

Performance Measures for Machine Learning

Performance Measures

- Accuracy
- Weighted (Cost-Sensitive) Accuracy
- Lift
- Precision/Recall
 - F
 - Break Even Point
- ROC
 - ROC Area

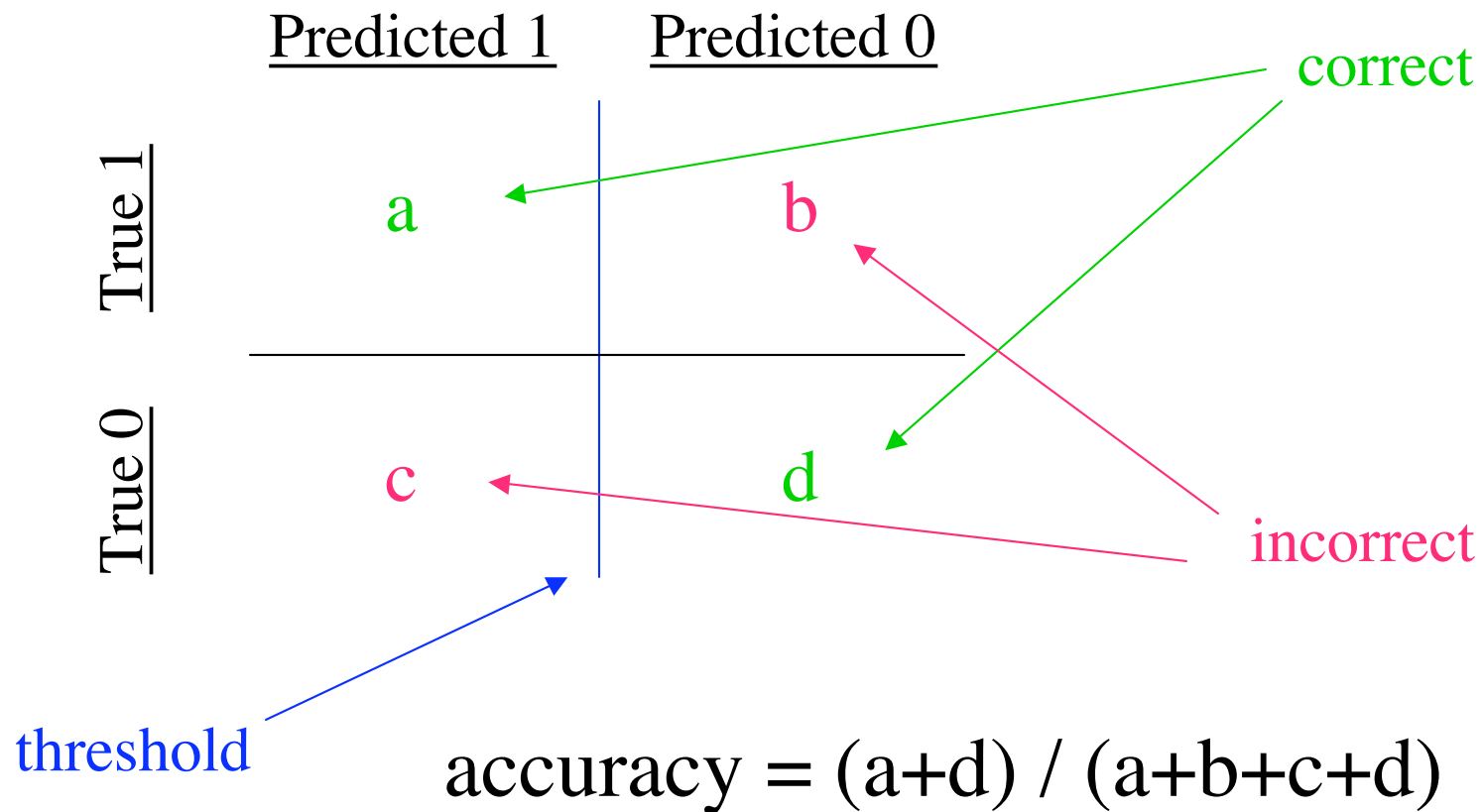
Accuracy

- Target: 0/1, -1/+1, True/False, ...
- Prediction = $f(\text{inputs}) = f(\mathbf{x})$: 0/1 or Real
- Threshold: $f(\mathbf{x}) > \text{thresh} \Rightarrow 1$, else $\Rightarrow 0$
- If $\text{threshold}(f(\mathbf{x}))$ and targets both 0/1:

$$\text{accuracy} = \frac{\sum_{i=1..N} \left(1 - |target_i - \text{threshold}(f(\vec{x}_i))|_{ABS} \right)}{N}$$

- #right / #total
- $p(\text{“correct”})$: $p(\text{threshold}(f(\mathbf{x})) = \text{target})$

Confusion Matrix



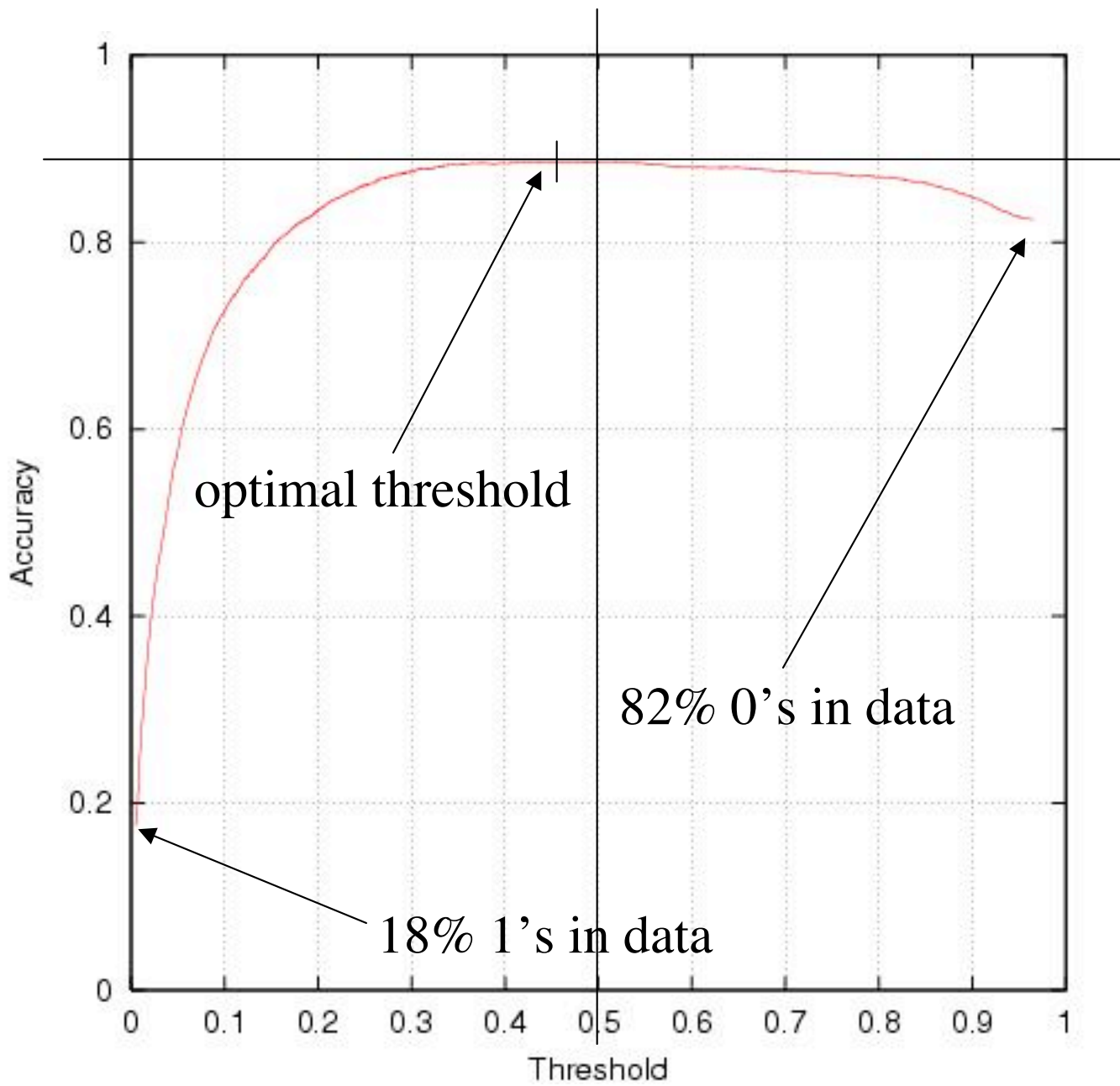
Prediction Threshold

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	0	b
<u>True 0</u>	0	d

- threshold $>$ MAX(f(x))
- all cases predicted 0
- (b+d) = total
- accuracy = %False = %0's

