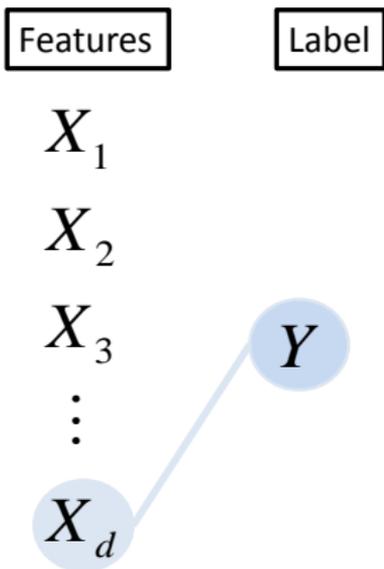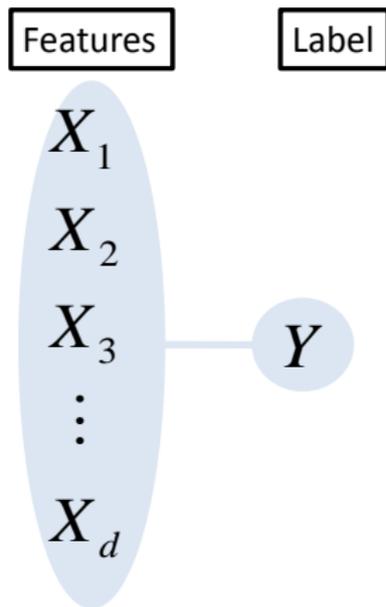


| Full: | Partial (i.e. pairwise dependencies): | None: |
|---|---|---|

$$J_{MIM}(X_k) = I(X_k; Y)$$

$$J_{JMI}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_{\theta_j} X_k; Y)$$

$$J_{CMI}(X_k) = I(X_k; Y | X_\theta)$$

- Next, extend to multi-label where $Y$ is $q$-dimensional



Man, Hat,
Person

- Next, extend to multi-label where $Y$ is $q$-dimensional



Man, Hat,
Person

- What independence assumptions can we make in label space?

# Extending Framework to Multi-label Setting

- Next, extend to multi-label where $Y$ is $q$-dimensional



Man, Hat,
Person

- What independence assumptions can we make in label space?
- In this work we examined:
  - ▶ Binary Relevance (BR) vs Label Powerset (LP)

- Label Powerset (LP): No independence among labels

# Multi-label Extension: LP Transformation

- Label Powerset (LP): No independence among labels
- Binary $q$-label problem $\Rightarrow$ 1 single-label, $2^q$-class problem

- Label Powerset (LP): No independence among labels
- Binary $q$-label problem $\Rightarrow$ 1 single-label, $2^q$-class problem

|  | Animal | Building | Vehicle |  | $y$ |
|---|---|---|---|---|---|
|  | 1 | 0 | 1 | $\Rightarrow$ | 1 0 1 |
|  | 0 | 1 | 1 | $\Rightarrow$ | 0 1 1 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

# Multi-label Extension: LP Transformation

- Label Powerset (LP): No independence among labels
- Binary $q$-label problem $\Rightarrow$ 1 single-label, $2^q$-class problem



|  | Animal | Building | Vehicle |  | $y$ |
|---|---|---|---|---|---|
|  | 1 | 0 | 1 | $\Rightarrow$ | 1 0 1 |
|  | 0 | 1 | 1 | $\Rightarrow$ | 0 1 1 |
|  | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

- Pros: dependencies among labels are accounted for

# Multi-label Extension: LP Transformation

- Label Powerset (LP): No independence among labels
- Binary $q$-label problem $\Rightarrow$ 1 single-label, $2^q$-class problem

|  | Animal | Building | Vehicle |  | $y$ |
|---|---|---|---|---|---|
|  | 1 | 0 | 1 | $\Rightarrow$ | 1 0 1 |
|  | 0 | 1 | 1 | $\Rightarrow$ | 0 1 1 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

- Pros: dependencies among labels are accounted for
- Cons: probability estimates unreliable (curse of dimensionality)

Feature space independence assumptions:
Full:

| Features | Labels |
|---|---|

$X_1$

$X_2$

$X_3$

$\vdots$

$X_d$

$Y_1$

$Y_2$

$\vdots$

$Y_q$

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

# Multi-label Extension: LP Transformation
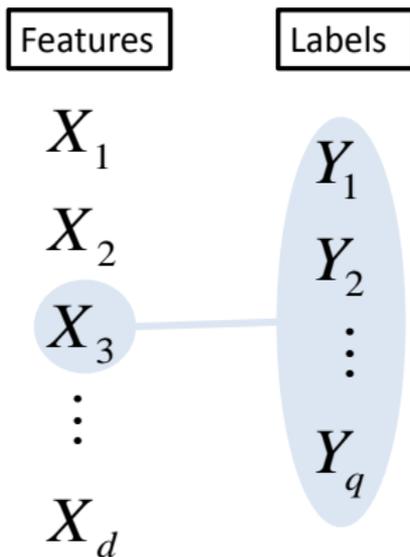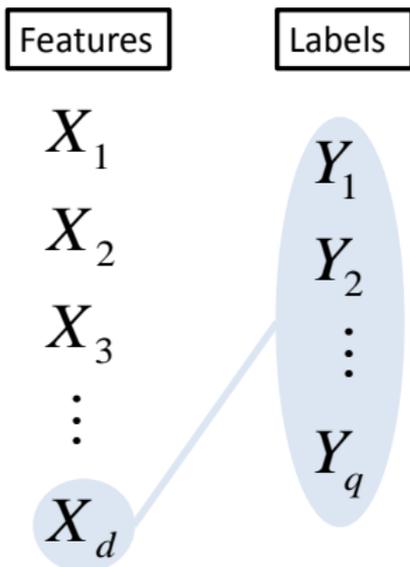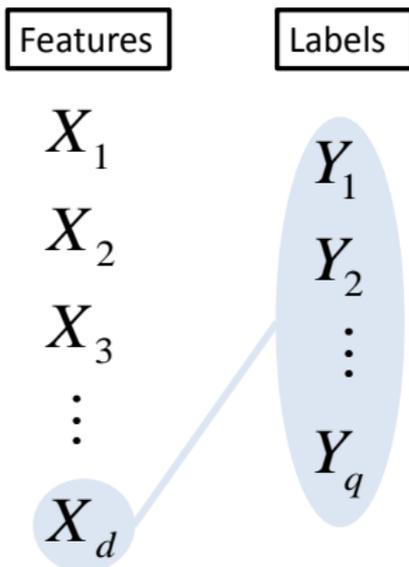
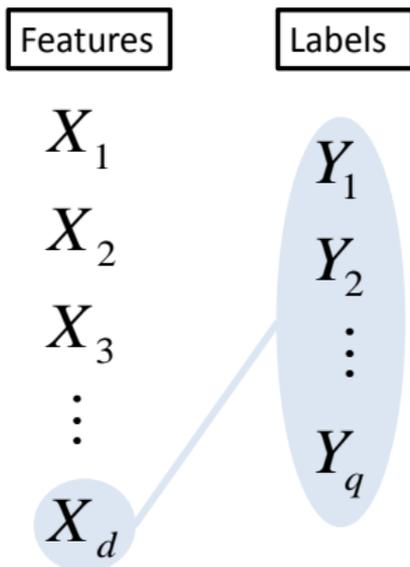Feature space independence assumptions:
Full:



