# IMPROVING GENERALIZATION BY DATA CATEGORIZATION

Ling Li, Amrit Pratap, Hsuan-Tien Lin, and Yaser Abu-Mostafa

Learning Systems Group, Caltech

ECML/PKDD, October 4, 2005

# EXAMPLES IN LEARNING

## A LEARNING SYSTEM

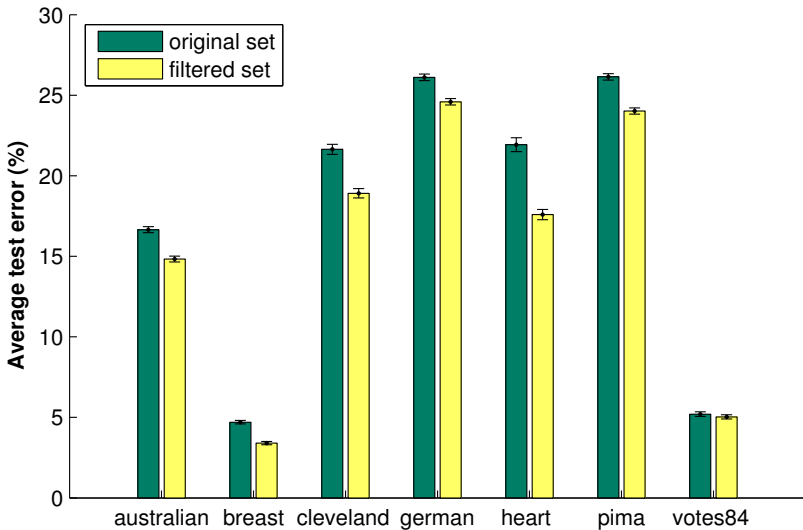Unknown Target $f$ $\longrightarrow$ Examples $\{(\mathbf{x}_i, y_i)\}_i$ $\longrightarrow$ Learner

Examples are essential since they act as the information gateway between the target and the learner.
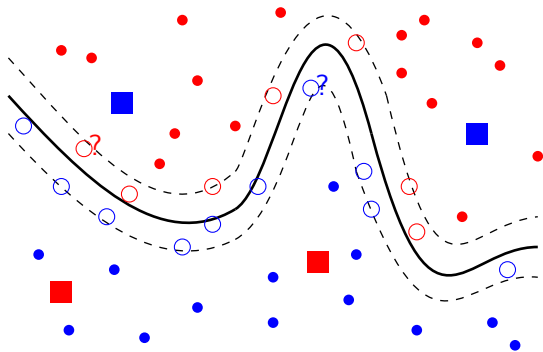
## NOT ALL EXAMPLES ARE EQUALLY USEFUL

1. Surprising examples carry more information     $\surd$

   - Garbage examples are also surprising (Guyon et al., 1996)     $\times$

2. Noisy examples and outliers     $\times$

3. Examples beyond the ability of the learner     $\times$

Can we improve learning by automatically categorizing examples?

INTRODUCTION
○●○
CONCEPTS
○○○
METHODS
○○○
EXPERIMENTS
○○○○○○○
CONCLUSION
○

## IMPROVED GENERALIZATION

INTRODUCTION
○○●

CONCEPTS
○○○

METHODS
○○○

EXPERIMENTS
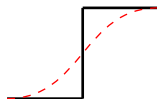○○○○○○○

CONCLUSION
○

# CATEGORIZE EXAMPLES



- Which examples are "bad"?
- Close-to-boundary examples are informative
- Three categories: typical, critical, and noisy

- The automatic data categorization is for better learning.
- The criteria are usually related with how useful or reliable the example is to learning, such as the margin.

