

Label Space Coding for Multi-label Classification

Hsuan-Tien Lin

National Taiwan University

Talk at IIS Sinica, April 26, 2012

joint works with

*Farbound Tai (MLD Workshop 2011, NC Journal 2012) &
Chun-Sung Ferng (ACML Conference 2011, NTU Thesis 2012)*



Which Fruit?



?



apple



orange



strawberry



kiwi

multi-class classification:
classify input (picture) to **one category** (label)



Which Fruits?



?: {orange, strawberry, kiwi}



apple



orange



strawberry



kiwi

multi-label classification:
classify input to **multiple (or no)** categories



Powerset: Multi-label Classification via Multi-class

Multi-class w/ $L = 4$ classes

4 possible outcomes

$\{a, o, s, k\}$

Multi-label w/ $L = 4$ classes

$2^4 = 16$ possible outcomes

$2^{\{a, o, s, k\}}$



$\{ \phi, a, o, s, k, ao, as, ak, os, ok, sk, aos, aok, ask, osk, aosk \}$

- **Powerset** approach: reduction to multi-class classification
- difficulties for large L :
 - **computation** (super-large 2^L)
 - hard to construct classifier
 - **sparsity** (no example for some of 2^L)
 - hard to discover hidden combination

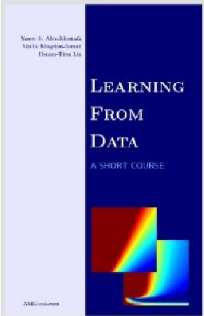
Powerset: feasible only for **small** L with enough examples for every combination




What Tags?

http://www.amazon.com/gp/product/1600490069

Amazon.com: Learning From D...



Learning From Data [Hardcover]
 Yaser S. Abu-Mostafa (Author), Malik Magdon-Ismael (Author),
 Hsuan-Tien Lin (Author)
 ★★★★★ (2 customer reviews) |  Liked (9)

Available from [these sellers](#).

1 new from \$28.00

?: { machine learning, data-structure, data mining, object oriented-programming, artificial intelligence, compiler, architecture, chemistry, textbook, children-book, ... etc. }

another **multi-label** classification problem:
tagging input to multiple categories



Binary Relevance: Multi-label Classification via Yes/No

Binary Classification

{yes, no}

Multi-label w/ L classes: L yes/no questions

machine learning (Y), data structure (N), data mining (Y), OOP (N), AI (Y), compiler (N), architecture (N), chemistry (N), textbook (Y), children book (N), *etc.*

- **Binary Relevance** approach:
reduction to **multiple isolated binary classification**
- disadvantages:
 - **isolation**—hidden relations not exploited (e.g. ML and DM **highly correlated**, ML **subset of** AI, textbook & children book **disjoint**)
 - **unbalancedness**—few **yes**, many **no**

Binary Relevance: simple (& good) benchmark with known disadvantages



Multi-label Classification Setup

Given

N examples (input \mathbf{x}_n , label-set $\mathcal{Y}_n \in \mathcal{X} \times 2^{\{1,2,\dots,L\}}$)

- fruits: $\mathcal{X} = \text{encoding}(\text{pictures})$, $\mathcal{Y}_n \subseteq \{1, 2, \dots, 4\}$
- tags: $\mathcal{X} = \text{encoding}(\text{merchandise})$, $\mathcal{Y}_n \subseteq \{1, 2, \dots, L\}$

Goal

a multi-label classifier $g(\mathbf{x})$ that **closely predicts** the label-set \mathcal{Y} associated with some **unseen** inputs \mathbf{x} (by **exploiting hidden relations/combinations between labels**)

- **0/1 loss**: any discrepancy $\llbracket g(\mathbf{x}) \neq \mathcal{Y} \rrbracket$
- **Hamming loss**: averaged symmetric difference $\frac{1}{L} |g(\mathbf{x}) \triangle \mathcal{Y}|$

multi-label classification: hot and important



Topics in this Talk

- ① **Coding/Geometric** View of Multi-label Classification
—**unify** existing algorithms w/ intuitive explanations
- ② **Compression** Coding
—**condense** for efficiency
—capture hidden correlation
- ③ **Error-correction** Coding
—**expand** for accuracy
—capture hidden combination