$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$
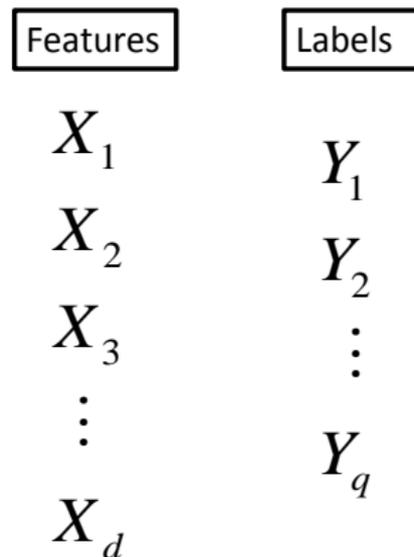
# Multi-label Extension: LP Transformation

Feature space independence assumptions:

Full:

Partial
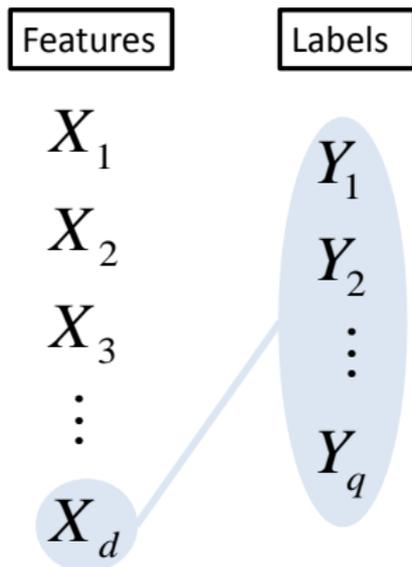(i.e. pairwise dependencies):



$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$
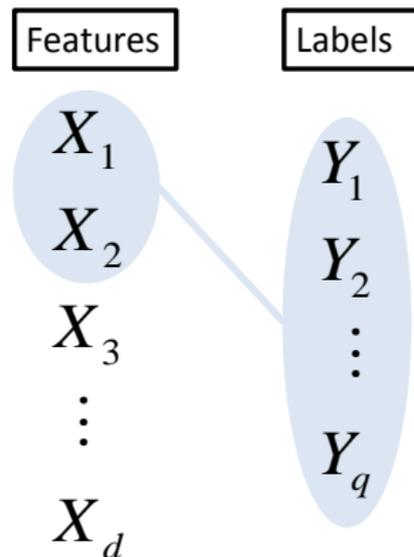
# Multi-label Extension: LP Transformation

Feature space independence assumptions:

Full:

Partial (i.e. pairwise dependencies):
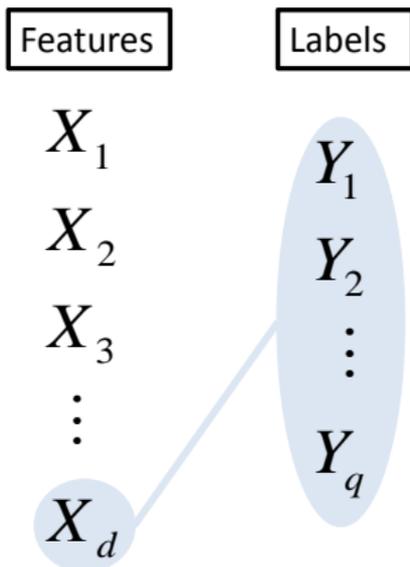


$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

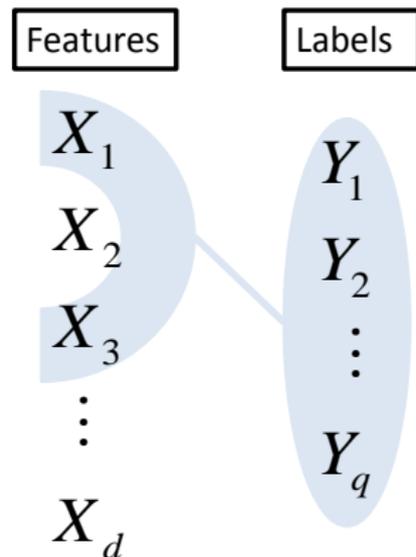# Multi-label Extension: LP Transformation

Feature space independence assumptions:

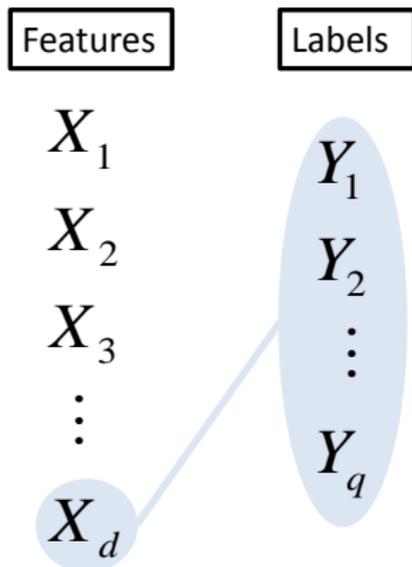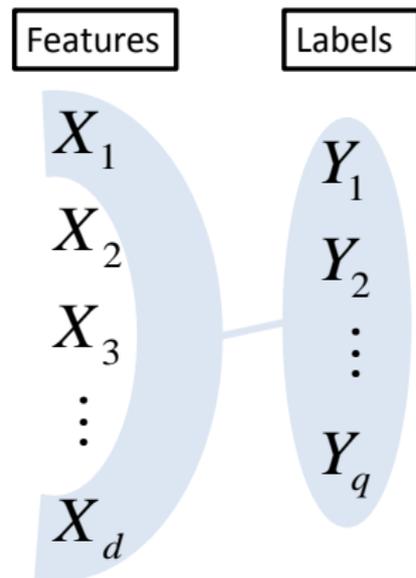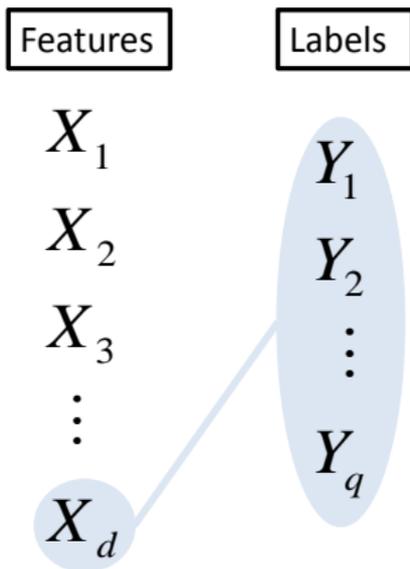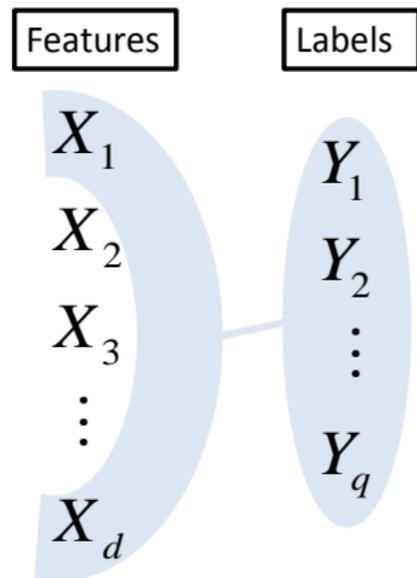| Full: | Partial (i.e. pairwise dependencies): |
|---|---|



$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

$$J_{JMI}^{LP}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_k X_{\theta_j}; Y_{1:q})$$

Feature space independence assumptions:



Full:

Partial (i.e. pairwise dependencies):

None:

$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

$$J_{JMI}^{LP}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_k X_{\theta_j}; Y_{1:q})$$

Feature space independence assumptions:



Full:

Partial
(i.e. pairwise dependencies):