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	a	0
<u>True 0</u>	c	0

- threshold $<$ MIN(f(x))
- all cases predicted 1
- (a+c) = total
- accuracy = %True = %1's



threshold demo

Problems with Accuracy

- Assumes equal cost for both kinds of errors
 - $\text{cost}(\text{b-type-error}) = \text{cost}(\text{c-type-error})$
- is 99% accuracy good?
 - can be excellent, good, mediocre, poor, terrible
 - depends on problem
- is 10% accuracy bad?
 - information retrieval
- BaseRate = accuracy of predicting predominant class
(on most problems obtaining BaseRate accuracy is easy)

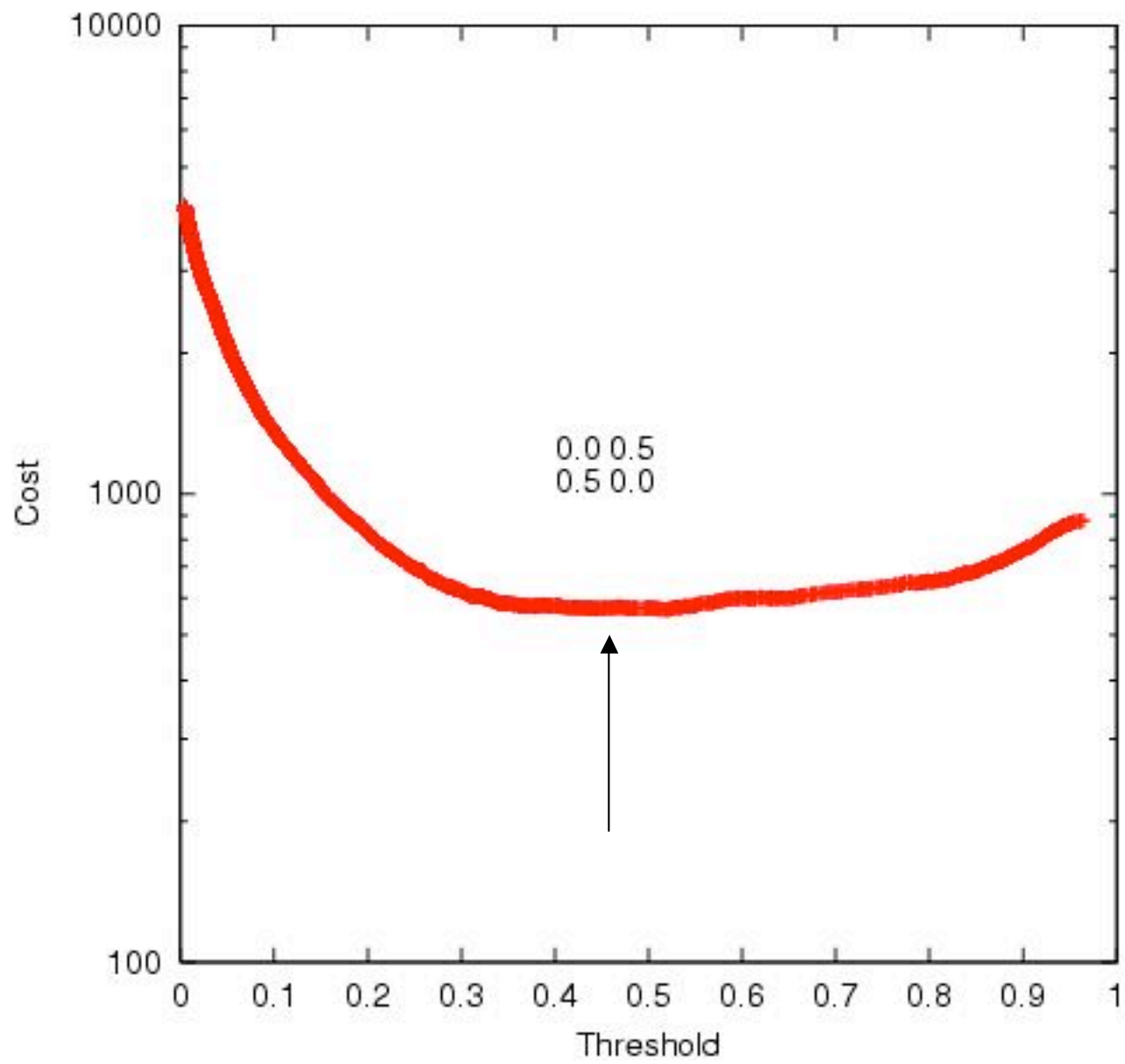
Percent Reduction in Error

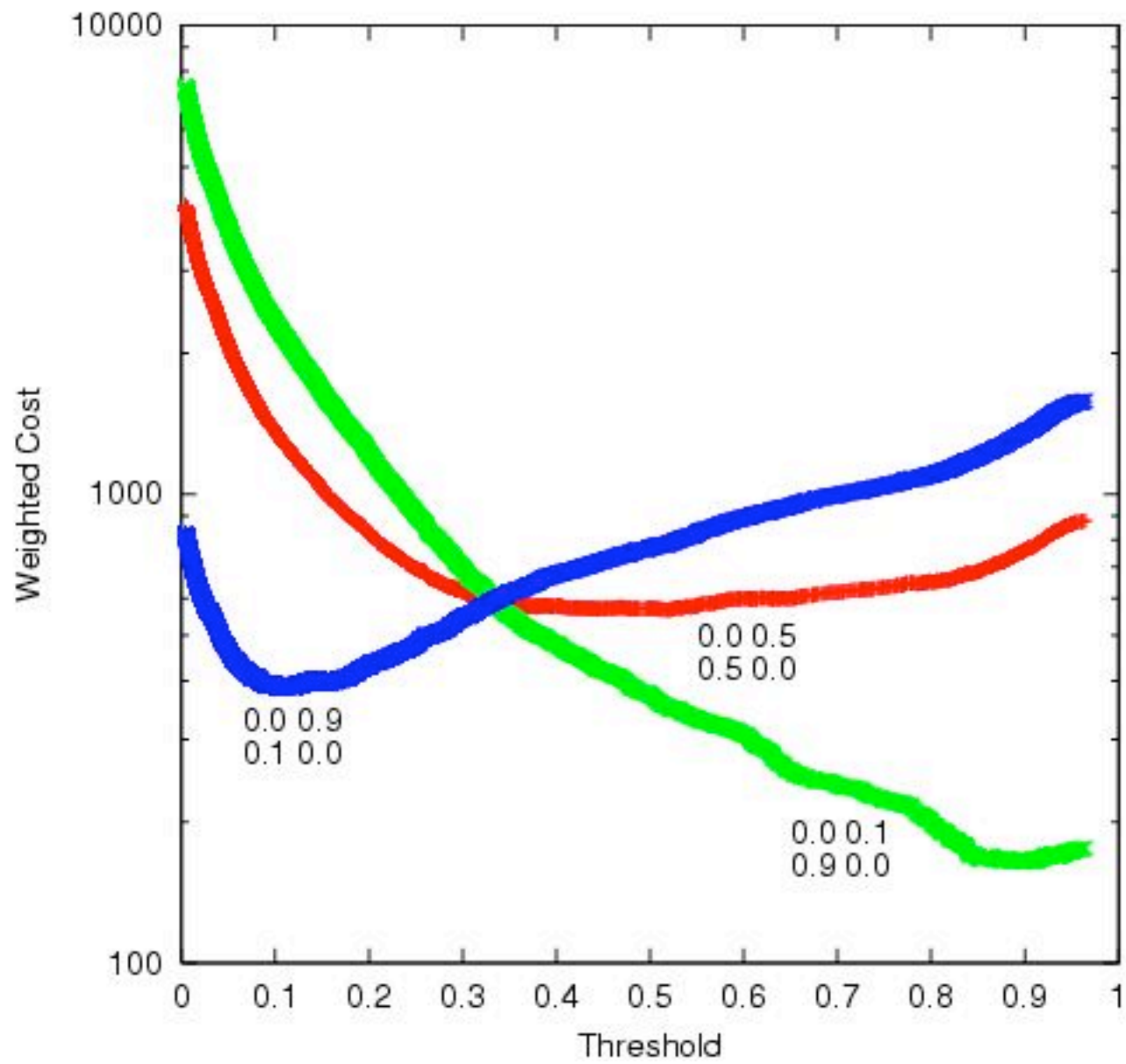
- 80% accuracy = 20% error
 - suppose learning increases accuracy from 80% to 90%
 - error reduced from 20% to 10%
 - 50% reduction in error
-
- 99.90% to 99.99% = 90% reduction in error
 - 50% to 75% = 50% reduction in error
 - can be applied to many other measures