Coding/Geometric View of Multi-label Classification



From Label-set to Coding View

	label set	apple	orange	strawberry	binary code
	$\mathcal{Y}_1 = \{o\}$	0 (N)	1 (Y)	0 (N)	$\mathbf{y}_1 = [0, 1, 0]$
	$\mathcal{Y}_2 = \{a, o\}$	1 (Y)	1 (Y)	0 (N)	$\mathbf{y}_2 = [1, 1, 0]$
	$\mathcal{Y}_3 = \{a, s\}$	1 (Y)	0 (N)	1 (Y)	$\mathbf{y}_3 = [1, 0, 1]$
	$\mathcal{Y}_4 = \{o\}$	0 (N)	1 (Y)	0 (N)	$\mathbf{y}_4 = [0, 1, 0]$
	$\mathcal{Y}_5 = \{\}$	0 (N)	0 (N)	0 (N)	$\mathbf{y}_5 = [0, 0, 0]$

subset \mathcal{Y} of $2^{\{1,2,\dots,L\}}$ \Leftrightarrow length- L binary code \mathbf{y}



Existing Approach: Compressive Sensing

General Compressive Sensing

sparse (many 0) binary vectors $\mathbf{y} \in \{0, 1\}^L$ can be **robustly compressed** by projecting to $M \ll L$ basis vectors $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M\}$

Compressive Sensing for Multi-label Classification (Hsu et al., 2009)

- 1 **compress**: transform $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ to $\{(\mathbf{x}_n, \mathbf{c}_n)\}$ by $\mathbf{c}_n = \mathbf{P}\mathbf{y}_n$ with some M by L **random** matrix $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M]^T$
- 2 **learn**: get **regression** function $\mathbf{r}(\mathbf{x})$ from \mathbf{x}_n to \mathbf{c}_n
- 3 **decode**: $g(\mathbf{x}) = \text{find closest sparse binary vector to } \mathbf{P}^T \mathbf{r}(\mathbf{x})$






Compressive Sensing:

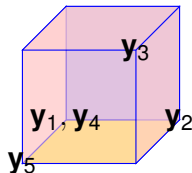
reduction to multi-output regression w/ **codewords \mathbf{c}**

- efficient in training: **random projection** w/ $M \ll L$
- inefficient in testing: **time-consuming decoding**



From Coding View to Geometric View

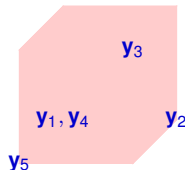
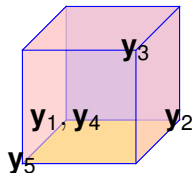
	label set	binary code
	$\mathcal{Y}_1 = \{o\}$	$\mathbf{y}_1 = [0, 1, 0]$
	$\mathcal{Y}_2 = \{a, o\}$	$\mathbf{y}_2 = [1, 1, 0]$
	$\mathcal{Y}_3 = \{a, s\}$	$\mathbf{y}_3 = [1, 0, 1]$
	$\mathcal{Y}_4 = \{o\}$	$\mathbf{y}_4 = [0, 1, 0]$
	$\mathcal{Y}_5 = \{\}$	$\mathbf{y}_5 = [0, 0, 0]$



length- L binary code \Leftrightarrow vertex of hypercube $\{0, 1\}^L$



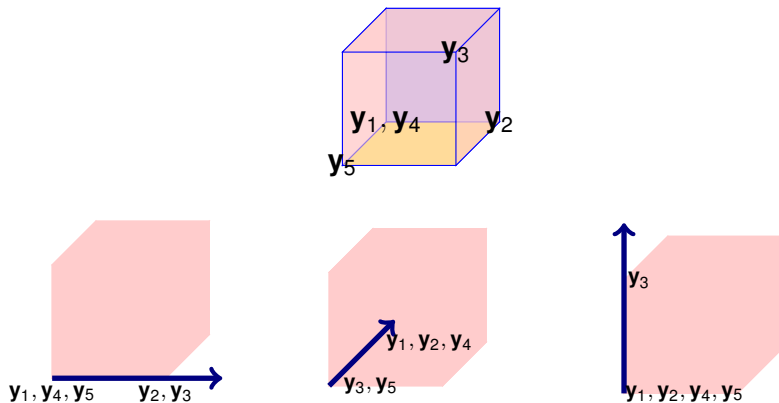
Geometric Interpretation of Powerset



Powerset: directly classify to the **vertices** of hypercube



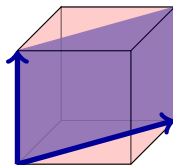
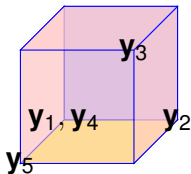
Geometric Interpretation of Binary Relevance



Binary Relevance: project to the **natural axes** & classify



Geometric Interpretation of Compressive Sensing



Compressive Sensing:

- project to **random flat** (linear subspace)
- learn “on” the flat; decode to **closest sparse vertex**

other (better) flat? other (faster) decoding?