None:

| Features | Labels |
| Features | Labels |
| Features | Labels |

$X_1$   $Y_1$

$X_2$   $Y_2$

$X_3$   $\vdots$

$\vdots$

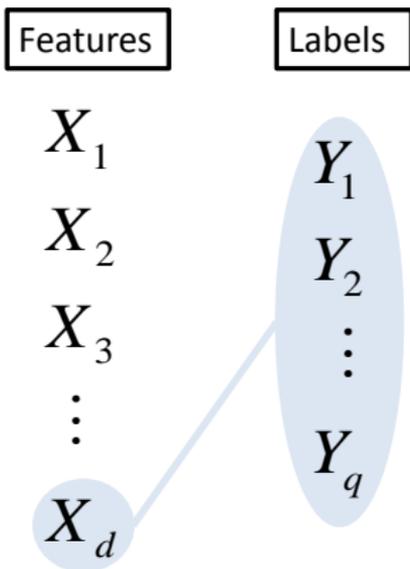$X_d$   $Y_q$

$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

$$J_{JMI}^{LP}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_k X_{\theta_j}; Y_{1:q})$$

# Multi-label Extension: LP Transformation

Feature space independence assumptions:



Full:

$$J_{MIM}^{LP}(X_k) = I(X_k; Y_{1:q})$$

Partial (i.e. pairwise dependencies):

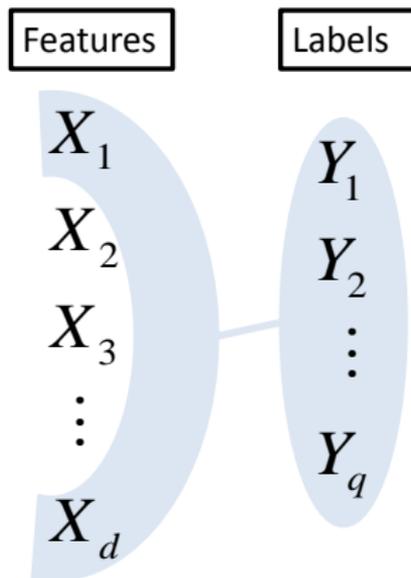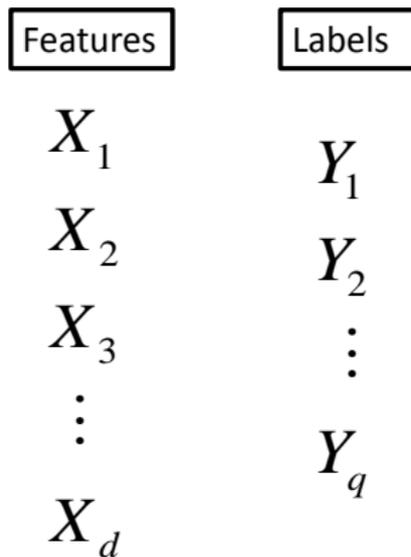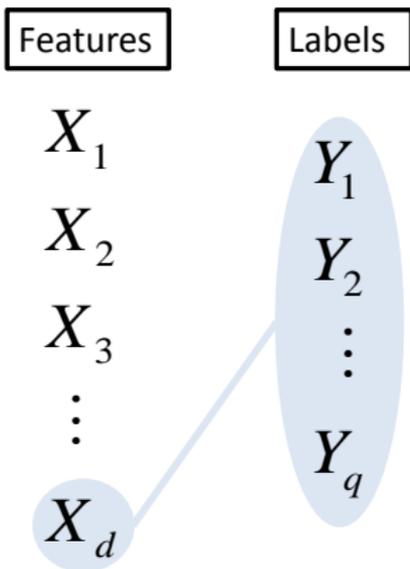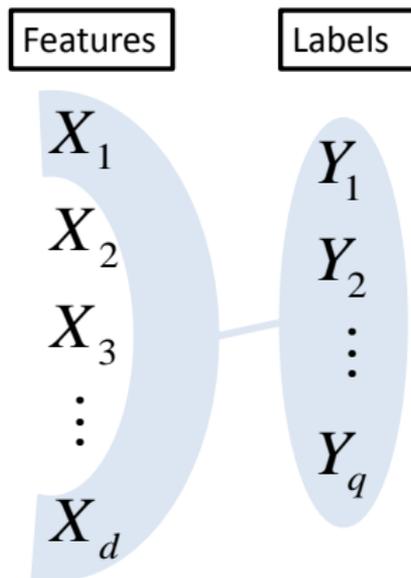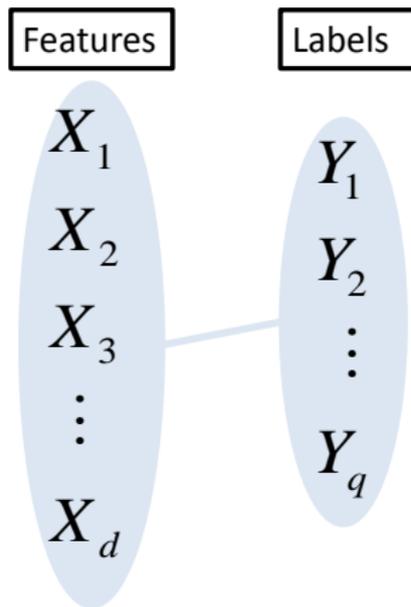$$J_{JMI}^{LP}(X_k) = \sum_{j=1}^{|X_\theta|} I(X_k X_{\theta_j}; Y_{1:q})$$

None:

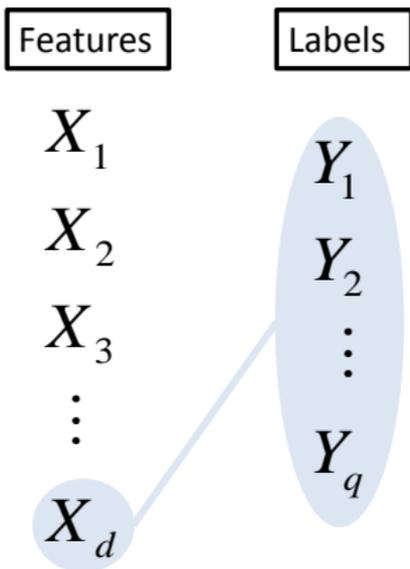$$J_{CMI}^{LP}(X_k) = I(X_k; Y_{1:q}|X_\theta)$$
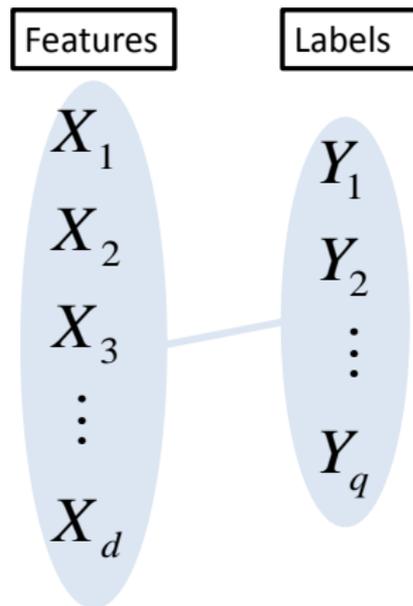
- Binary Relevance (BR): Full Independence among labels

- Binary Relevance (BR): Full Independence among labels
- Binary $q$-label problem $\Rightarrow$ $q$ independent single-label, binary problems

# Multi-label Extension: BR Transformation

- Binary Relevance (BR): Full Independence among labels
- Binary $q$-label problem $\Rightarrow$ $q$ independent single-label, binary problems

| | Animal | Building | Vehicle | | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|---|---|---|---|
|  | 1 | 0 | 1 | $\Rightarrow$ | 1 | 0 | 1 |
|  | 0 | 1 | 1 | $\Rightarrow$ | 0 | 1 | 1 |
| | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ |

# Multi-label Extension: BR Transformation

- Binary Relevance (BR): Full Independence among labels
- Binary $q$-label problem $\Rightarrow q$ independent single-label, binary problems



| | Animal | Building | Vehicle | | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|---|---|---|---|
| | 1 | 0 | 1 | $\Rightarrow$ | 1 | 0 | 1 |
| | 0 | 1 | 1 | $\Rightarrow$ | 0 | 1 | 1 |
| | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ | $\vdots$ |

- Pros: more reliable probability estimates

# Multi-label Extension: BR Transformation

- Binary Relevance (BR): Full Independence among labels
- Binary $q$-label problem $\Rightarrow$ $q$ independent single-label, binary problems



| | Animal | Building | Vehicle | | $y_1$ | $y_2$ | $y_3$ |
|---|---|---|---|---|---|---|---|
| | 1 | 0 | 1 | $\Rightarrow$ | 1 | 0 | 1 |
| | 0 | 1 | 1 | $\Rightarrow$ | 0 | 1 | 1 |
| $\vdots$ | $\vdots$ | $\vdots$ | $\vdots$ | | $\vdots$ | $\vdots$ | $\vdots$ |

- Pros: more reliable probability estimates
- Cons: dependencies among labels are not accounted for

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

Feature space independence assumptions:
Full:

# Multi-label Extension: BR Transformation

Feature space independence assumptions:
Full:

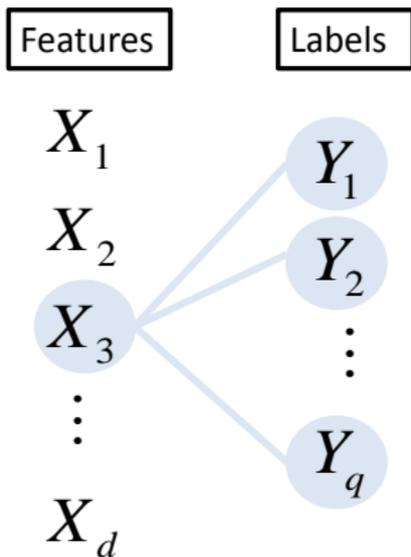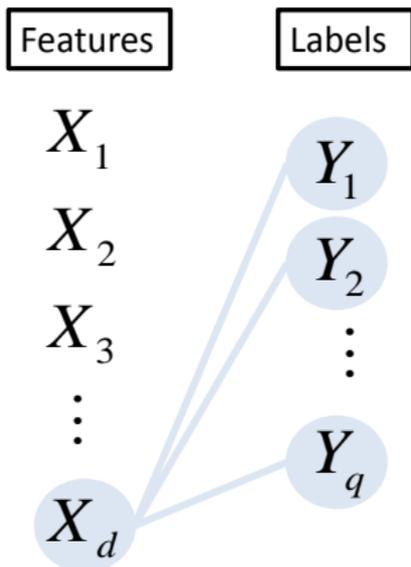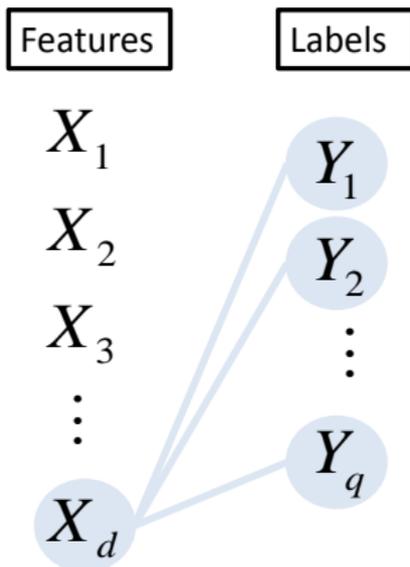

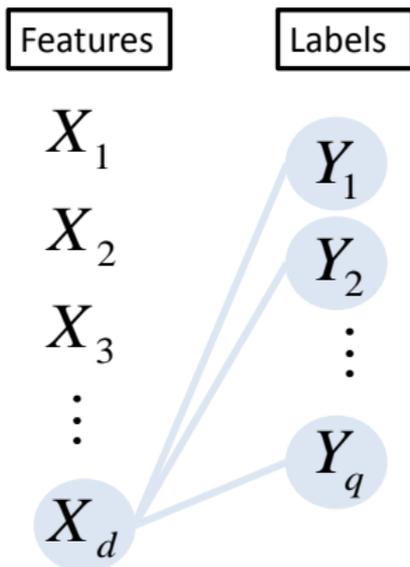$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

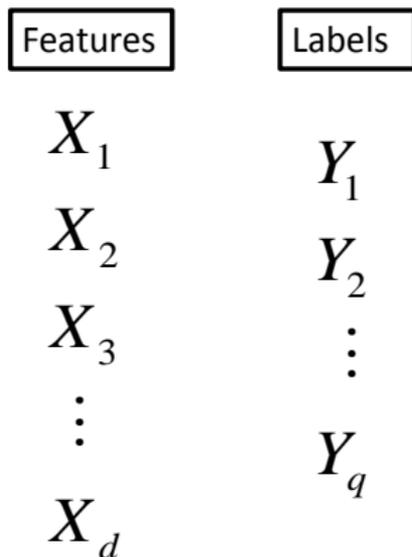Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

Feature space independence assumptions:

Full:

Partial
(i.e. pairwise dependencies):



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$
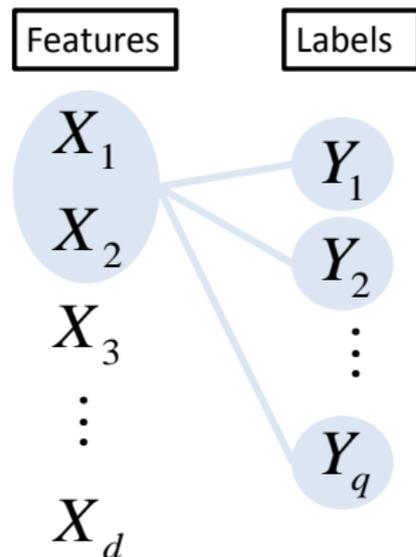
$$J_{JMI}^{BR}(X_k) = \sum_{j=1}^{|X_\theta|} \sum_{l=1}^{q} I(X_k X_{\theta_j}; Y_l)$$

# Multi-label Extension: BR Transformation

Feature space independence assumptions:

**Full:**

| Features | Labels |

$X_1$
$X_2$
$X_3$
$\vdots$
$X_d$

$Y_1$
$Y_2$
$\vdots$
$Y_q$

$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

**Partial**
(i.e. pairwise dependencies):

| Features | Labels |

$X_1$
$X_2$
$X_3$
$\vdots$
$X_d$

$Y_1$
$Y_2$
$\vdots$
$Y_q$

$$J_{JMI}^{BR}(X_k) = \sum_{j=1}^{|X_\theta|} \sum_{l=1}^{q} I(X_k X_{\theta_j}; Y_l)$$

**None:**

| Features | Labels |

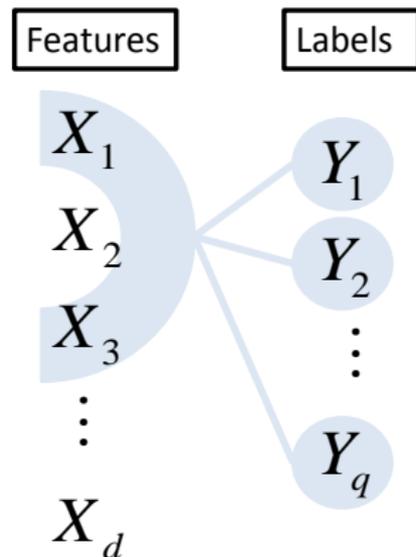$X_1$
$X_2$
$X_3$
$\vdots$
$X_d$

$Y_1$
$Y_2$
$\vdots$
$Y_q$

Feature space independence assumptions:



Full:

Partial
(i.e. pairwise dependencies):
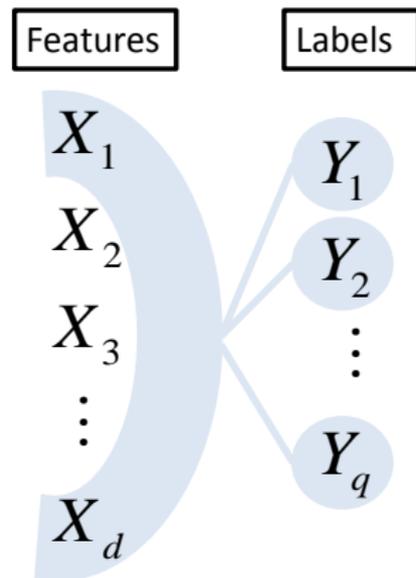
None:

$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$
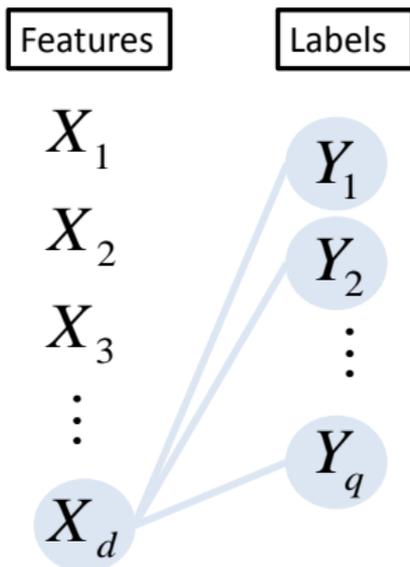
$$J_{JMI}^{BR}(X_k) = \sum_{j=1}^{|X_\theta|} \sum_{l=1}^{q} I(X_k X_{\theta_j}; Y_l)$$

# Multi-label Extension: BR Transformation
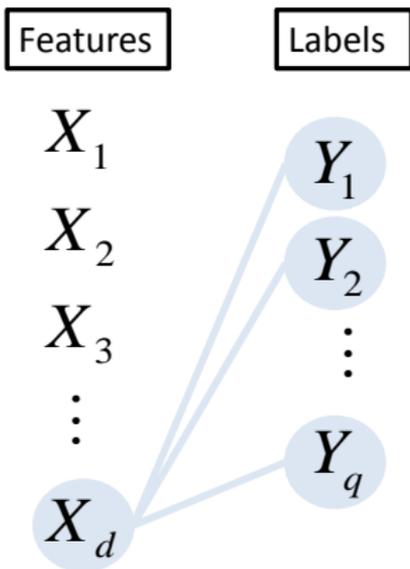
Feature space independence assumptions:



$$J_{MIM}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l)$$

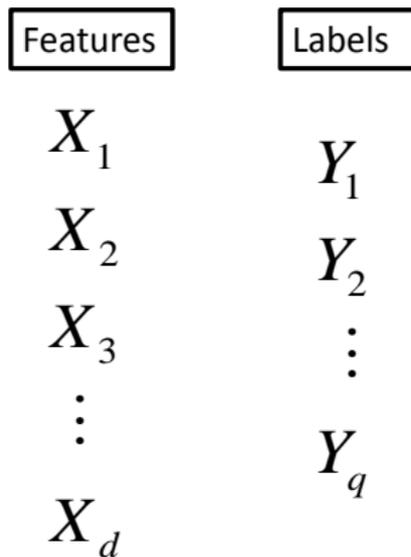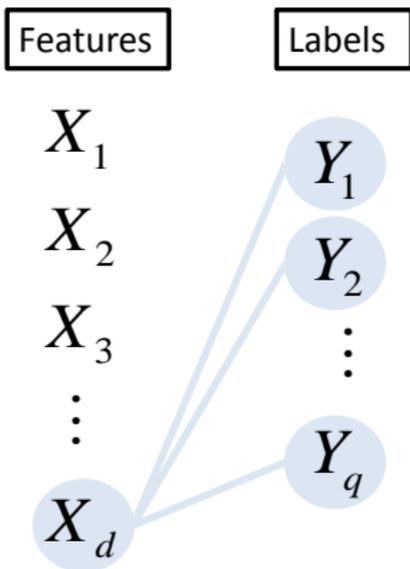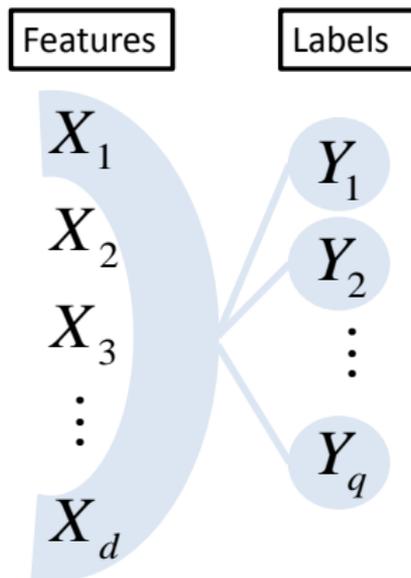$$J_{JMI}^{BR}(X_k) = \sum_{j=1}^{|X_\theta|} \sum_{l=1}^{q} I(X_k X_{\theta_j}; Y_l)$$

$$J_{CMI}^{BR}(X_k) = \sum_{l=1}^{q} I(X_k; Y_l | X_\theta)$$

- Summarizing, the criteria based on feature space $X$ and label space $Y$ independence assumptions:

|  | | Feature space independence assumptions | | |
|---|---|---|---|---|
|  | | *CMI* (none) | *JMI* (partial) | *MIM* (full) |
| Label space independence assumptions | Label Powerset (none) | $J_{\text{X:none}}^{\text{Y:none}}$ | $J_{\text{X:partial}}^{\text{Y:none}}$ | $J_{\text{X:full}}^{\text{Y:none}}$ |
|  | Binary Relevance (full) | $J_{\text{X:none}}^{\text{Y:full}}$ | $J_{\text{X:partial}}^{\text{Y:full}}$ | $J_{\text{X:full}}^{\text{Y:full}}$ |

- Summarizing, the criteria based on feature space $X$ and label space $Y$ independence assumptions:

| | | Feature space independence assumptions | | |
|---|---|---|---|---|
| | | CMI (none) | JMI (partial) | MIM (full) |
| Label space | Label Powerset (none) | *Doquire & Verleysen (2013)* | $J_{\text{X:partial}}^{\text{Y:none}}$ | $J_{\text{X:full}}^{\text{Y:none}}$ |
| independence | | | | |
| assumptions | Binary Relevance (full) | $J_{\text{X:none}}^{\text{Y:full}}$ | $J_{\text{X:partial}}^{\text{Y:full}}$ | $J_{\text{X:full}}^{\text{Y:full}}$ |

- Summarizing, the criteria based on feature space $X$ and label space $Y$ independence assumptions:

|  | Feature space independence assumptions | | |
|---|---|---|---|
| Label space | *CMI* (none) | *JMI* (partial) | *MIM* (full) |
| independence | Label Powerset (none) | *Doquire & Verleysen (2013)* | $J_{\mathrm{X:partial}}^{\mathrm{Y:none}}$ | *Spolaôr et al. (2013)* |
| assumptions | Binary Relevance (full) | $J_{\mathrm{X:none}}^{\mathrm{Y:full}}$ | $J_{\mathrm{X:partial}}^{\mathrm{Y:full}}$ | $J_{\mathrm{X:full}}^{\mathrm{Y:full}}$ |