Costs (Error Weights)

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	W_a	W_b
<u>True 0</u>	W_c	W_d

- Often $W_a = W_d = \mathbf{zero}$ and $W_b \neq W_c \neq \mathbf{zero}$





Lift

- not interested in accuracy on entire dataset
- want accurate predictions for 5%, 10%, or 20% of dataset
- don't care about remaining 95%, 90%, 80%, resp.
- typical application: marketing

$$\text{lift}(\text{threshold}) = \frac{\% \text{positives} > \text{threshold}}{\% \text{dataset} > \text{threshold}}$$

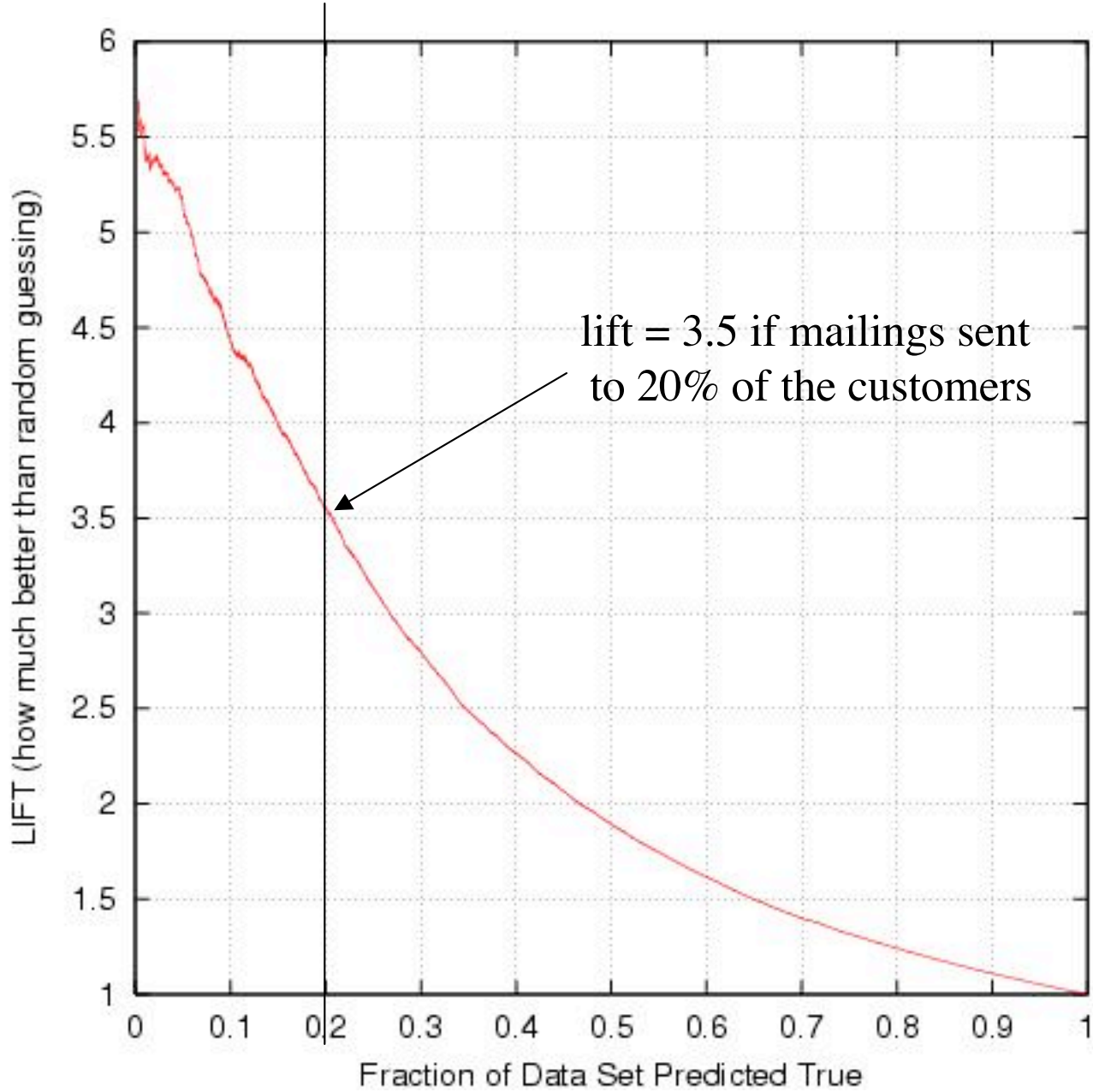
- how much better than random prediction on the fraction of the dataset predicted true ($f(x) > \text{threshold}$)