Our Contributions (First Part)

Compression Coding (Using Geometry)

A Novel Approach for Label Space Compression

- algorithmic: scheme for **fast decoding**
- theoretical: justification for **best flat**
- practical: **significantly better performance** than compressive sensing (& binary relevance)

will now introduce the key ideas behind the approach



Faster Decoding: Round-based

Compressive Sensing Revisited

- ① **compress**: transform $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ to $\{(\mathbf{x}_n, \mathbf{c}_n)\}$ by $\mathbf{c}_n = \mathbf{P}\mathbf{y}_n$ with some M by L **random** matrix \mathbf{P}
- ② **learn**: get **regression** function $\mathbf{r}(\mathbf{x})$ from \mathbf{x}_n to \mathbf{c}_n
- ③ **decode**: $g(\mathbf{x}) =$ find **closest sparse binary vector** to $\mathbf{P}^T \mathbf{r}(\mathbf{x})$

- find closest **sparse** binary vector to $\tilde{\mathbf{y}}$: **slow**
optimization of ℓ_1 -**regularized** objective
- find closest **any** binary vector to $\tilde{\mathbf{y}}$: **fast**

$$g(\mathbf{x}) = \text{round}(\mathbf{y})$$

round-based decoding: simple & faster alternative



Better Flat: Principal Directions

Compressive Sensing Revisited

- 1 **compress**: transform $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ to $\{(\mathbf{x}_n, \mathbf{c}_n)\}$ by $\mathbf{c}_n = \mathbf{P}\mathbf{y}_n$ with some M by L **random** matrix \mathbf{P}
- 2 **learn**: get **regression** function $\mathbf{r}(\mathbf{x})$ from \mathbf{x}_n to \mathbf{c}_n
- 3 **decode**: $g(\mathbf{x}) =$ find **closest sparse binary vector** to $\mathbf{P}^T \mathbf{r}(\mathbf{x})$

- **random** flat: **arbitrary** directions
- **best** flat: **principal** directions

principal directions/flat: best approximation to vertices \mathbf{y}_n during **compression** (**why?**)



Novel Theoretical Guarantee

Linear Transform + Regress + Round-based Decoding

Theorem (Tai and Lin, 2012)

If $g(\mathbf{x}) = \text{round}(\mathbf{P}^T \mathbf{r}(\mathbf{x}))$,

$$\underbrace{\frac{1}{L} |g(\mathbf{x}) \triangle \mathcal{Y}|}_{\text{Hamming loss}} \leq \text{const} \cdot \left(\underbrace{\|\mathbf{r}(\mathbf{x}) - \overbrace{\mathbf{P}\mathbf{y}}^{\mathbf{c}}\|^2}_{\text{learn}} + \underbrace{\|\mathbf{y} - \mathbf{P}^T \overbrace{\mathbf{P}\mathbf{y}}^{\mathbf{c}}\|^2}_{\text{compress}} \right)$$

- $\|\mathbf{r}(\mathbf{x}) - \mathbf{c}\|^2$: prediction error from input to codeword
- $\|\mathbf{y} - \mathbf{P}^T \mathbf{c}\|^2$: encoding error from vertex to codeword

principal directions/flat: best approximation to vertices \mathbf{y}_n during compression (indeed)



Proposed Approach: Principal Label Space Transform

From Compressive Sensing to PLST

- 1 **compress**: transform $\{(\mathbf{x}_n, \mathbf{y}_n)\}$ to $\{(\mathbf{x}_n, \mathbf{c}_n)\}$ by $\mathbf{c}_n = \mathbf{P}(\mathbf{y}_n - \mathbf{o})$ with the M by L **principal** matrix \mathbf{P} and **some reference point \mathbf{o}**
- 2 **learn**: get regression function $\mathbf{r}(\mathbf{x})$ from \mathbf{x}_n to \mathbf{c}_n
- 3 **decode**: $g(\mathbf{x}) = \text{round}(\mathbf{P}^T \mathbf{r}(\mathbf{x}) + \mathbf{o})$