## INTRINSIC FUNCTION

The target $f \colon \mathcal{X} \to \{-1, 1\}$ comes from thresholding an
intrinsic function $f_r \colon \mathcal{X} \to \mathbb{R}$. That is

$$f(\mathbf{x}) = \text{sign}\left(f_r(\mathbf{x})\right).$$
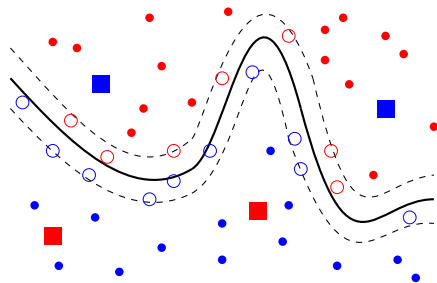
### EXAMPLES OF $f_r(\mathbf{x})$

1. The credit score of the applicant $\mathbf{x}$ minus some threshold
2. The signed Euclidean distance of $\mathbf{x}$ to the boundary
3. The probability of $\mathbf{x}$ belonging to class 1 minus 0.5

### PROPERTIES

- Problem-dependent (e.g., the knowledge of experts)
- Tells the usefulness or reliability of an example
- Unknown

# INTRINSIC MARGIN AND DATA CATEGORIZATION

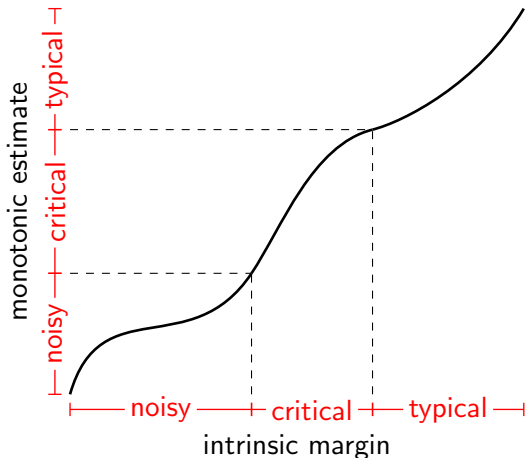For an example $(\mathbf{x}, y)$, its intrinsic margin is $yf_r(\mathbf{x})$.



The intrinsic margin $yf_r(\mathbf{x})$ can be treated as a measure of how close $\mathbf{x}$ is to the decision boundary.

- Small positive: near the boundary                    critical
- Large positive: deep in the class territory          typical
- Negative: mislabeled         noisy
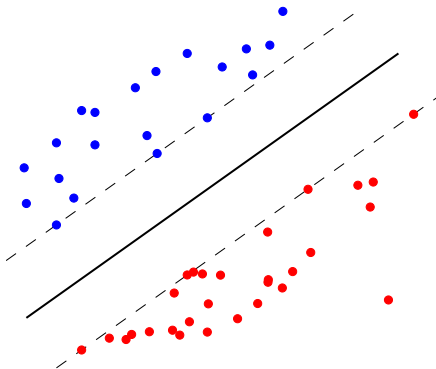
## MONOTONIC ESTIMATE

However, the intrinsic margin is unknown.



- a monotonic estimate of the intrinsic margin
- two proper thresholds
- three categories

## SELECTION COST

For an example $(\mathbf{x}, y)$, a hypothesis $g$ may classify it either wrongly or correctly. Consider the expected out-of-sample errors.

$$\mathsf{E}_g[\pi(g) \mid \underbrace{g(\mathbf{x}) \neq y}_{\text{wrongly}}] - \mathsf{E}_g[\pi(g) \mid \underbrace{g(\mathbf{x}) = y}_{\text{correctly}}]$$

We may select to trust $(\mathbf{x}, y)$, or not. The difference is the cost we pay when we make the selection. We call it the selection cost.



We actually estimate a scaled version of the selection cost (Nicholson, 2002).

The model for learning should be also used for the estimation.

# SVM CONFIDENCE MARGIN

The soft-margin support vector machine (SVM) (Vapnik, 1995) finds a large-confidence hyperplane classifier in the feature space.



- The confidence margin is a meaningful estimate of the intrinsic margin.
- Better than the one used in (Guyon et al., 1996).
- Confidence margin $\leq 1$: support vectors          critical
- Negative margin                          noisy