Lift

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	a	b
<u>True 0</u>	c	d

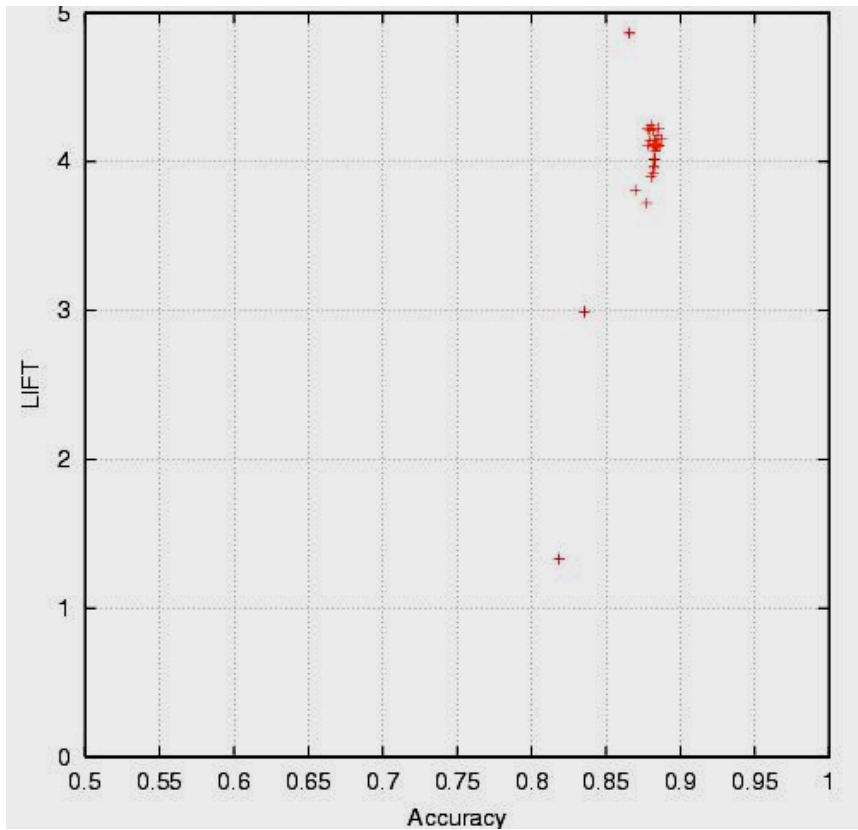
threshold

$$lift = \frac{a/(a + b)}{(a + c)/(a + b + c + d)}$$

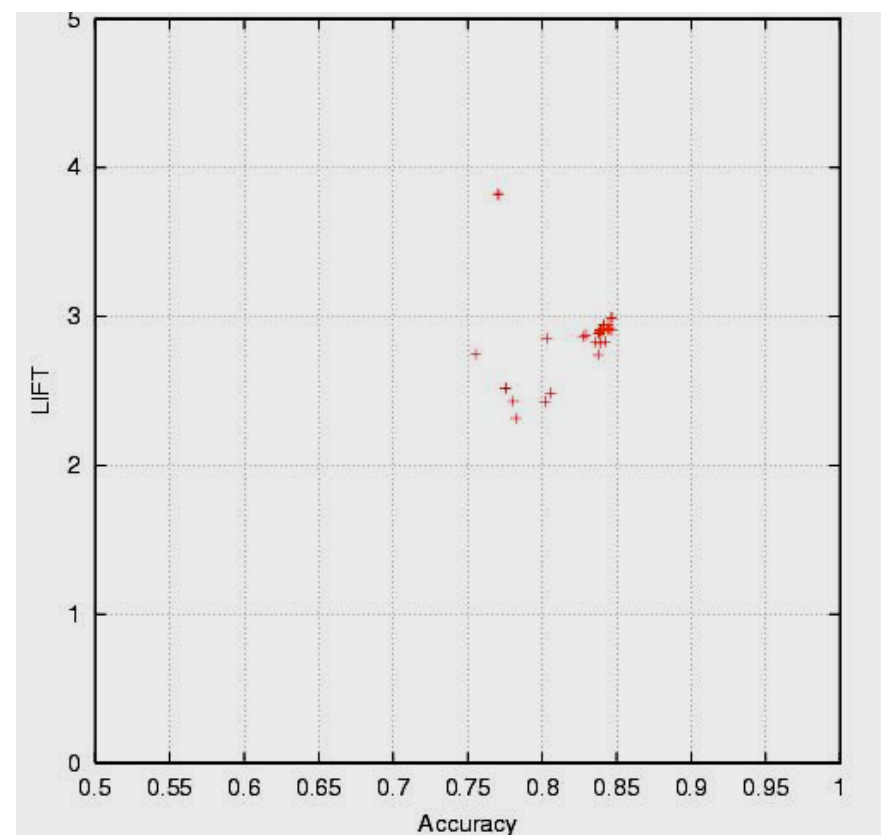


Lift and Accuracy do not always correlate well

Problem 1



Problem 2

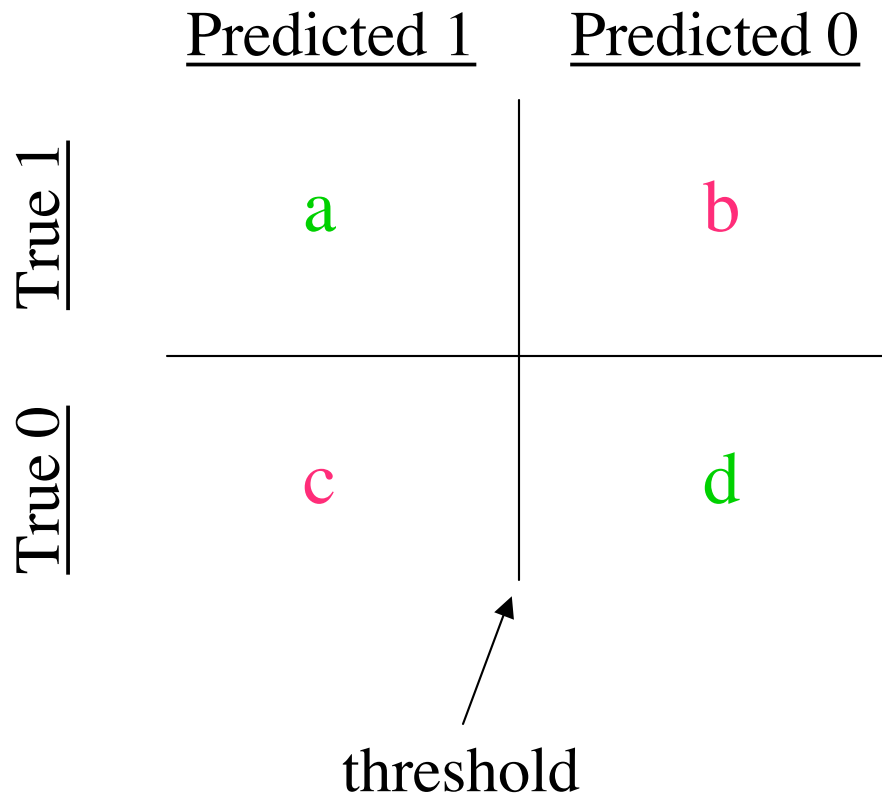


(thresholds arbitrarily set at 0.5 for both lift and accuracy)

Precision and Recall

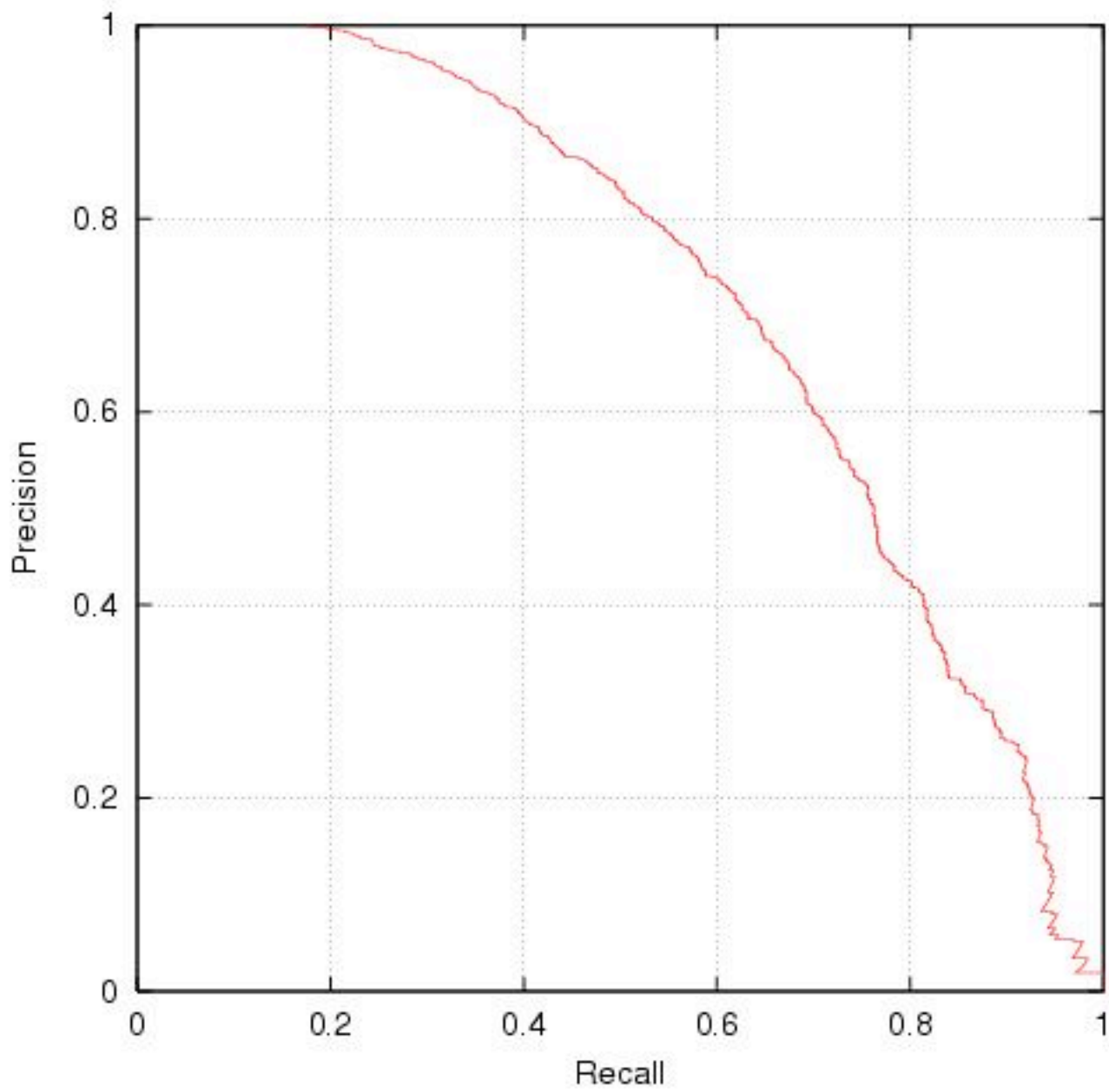
- typically used in document retrieval
- Precision:
 - how many of the returned documents are correct
 - $\text{precision}(\text{threshold})$
- Recall:
 - how many of the positives does the model return
 - $\text{recall}(\text{threshold})$
- Precision/Recall Curve: sweep thresholds