- Summarizing, the criteria based on feature space $X$ and label space $Y$ independence assumptions:

| Label space | | Feature space independence assumptions | | |
| --- | --- | --- | --- | --- |
| | | *CMI* (none) | *JMI* (partial) | *MIM* (full) |
| independence | Label Powerset (none) | *Doquire & Verleysen (2013)* | $J_{\text{X:partial}}^{\text{Y:none}}$ | *Spolaôr et al. (2013)* |
| assumptions | Binary Relevance (full) | $J_{\text{X:none}}^{\text{Y:full}}$ | $J_{\text{X:partial}}^{\text{Y:full}}$ | *Young & Pedersen (1997), Trohidis et al. (2008), ...* |

- Compare

- Compare
  - effect of label space assumptions

- Compare
    - effect of label space assumptions
    - effect of feature space assumptions

- Compare
  - ▶ effect of label space assumptions
  - ▶ effect of feature space assumptions
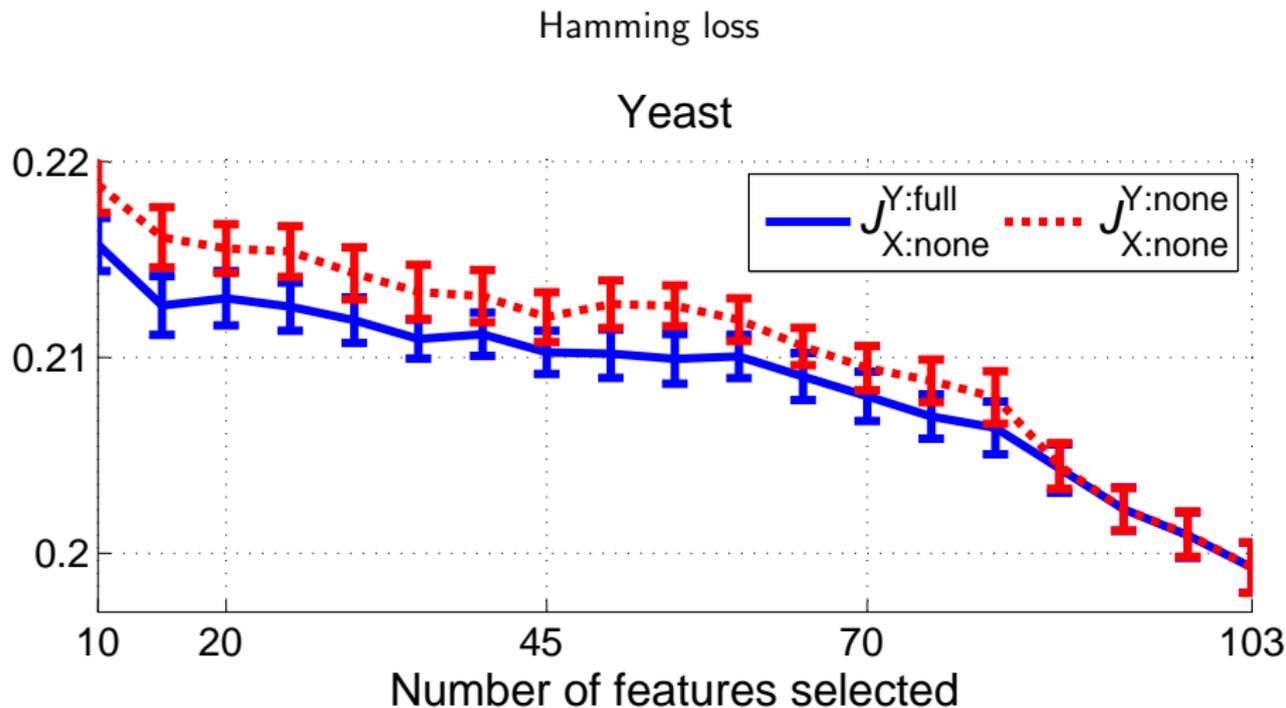  - ▶ our best criterion vs. state-of-the-art

- Compare
  - effect of label space assumptions
  - effect of feature space assumptions
  - our best criterion vs. state-of-the-art
- Procedure: Select $M$ top features under each criterion, classify, evaluate; vary $M$

- Compare
  - effect of label space assumptions
  - effect of feature space assumptions
  - our best criterion vs. state-of-the-art
- Procedure: Select $M$ top features under each criterion, classify, evaluate; vary $M$
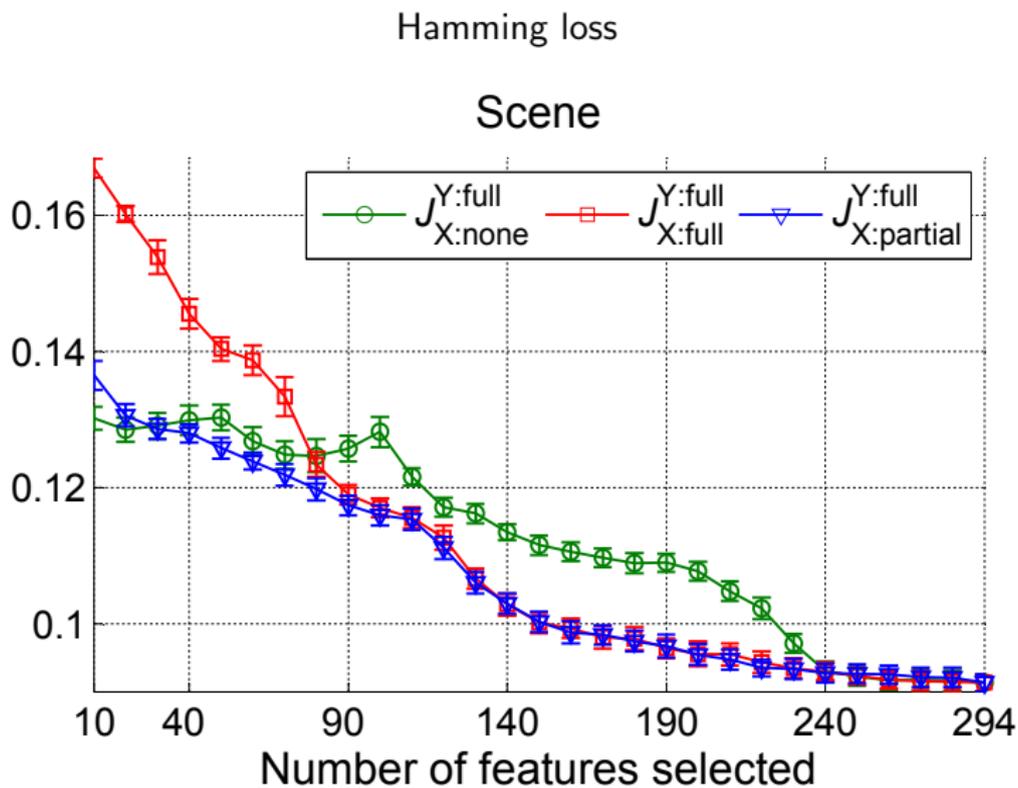- Datasets: scene and yeast

- Compare
  - effect of label space assumptions
  - effect of feature space assumptions
  - our best criterion vs. state-of-the-art
- Procedure: Select $M$ top features under each criterion, classify, evaluate; vary $M$
- Datasets: scene and yeast
- Classification: ML-$k$NN, $k = 7$