## ADABOOST SAMPLE WEIGHT

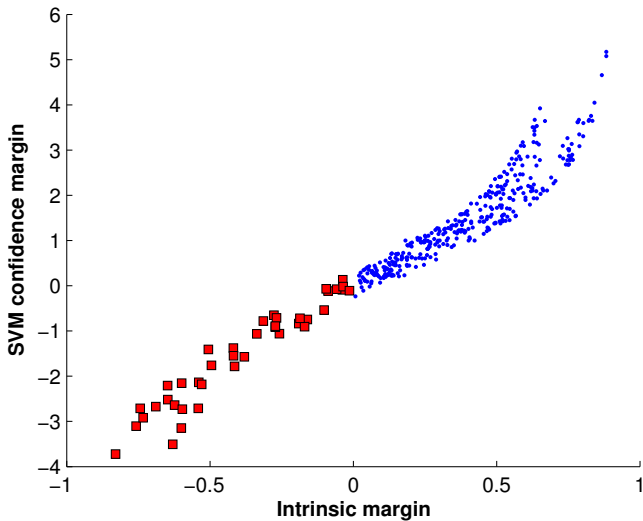AdaBoost (Freund & Schapire, 1996) is an algorithm to improve the accuracy of a base learner.

- It iteratively generates an ensemble of base hypotheses.
- It gradually forces the base learner to focus on "hard" examples by giving erroneous examples higher sample weight.

The sample weight is actually a consensus among the base hypotheses on the "hardness" of the example.

- If an example is too "hard", it is probably noisy.
- If an example is too "easy", it is probably typical.
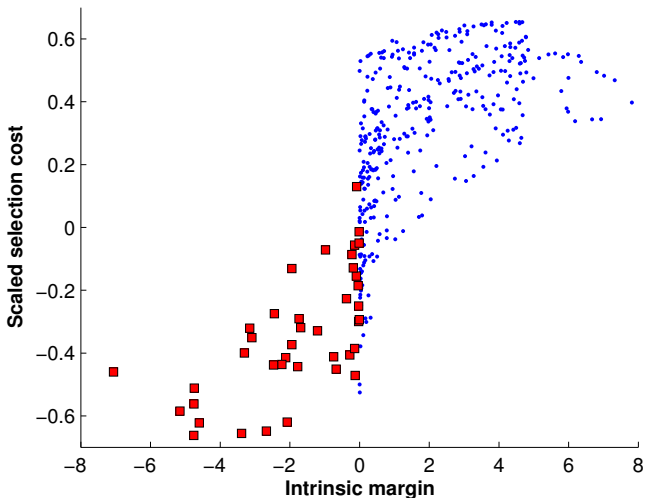- The negative average sample weight over different iterations is a robust estimate of the intrinsic margin.
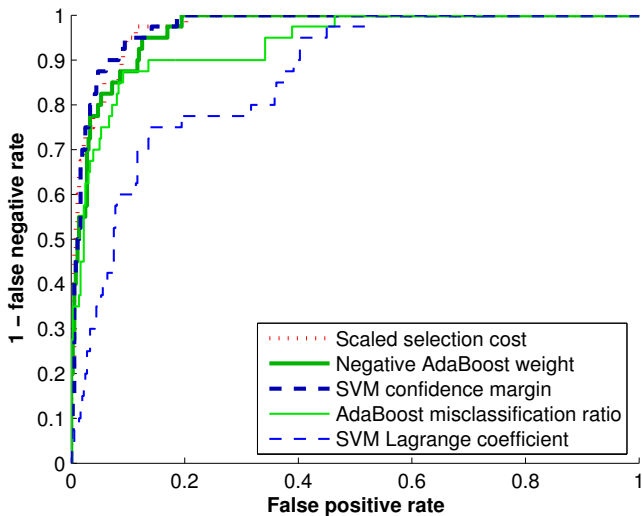
# SCATTER PLOT
3-5-1 NNET