Precision/Recall



$$PRECISION = a / (a + c)$$

$$RECALL = a / (a + b)$$



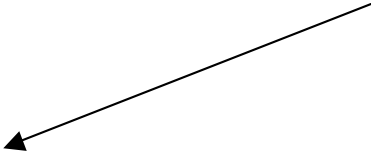
Summary Stats: F & BreakEvenPt

$$PRECISION = a / (a + c)$$

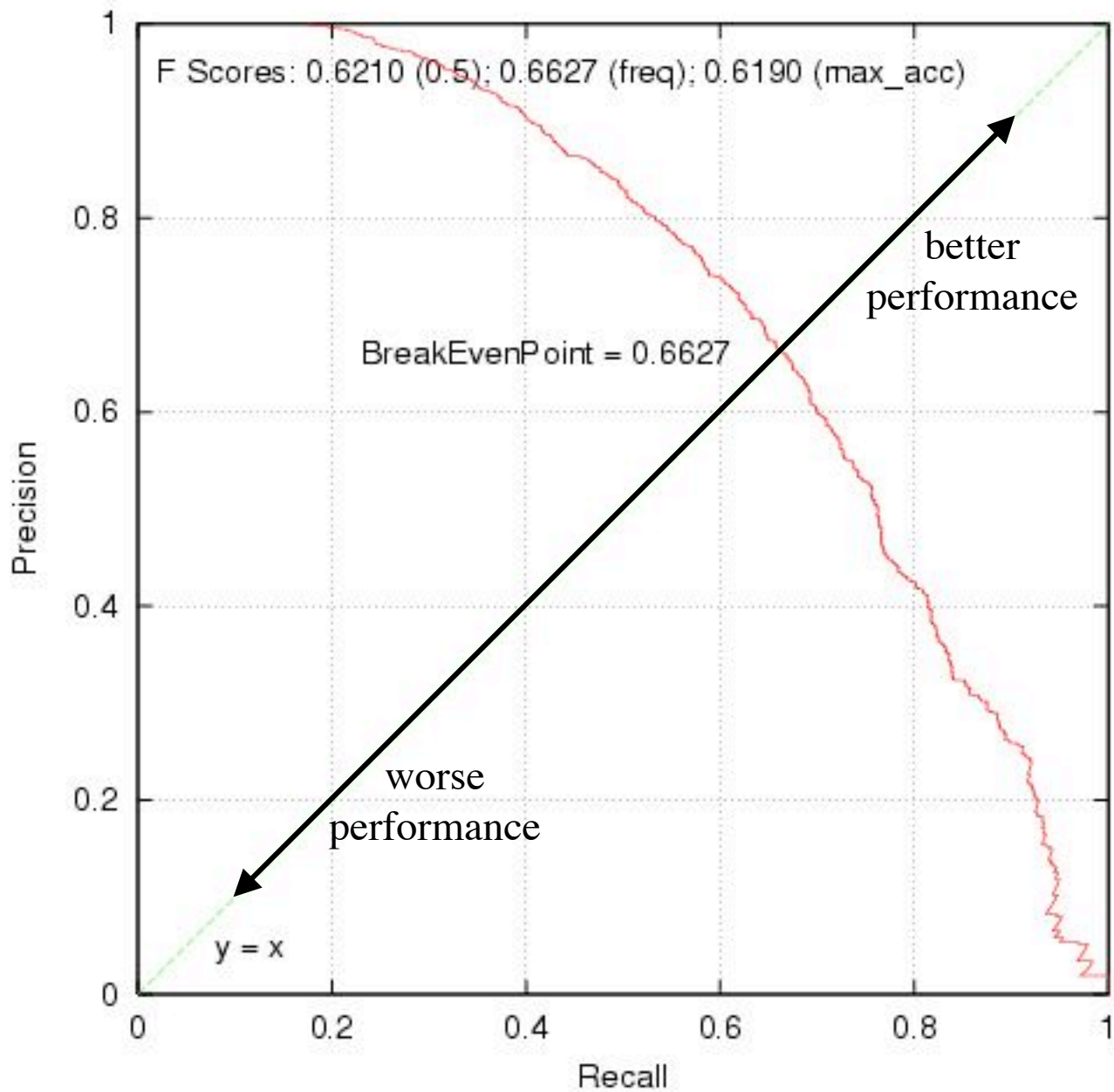
$$RECALL = a / (a + b)$$

$$F = \frac{2 * (PRECISION * RECALL)}{(PRECISION + RECALL)}$$

harmonic average of
precision and recall

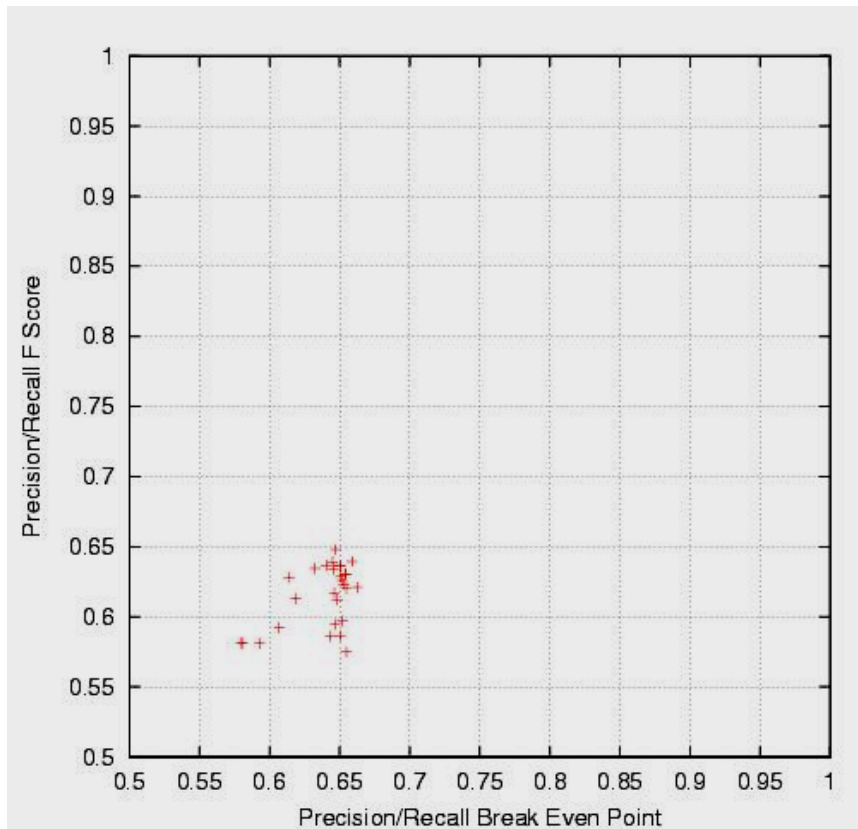


$$BreakEvenPoint = PRECISION = RECALL$$

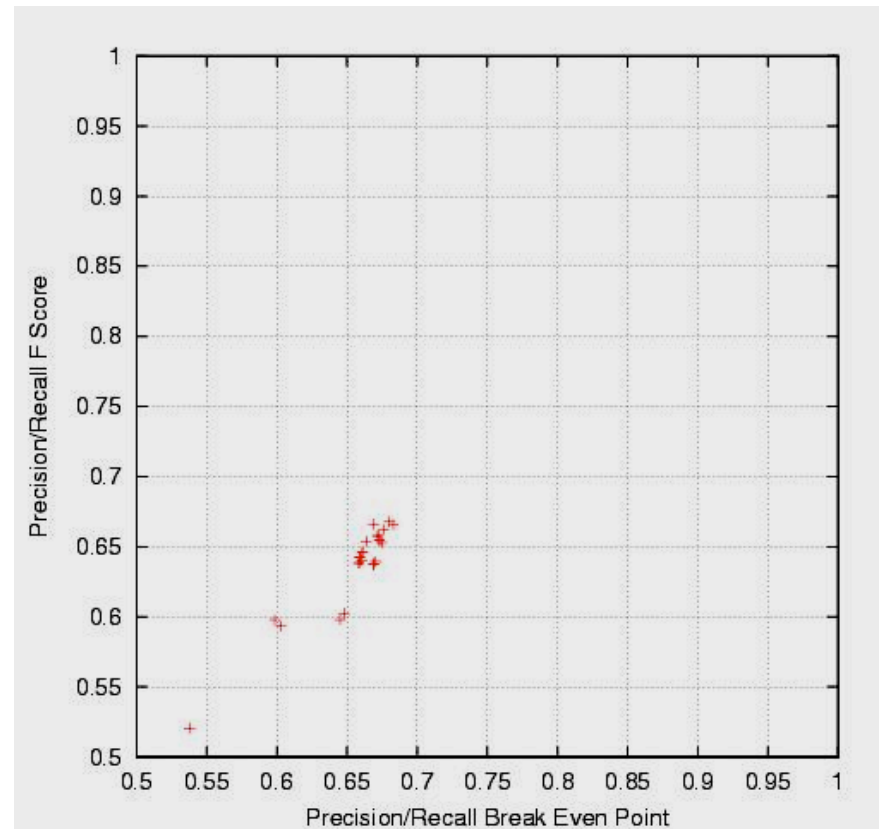


F and BreakEvenPoint do not always correlate well

Problem 1



Problem 2



	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	true positive	false negative
<u>True 0</u>	false positive	true negative