- Compare
  - effect of label space assumptions
  - effect of feature space assumptions
  - our best criterion vs. state-of-the-art
- Procedure: Select $M$ top features under each criterion, classify, evaluate; vary $M$
- Datasets: scene and yeast
- Classification: ML-$k$NN, $k = 7$
- Evaluation: Hamming Loss (shown here) and Ranking Loss (similar)

Hamming loss

Yeast

Hamming loss
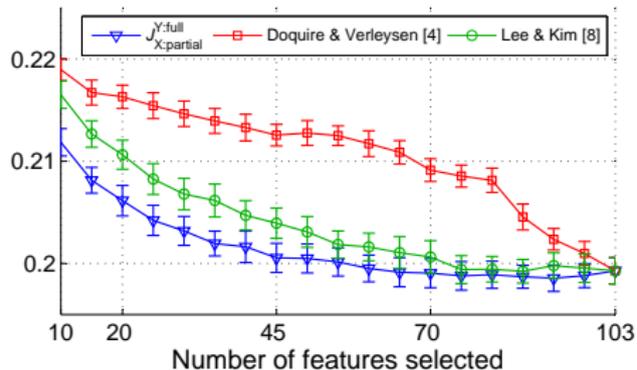
Scene

Hamming loss

Hamming loss



Yeast

Scene

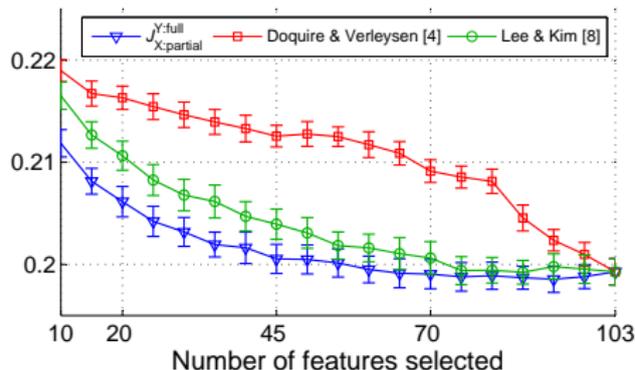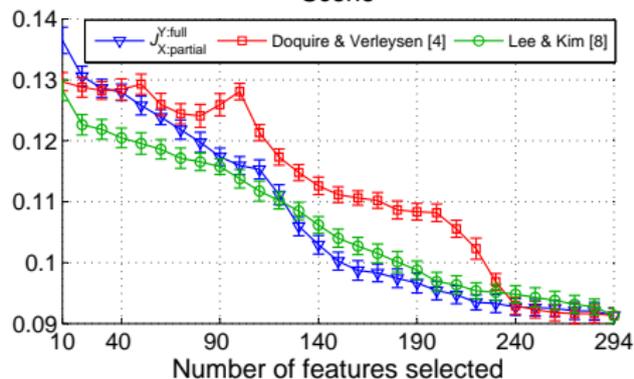- Doquire & Verleysen (2013) : $J^{Y:none}_{X:none}$ with pruning of rare cases

Hamming loss

Yeast — Scene

- Doquire & Verleysen (2013) : $J_{\text{X:none}}^{\text{Y:none}}$ with pruning of rare cases
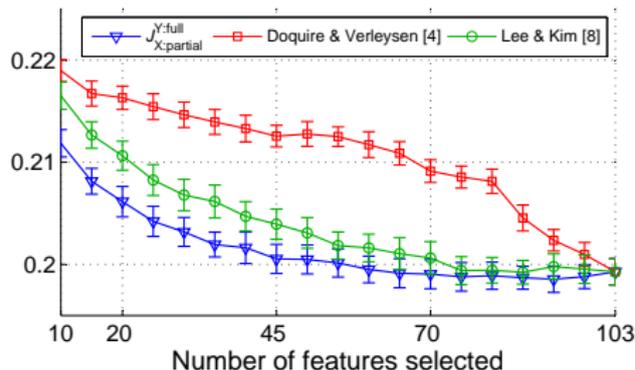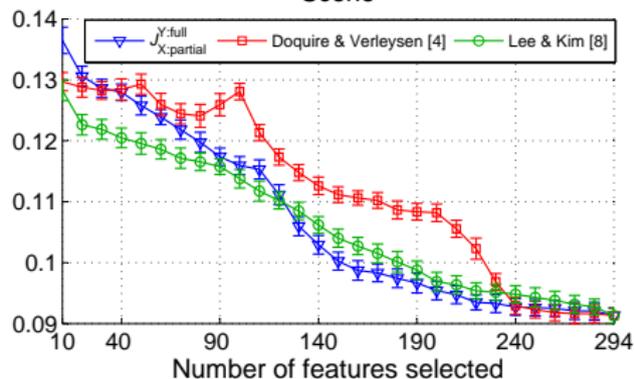- Lee & Kim (2013) : Multivariate Mutual Information

Hamming loss



- Doquire & Verleysen (2013) : $J_{\text{X:none}}^{\text{Y:none}}$ with pruning of rare cases
- Lee & Kim (2013) : Multivariate Mutual Information
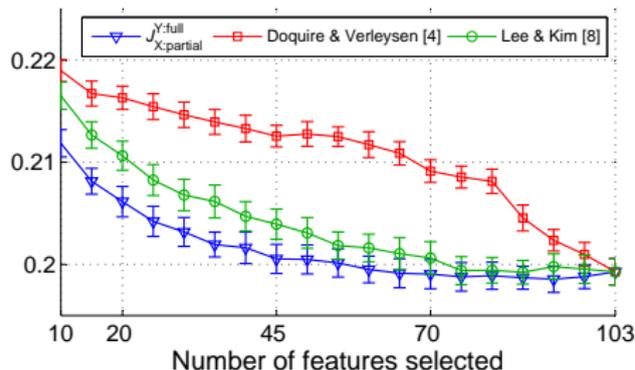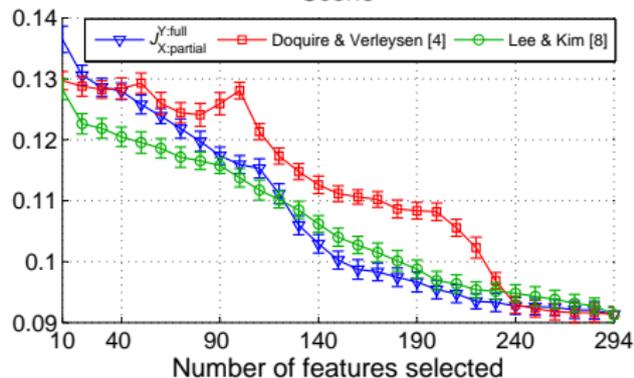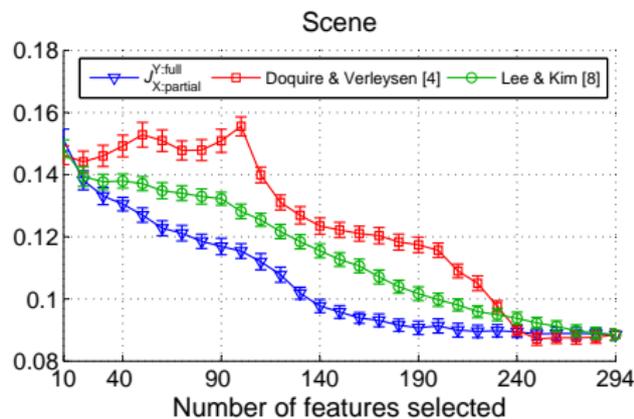- $J_{\text{X:partial}}^{\text{Y:full}}$ tends to outperform state-of-the-art criteria