INTRODUCTION
000

CONCEPTS
000

METHODS
000

EXPERIMENTS
0●00000
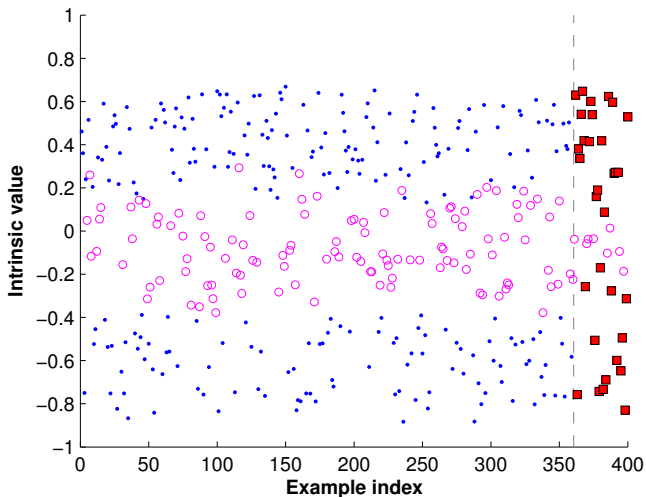
CONCLUSION
0
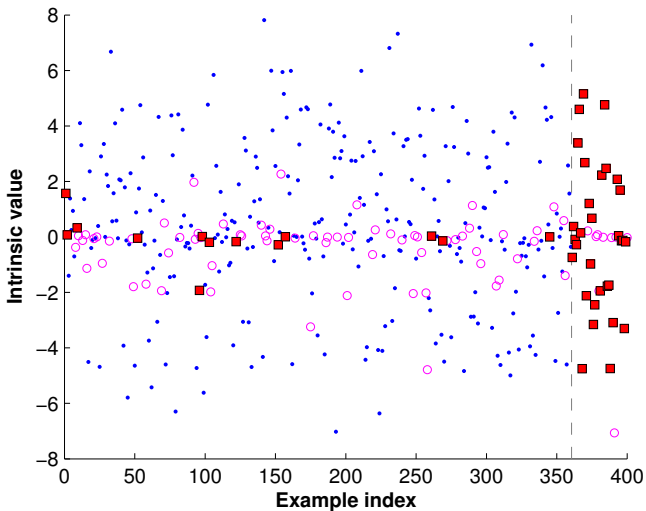
# SCATTER PLOT
## SIN (MERLER ET AL., 2004)

# ROC CURVES
SIN

# FINGERPRINT PLOT
3-5-1 NNET

# FINGERPRINT PLOT
## SIN

## 2-D PLOT
YIN-YANG (HTTP://WWW.WORK.CALTECH.EDU/LING/DATA/YINYANG.HTML)



- typical
- critical, clean
- critical, mislabeled
- noisy, mislabeled
- noisy, clean

# REAL-WORLD DATA

### UTILIZE DATA CATEGORIZATION

It is now possible to treat different categories differently.

- Noisy examples: remove
- Critical examples: emphasize
- Typical examples: reduce

| dataset | orig. dataset | selection cost | SVM margin | AdaBoost weight |
|---------|---------------|----------------|------------|-----------------|
| australian | $16.65 \pm 0.19$ | $15.23 \pm 0.20$ | $14.83 \pm 0.18$ | $13.92 \pm 0.16$ |
| breast | $4.70 \pm 0.11$ | $6.44 \pm 0.13$ | $3.40 \pm 0.10$ | $3.32 \pm 0.10$ |
| cleveland | $21.64 \pm 0.31$ | $18.24 \pm 0.30$ | $18.91 \pm 0.29$ | $18.56 \pm 0.30$ |
| german | $26.11 \pm 0.20$ | $30.12 \pm 0.15$ | $24.59 \pm 0.20$ | $24.68 \pm 0.22$ |
| heart | $21.93 \pm 0.43$ | $17.33 \pm 0.34$ | $17.59 \pm 0.32$ | $18.52 \pm 0.37$ |
| pima | $26.14 \pm 0.20$ | $35.16 \pm 0.20$ | $24.02 \pm 0.19$ | $25.15 \pm 0.20$ |
| votes84 | $5.20 \pm 0.14$ | $6.45 \pm 0.17$ | $5.03 \pm 0.13$ | $4.91 \pm 0.13$ |

# CONCLUSION

## CONTRIBUTIONS

1. Proposed 3 methods for automatically categorizing examples.
   - The methods are from different parts of learning theory.
   - They all gave reasonable categorization results.

2. Tested learning with categorized data.
   - A simple strategy is enough to improve learning.
   - The categorization results can be used in conjunction with a large variety of learning algorithms.

3. Showed experimentally data categorization is powerful.

## FUTURE WORK

- Estimate the optimal thresholds (say, using a validation set)
- Better utilize the categorization in learning
- Extend the framework to problems other than classification