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	TP	FN
<u>True 0</u>	FP	TN

	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	hits	misses
<u>True 0</u>	false alarms	correct rejections

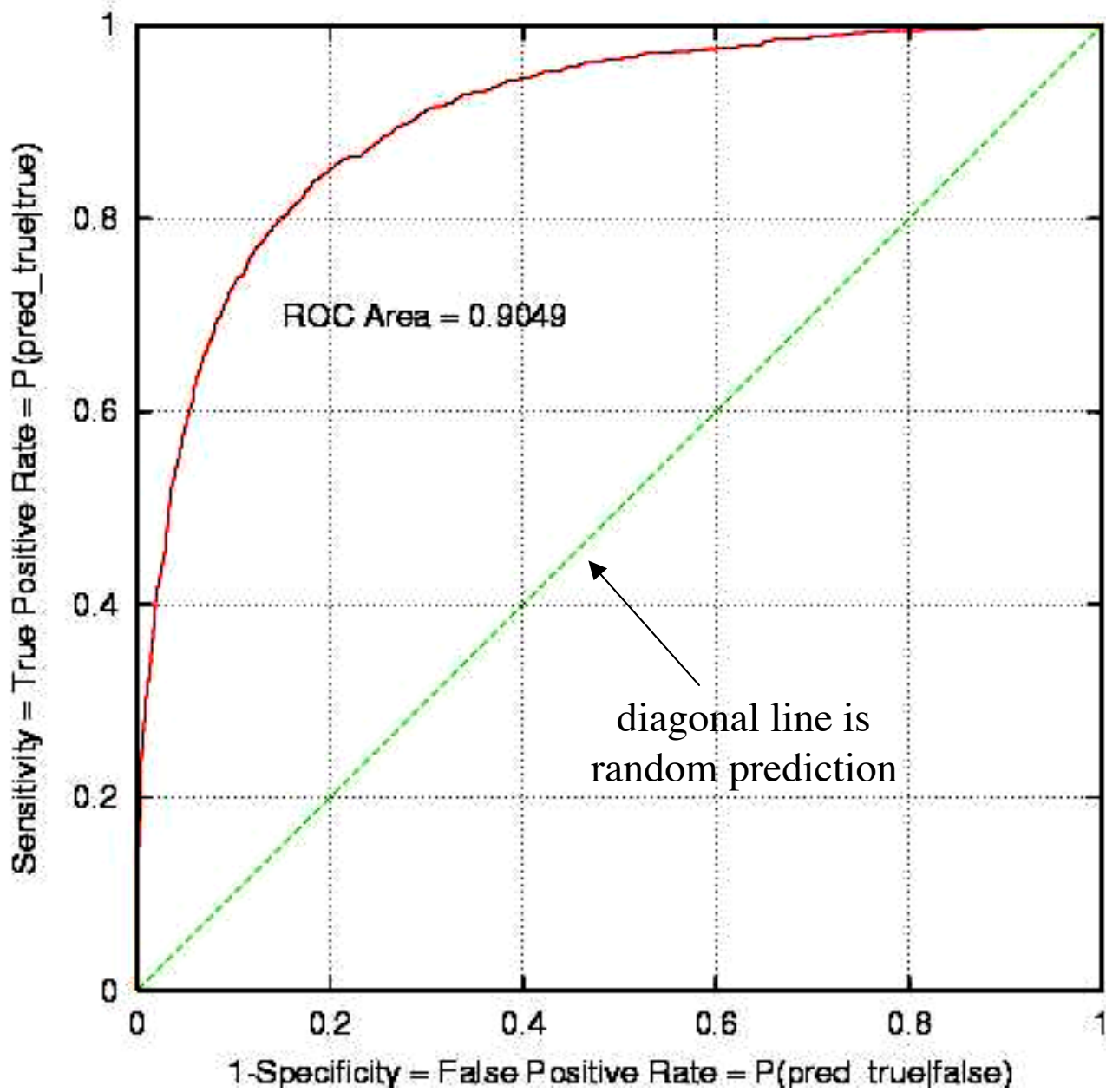
	<u>Predicted 1</u>	<u>Predicted 0</u>
<u>True 1</u>	$P(\text{pr1 tr1})$	$P(\text{pr0 tr1})$
<u>True 0</u>	$P(\text{pr1 tr0})$	$P(\text{pr0 tr0})$

ROC Plot and ROC Area

- Receiver Operator Characteristic
- Developed in WWII to statistically model false positive and false negative detections of radar operators
- Better statistical foundations than most other measures
- Standard measure in medicine and biology
- Becoming more popular in ML

ROC Plot

- Sweep threshold and plot
 - TPR vs. FPR
 - Sensitivity vs. 1-Specificity
 - $P(\text{true}|\text{true})$ vs. $P(\text{true}|\text{false})$
- Sensitivity = $a/(a+b)$ = Recall = LIFT numerator
- 1 - Specificity = $1 - d/(c+d)$



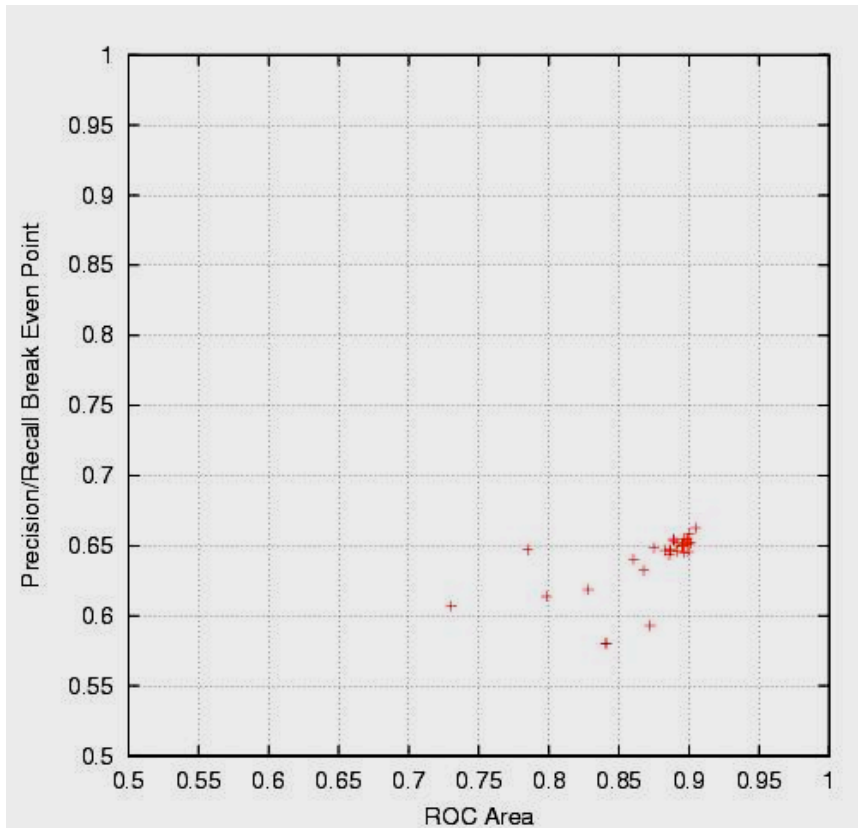
Properties of ROC

- ROC Area:
 - 1.0: perfect prediction
 - 0.9: excellent prediction
 - 0.8: good prediction
 - 0.7: mediocre prediction
 - 0.6: poor prediction
 - 0.5: random prediction
 - <0.5 : something wrong!