- reference point \mathbf{o} : allow flat **not passing the origin**
- best \mathbf{o} and \mathbf{P} :

$$\begin{aligned}
 & \min_{\mathbf{o}, \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M} \sum_{n=1}^N \left\| \mathbf{y}_n - \mathbf{o} - \mathbf{P}^T \mathbf{P}(\mathbf{y}_n - \mathbf{o}) \right\|^2 \\
 & \text{subject to} \quad \text{orthonormal vectors } \mathbf{p}_m
 \end{aligned}$$



Solving for Principal Directions

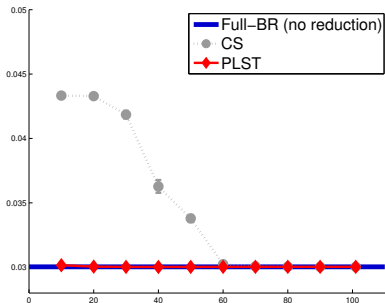
$$\begin{aligned} \min_{\mathbf{o}, \mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M} \quad & \sum_{n=1}^N \left\| \mathbf{y}_n - \mathbf{o} - \mathbf{P}^T \mathbf{P} (\mathbf{y}_n - \mathbf{o}) \right\|^2 \\ \text{subject to} \quad & \text{orthonormal vectors } \mathbf{p}_m \end{aligned}$$

- solution: **Principal Component Analysis** on $\{\mathbf{y}_n\}_{n=1}^N$
- best \mathbf{o} : $\frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$
- best \mathbf{p}_m : top eigenvectors of $\sum_{n=1}^N (\mathbf{y}_n - \mathbf{o})(\mathbf{y}_n - \mathbf{o})^T$
- physical meaning behind \mathbf{p}_m :
key (linear) **label correlations** (e.g. like eigenface in face recognition)

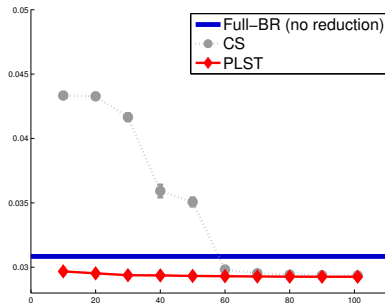
PLST: reduction to multi-output regression
by projecting to **key correlations**



Hamming Loss Comparison: Full-BR, PLST & CS



mediamill (Linear Regression)



mediamill (Decision Tree)

- **PLST** better than **Full-BR**: fewer dimensions, similar (or **better**) performance
- **PLST** better than **CS**: faster, **better** performance
- similar findings across **data sets** and **regression algorithms**



Semi-summary on PLST

- reduction to **multi-output regression**
- project to **principal directions** and capture key correlations
- efficient learning (**label space compression**)
- efficient decoding (**round-based**)
- sound theoretical guarantee + **good practical performance** (better than CS & BR)

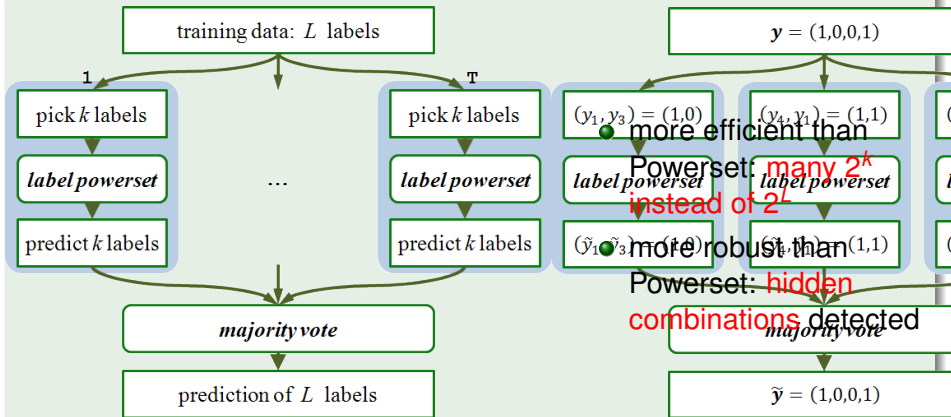
expansion (channel coding) instead of compression (“lossy” source coding)? **YES!**

will start by reviewing an existing algorithm



Random k -labelsets

Random k -labelsets (Tsoumakas & Vlahavas, 2007)