Ranking loss

- Doquire & Verleysen (2013) : $J_{X:none}^{Y:none}$ with pruning of rare cases
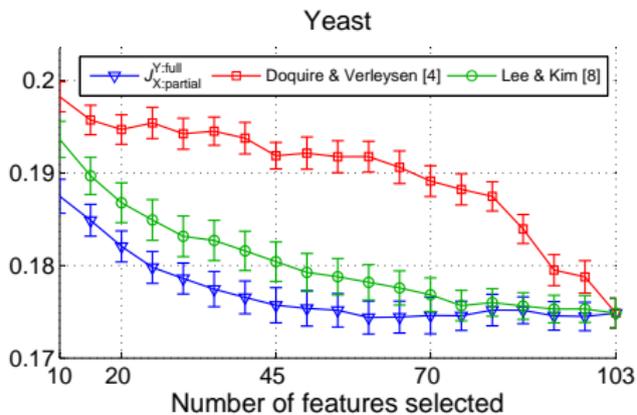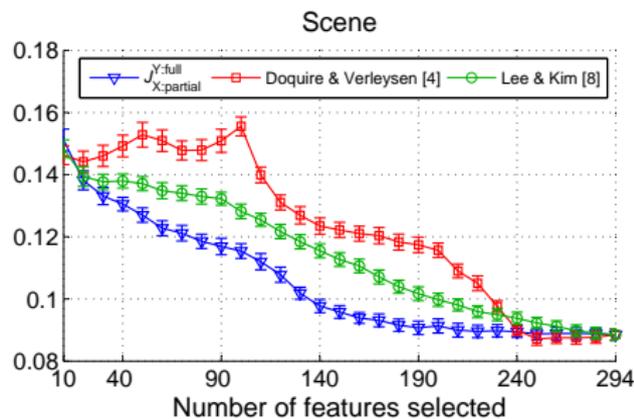- Lee & Kim (2013) : Multivariate Mutual Information

Ranking loss



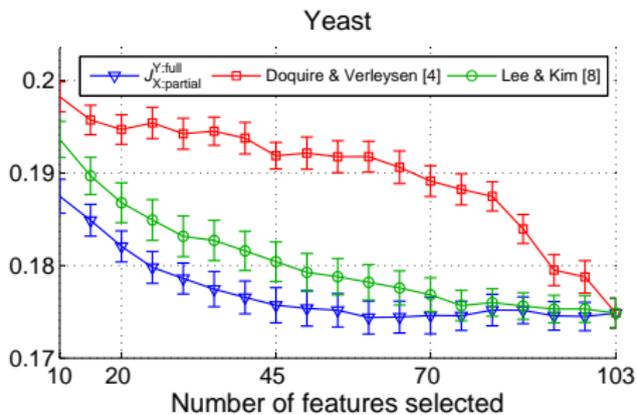- Doquire & Verleysen (2013) : $J_{X:none}^{Y:none}$ with pruning of rare cases
- Lee & Kim (2013) : Multivariate Mutual Information
- $J_{X:partial}^{Y:full}$ dominates state-of-the-art criteria

- Caution: Only 2 datasets! But based on them it appears that...

- Caution: Only 2 datasets! But based on them it appears that...
- ...independence assumptions in label space matter less than in feature space

# Empirical Observations

- Caution: Only 2 datasets! But based on them it appears that...
- ...independence assumptions in label space matter less than in feature space
  - Agrees with Gharroudi et al. (2014) for multilabel-label wrappers

- Caution: Only 2 datasets! But based on them it appears that...
- ...independence assumptions in label space matter less than in feature space
  - ▶ Agrees with Gharroudi et al. (2014) for multilabel-label wrappers
- ...in feature space, JMI gives best results

- Caution: Only 2 datasets! But based on them it appears that...
- ...independence assumptions in label space matter less than in feature space
  - ▶ Agrees with Gharroudi et al. (2014) for multilabel-label wrappers
- ...in feature space, JMI gives best results
  - ▶ Examining pairwise interactions seems a good compromise between capturing interdependencies vs obtaining reliable estimates...

# Empirical Observations

- Caution: Only 2 datasets! But based on them it appears that...
- ...independence assumptions in label space matter less than in feature space
  - ▸ Agrees with Gharroudi et al. (2014) for multilabel-label wrappers
- ...in feature space, JMI gives best results
  - ▸ Examining pairwise interactions seems a good compromise between capturing interdependencies vs obtaining reliable estimates...
  - ▸ Agrees with Brown et al. (JMLR 2012) findings in single-label filters

- Probabilistic framework allows explicit incorporation of domain knowledge...

- Probabilistic framework allows explicit incorporation of domain knowledge...

- ...as informative priors $P(X)$ or $P(Y)$
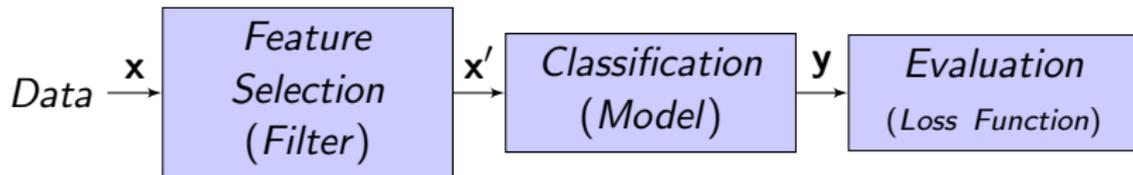
# Future Work: Incorporating Domain Knowledge

- Probabilistic framework allows explicit incorporation of domain knowledge...

- ...as informative priors $P(X)$ or $P(Y)$

- ...as to how the distribution $P(Y|X)$ is factored

# Future Work: Incorporating Domain Knowledge

- Probabilistic framework allows explicit incorporation of domain knowledge...

- ...as informative priors $P(X)$ or $P(Y)$

- ...as to how the distribution $P(Y|X)$ is factored

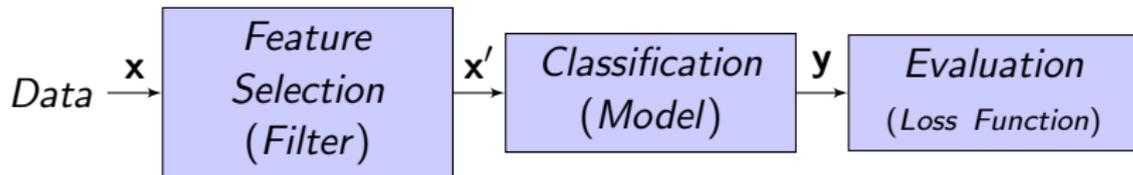- Thus constructing more problem specific filters

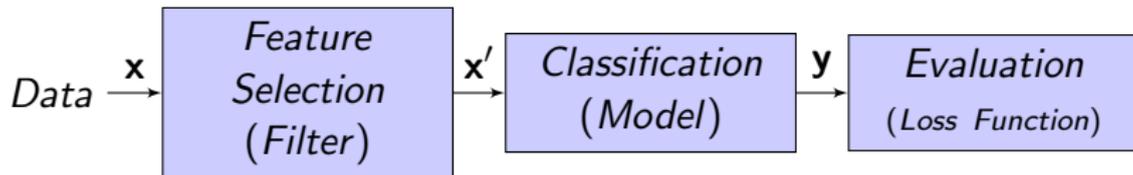- A typical machine learning pipeline

- A typical machine learning pipeline



- Assumptions in every step, often conflicting...

# Future Work: The Bigger Picture

- A typical machine learning pipeline



- Assumptions in every step, often conflicting...
- ...should investigate interplay between model, filter & loss function

# Thank you!
# Kiitos!