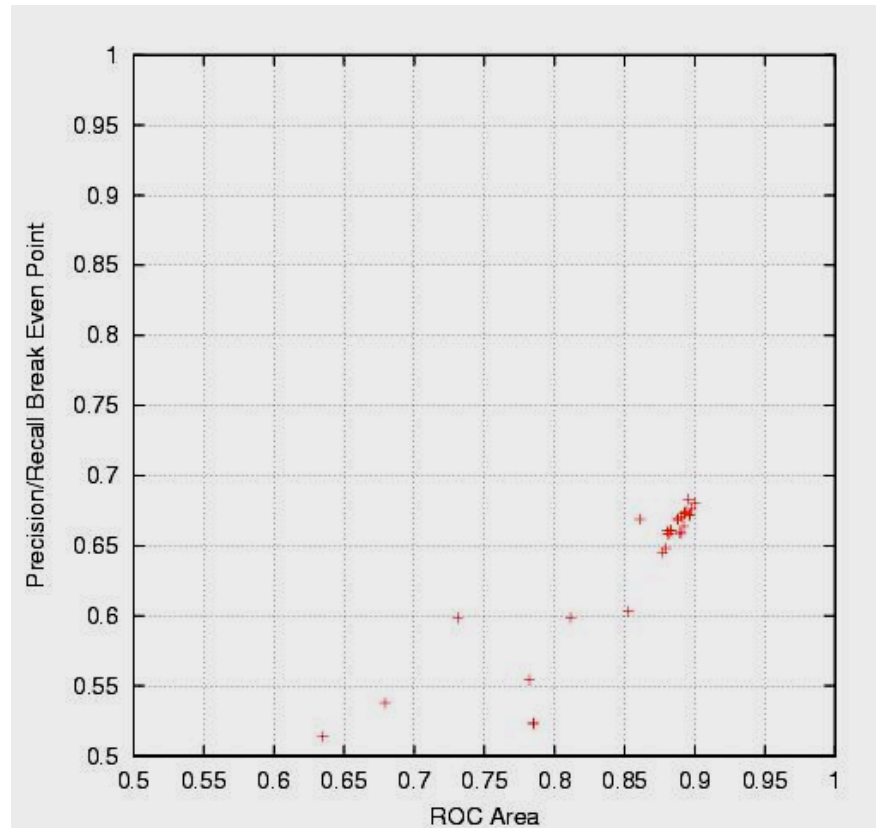
Properties of ROC

- Slope is non-increasing
- Each point on ROC represents different tradeoff (cost ratio) between false positives and false negatives
- Slope of line tangent to curve defines the cost ratio
- ROC Area represents performance averaged over all possible cost ratios
- If two ROC curves do not intersect, one method dominates the other
- If two ROC curves intersect, one method is better for some cost ratios, and other method is better for other cost ratios