RAkEL:

reduction to many 2^k -category classification tasks



Random k -labelsets from Coding View

RAkELs (Tsoumakas & Vlahavas, 2007)

$$\mathbf{y} = (1,0,0,1)$$

$$\mathbf{b} = (y_1, y_3, y_4, y_1, y_2, y_1) = (1,0,1,1,0,1)$$

2-label powerset

$$\tilde{\mathbf{b}} = (y_1, y_3, y_4, y_1, y_2, y_1) = (1,0,1,1,0,0)$$

majority vote

$$\tilde{\mathbf{y}} = (1,0,0,1)$$

- **encode ($\mathbf{y} \rightarrow \mathbf{b}$):**
repetition & permutation
- **learn:** k -label powerset,
i.e. run Powerset on
size- k chunk of bits
- **decode ($\tilde{\mathbf{b}} \rightarrow \tilde{\mathbf{y}}$):**
majority vote

RAkEL: encode + learn + decode



Our Contributions (Second Part)

Error-correction Coding

A Novel Framework for Label Space Error-correction

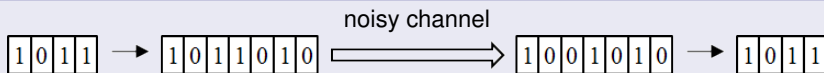
- algorithmic: generalize RAKEL and explain through **coding view**
- theoretical: link learning performance to **error-correcting ability**
- practical: explore **choices of error-correcting code** and obtain **better performance** than RAKEL (& binary relevance)

will now introduce the framework



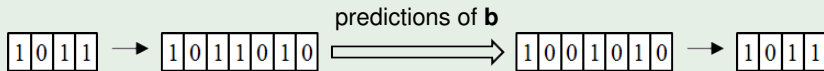
Key Idea: Redundant Information

General Error-correcting Codes (ECC)



- commonly used in communication systems
- detect & correct errors after transmitting data over a noisy channel
- encode data **redundantly**

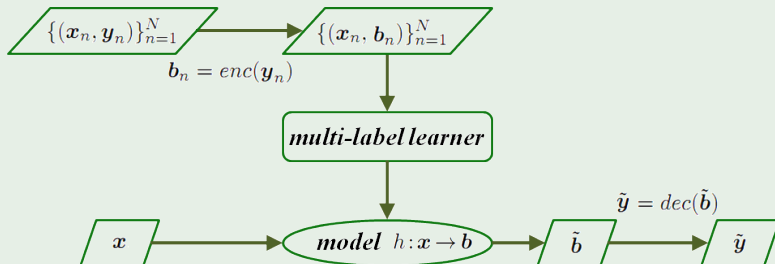
ECC for Machine Learning (successful for multi-class classification)



learn **redundant bits** \Rightarrow **correct** prediction **errors**



Proposed Framework: Multi-labeling with ECC



- **encode** to add redundant information $enc(\cdot): \{0, 1\}^L \rightarrow \{0, 1\}^M$
- **decode** to locate most possible binary vector $dec(\cdot): \{0, 1\}^M \rightarrow \{0, 1\}^L$
- reduction to **larger multi-label classification** with labels **b**

PLST: $M \ll L$ (works for large L);
MLECC: $M > L$ (works for small L)



Simple Theoretical Guarantee

ECC encode + Larger Multi-label Learning + **ECC decode**

Theorem

Let $g(\mathbf{x}) = \text{dec}(\tilde{\mathbf{b}})$ with $\tilde{\mathbf{b}} = h(\mathbf{x})$. Then,

$$\underbrace{\mathbb{I}[g(\mathbf{x}) \neq \mathcal{Y}]}_{0/1 \text{ loss}} \leq \text{const.} \cdot \frac{\text{Hamming loss of } h(\mathbf{x})}{\text{ECC strength} + 1}.$$

PLST: **principal directions** + decent regression
MLECC: which ECC balances **strength** & **difficulty**?



Simplest ECC: Repetition Code

encoding: $\mathbf{y} \in \{0, 1\}^L \rightarrow \mathbf{b} \in \{0, 1\}^M$

- **repeat** each bit $\frac{M}{L}$ times

$$L = 4, M = 28 : 1010 \longrightarrow \underbrace{1111111}_{\frac{28}{4}=7} 000000011111110000000$$

- permute the bits randomly

decoding: $\tilde{\mathbf{b}} \in \{0, 1\}^M \rightarrow \tilde{\mathbf{y}} \in \{0, 1\}^L$

- **majority vote** on each original bit

$L = 4, M = 28$: strength of repetition code (REP) = 3

RAkEL = REP + k -label powerset



Slightly More Sophisticated: Hamming Code

HAM(7, 4) Code

- $\{0, 1\}^4 \rightarrow \{0, 1\}^7$ via adding 3 **parity bits**
—physical meaning: **label combinations**
- $b_4 = y_0 \oplus y_1 \oplus y_3$, $b_5 = y_0 \oplus y_2 \oplus y_3$, $b_6 = y_1 \oplus y_2 \oplus y_3$
- e.g. 1011 \rightarrow 1011010
- strength = 1 (weak)

Our Proposed Code: Hamming on Repetition (HAMR)

$$\{0, 1\}^L \xrightarrow{\text{REP}} \{0, 1\}^{\frac{4M}{7}} \xrightarrow{\text{HAM}(7, 4) \text{ on each 4-bit block}} \{0, 1\}^{\frac{7M}{7}}$$

$L = 4$, $M = 28$: strength of HAMR = 4 **better** than REP!

HAMR + k -label powerset:
improvement of RAKEL on **code strength**



Even More Sophisticated Codes

Bose-Chaudhuri-Hocquenghem Code (BCH)

- modern code in **CD players**
- sophisticated extension of Hamming, with **more parity bits**
- codeword length $M = 2^p - 1$ for $p \in \mathbb{N}$
- $L = 4$, $M = 31$, strength of BCH = 5

Low-density Parity-check Code (LDPC)

- modern code for **satellite communication**
- connect ECC and Bayesian learning
- approach the theoretical limit in some cases

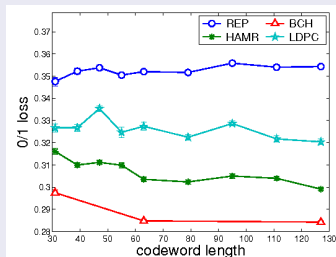
let's compare!



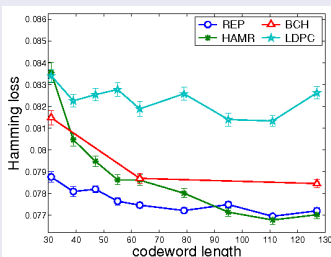
Different ECCs on 3-label Powerset (scene data set w/ $L = 6$)

- learner: 3-label powerset with Random Forests
- REP + 3-label powerset \approx RAKEL

0/1 loss



Hamming loss



Comparing to RAKEL (on most of data sets),

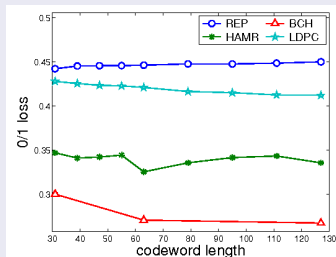
- HAMR: **better 0/1 loss**, similar Hamming loss
- BCH: **even better 0/1 loss**, pay for Hamming loss



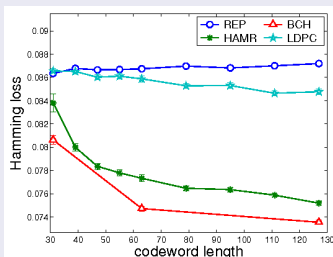
Different ECCs on Binary Relevance (scene data set w/ $L = 6$)

- Binary Relevance: simply 1-label powerset
- REP + Binary Relevance \approx Binary Relevance (with aggregation)

0/1 loss



Hamming loss



Comparing to BR (on most of data sets),

- BCH/HAMR + BR: **better 0/1 loss, better Hamming loss**



Semi-summary on MLECC

- reduction to **larger multi-label classification**
- encode via **error-correcting code** and capture label combinations (parity bits)
- effective decoding (**error-correcting**)
- simple theoretical guarantee + **good practical performance**
 - to **improve RAKE**, replace REP by
 - HAMR \Rightarrow lower 0/1 loss, similar Hamming loss
 - BCH \Rightarrow even lower 0/1 loss, but higher Hamming loss
 - to **improve Binary Relevance**, use
 - HAMR or BCH \Rightarrow lower 0/1 loss, lower Hamming loss



Conclusion

- 1 **Coding/Geometric** View of Multi-label Classification
—useful in linking to **Information Theory** & visualizing
- 2 **Compression** Coding
—**condense** for efficiency: better approach PLST
- 3 **Error-correction** Coding
—**expand** for accuracy: better code HAMR or BCH

More.....

- more geometric explanations (Tai & Lin, NC Journal 2012)
- beyond standard ECC-decoding (Feng, NTU Thesis 2012)
- improved PLST (Chen, NTU Thesis 2012)
- dynamic instead of static coding (...), combine ML-ECC & PLST (...)

Thank you. Questions?