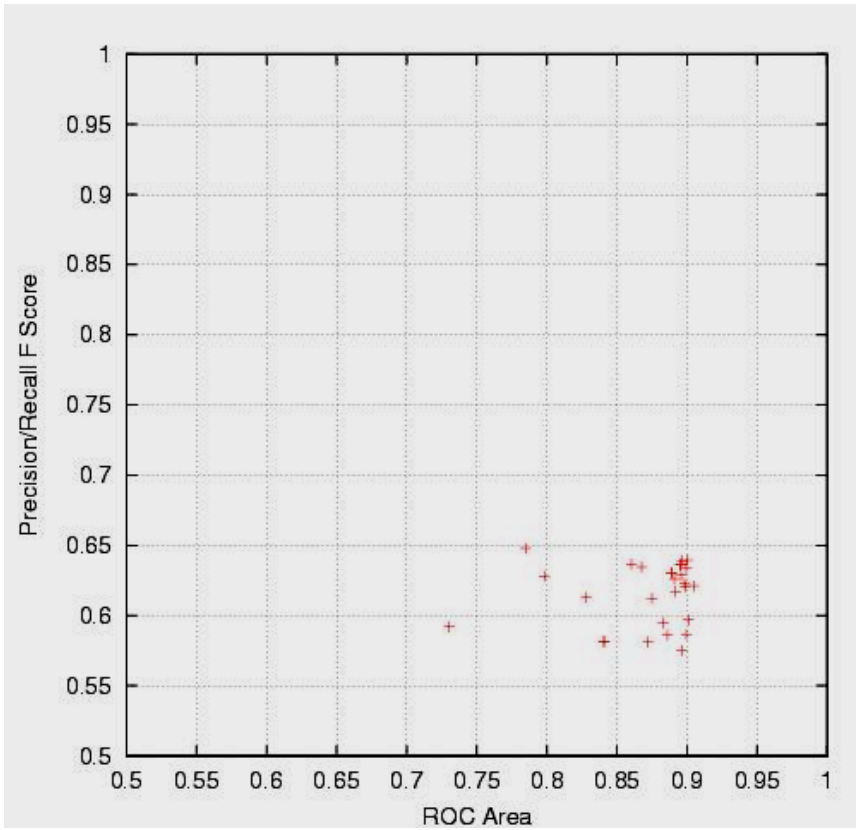
Problem 1



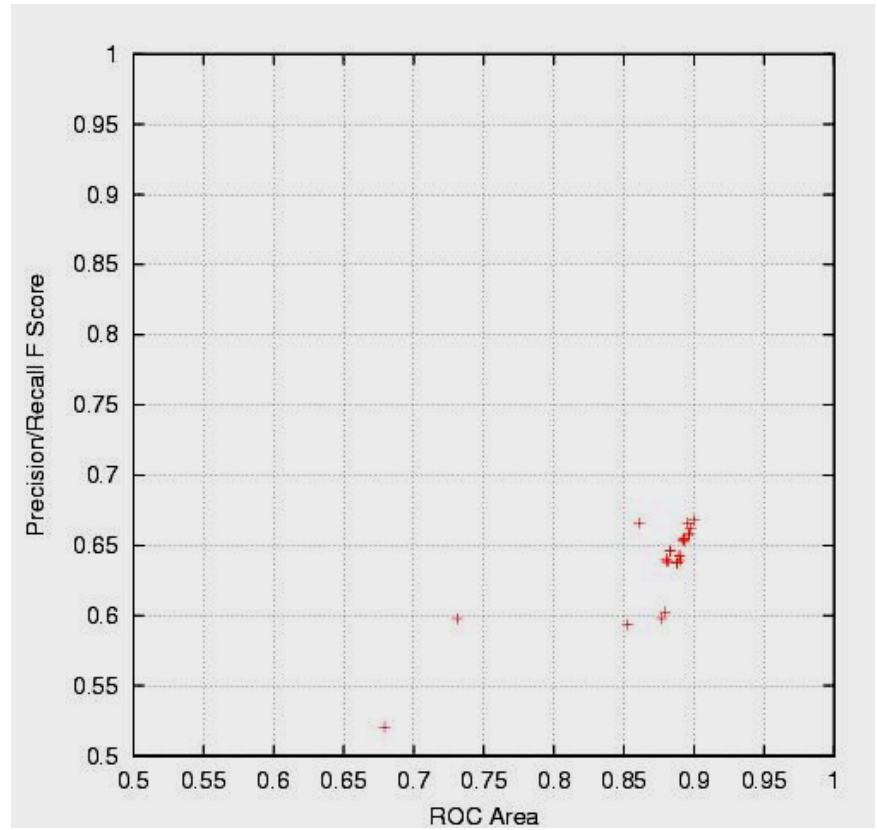
Problem 2



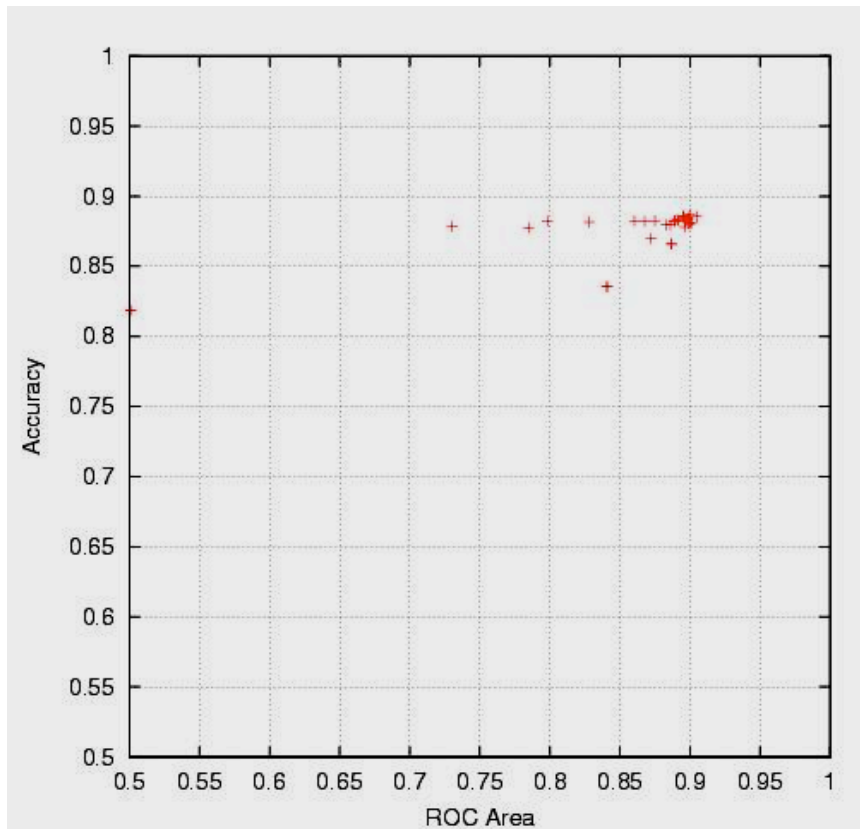
Problem 1



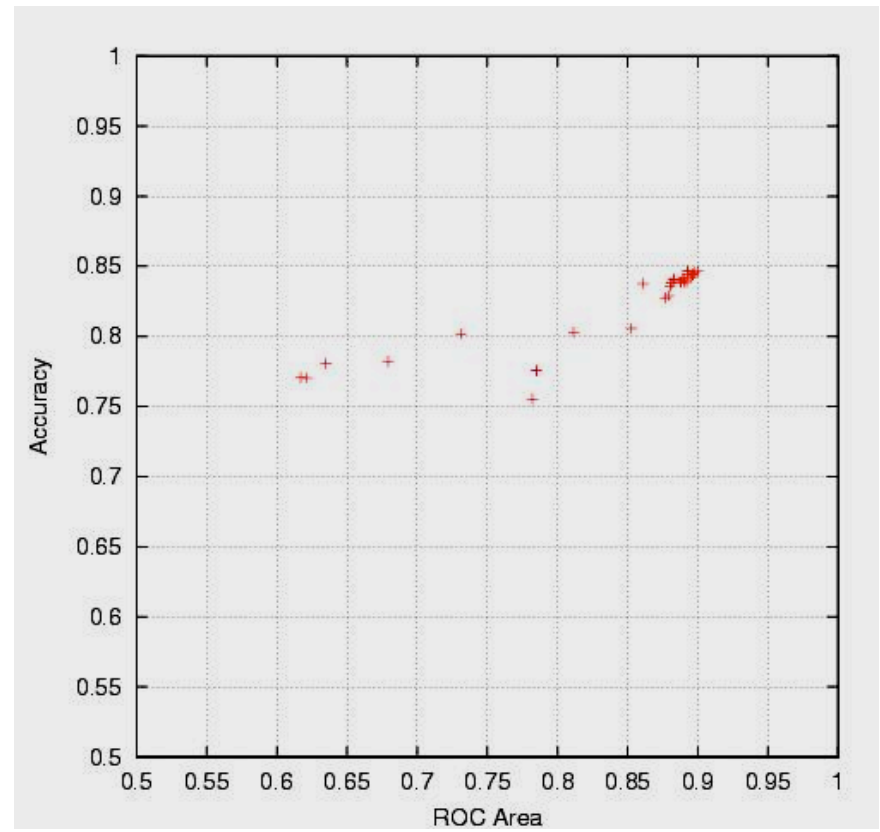
Problem 2



Problem 1



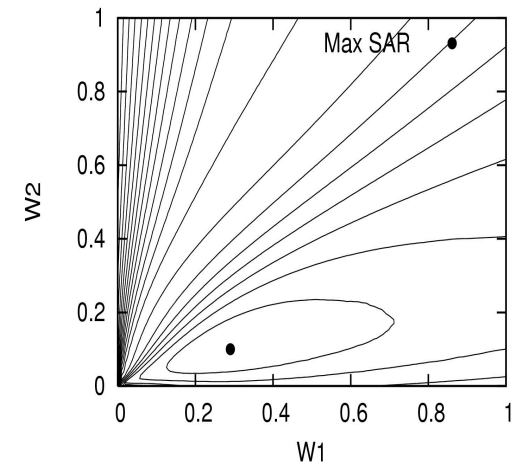
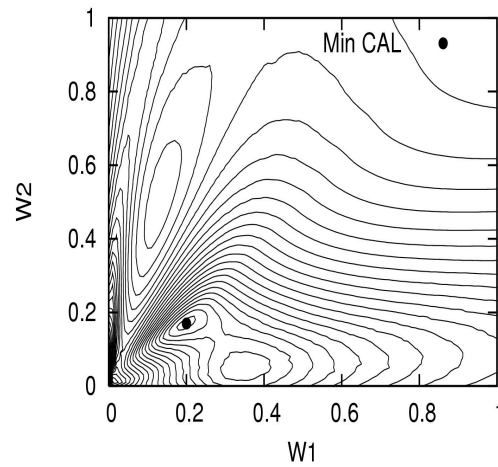
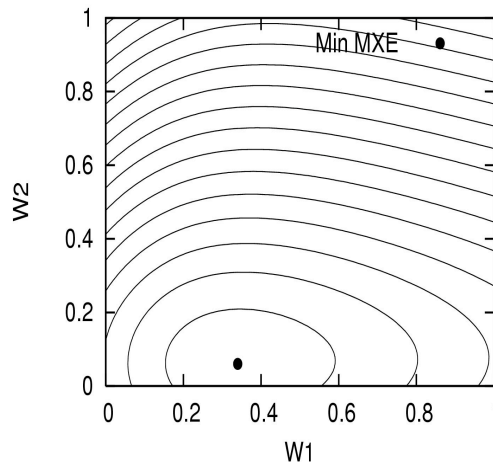
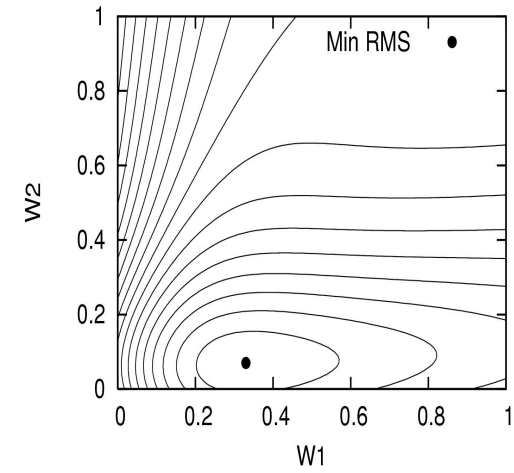
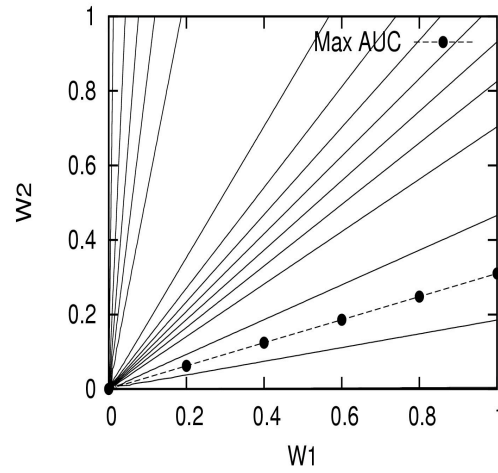
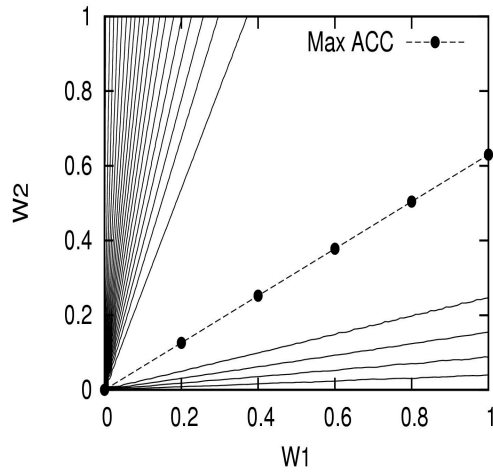
Problem 2



Summary

- the measure you optimize to makes a difference
- the measure you report makes a difference
- use measure appropriate for problem/community
- accuracy often is not sufficient/appropriate
- ROC is gaining popularity in the ML community
- only accuracy generalizes to >2 classes!

Different Models Best on Different Metrics



2-D Multi-Dimensional Scaling

