# Practical ML Advice

Rishabh Iyer

## Announcements

1. Assignment 4 is uploaded (elearning)
   — Nov 30th Due.

2. Finals on Dec 2nd.

3. Course - Evaluation.

# Practical ML

**Proper Experimental Methodology Can Have a Huge Impact:**

*Pure Sciences*
↓

A 2002 paper in *Nature* (a major journal) needed to be corrected due to "training on the testing set"
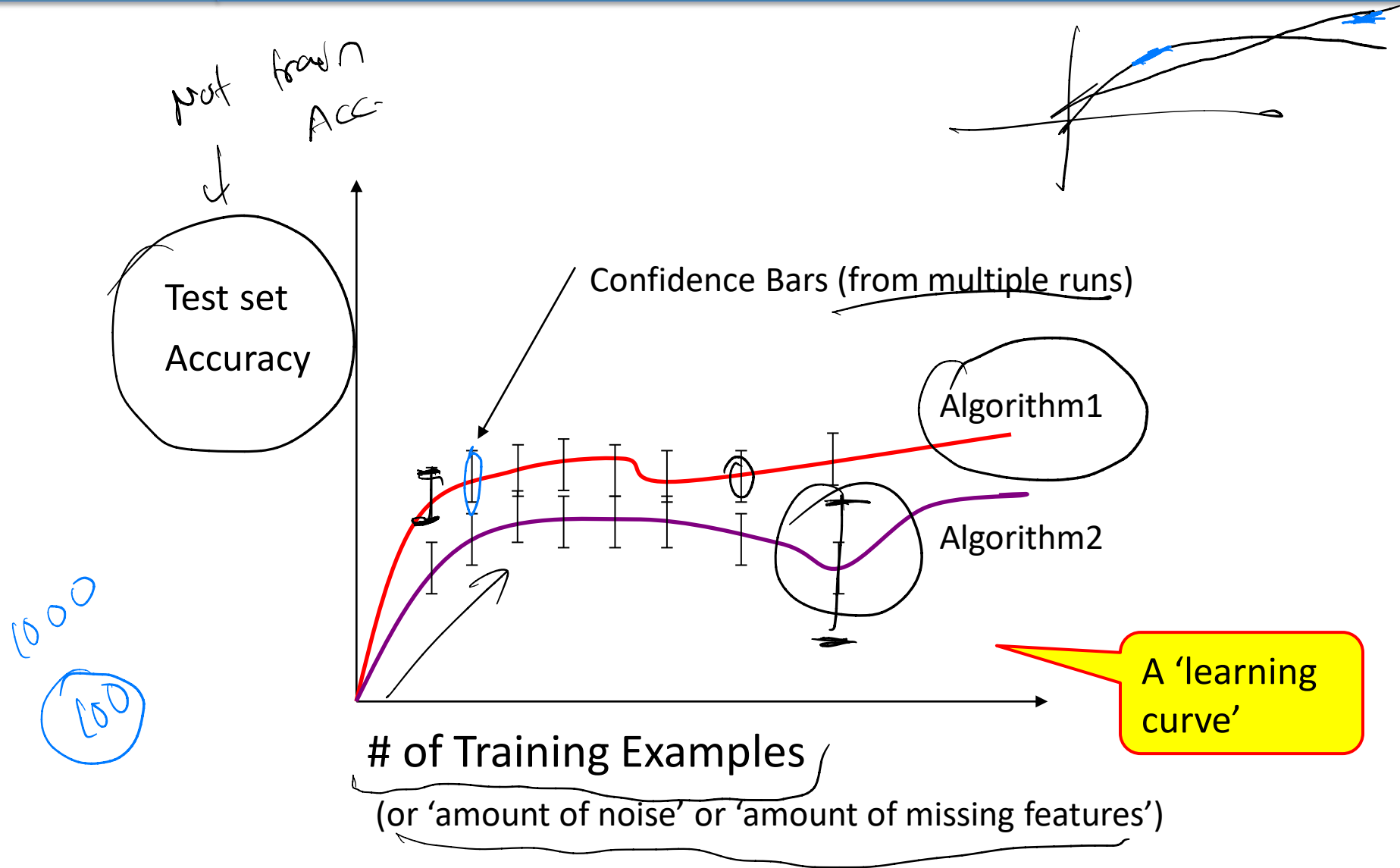
Original report : 95% accuracy (5% error rate)

Corrected report (which still is buggy):

73% accuracy (27% error rate)

Error rate increased over 400%!!!

# Some Typical ML Experiments



Not train Acc.

Test set Accuracy

Confidence Bars (from multiple runs)

Algorithm1

Algorithm2

A 'learning curve'

\# of Training Examples

(or 'amount of noise' or 'amount of missing features')

# Typical Experiments

Ablation Study

| | Test Set Performance |
|---|---|
| Full System $(A, B, C, ...)$ | 80% |
| Without Module A | 75% |
| Without Module B | 62% |
| Without Module C | 29% |

# Experimental Methodology

1) Start with a dataset of labeled examples

2) Randomly partition into *N* groups

3a) *N* times, combine *N* -1 groups into a train set

3b) Provide training set to learning system

3c) Measure accuracy on "left out" group (the test set)

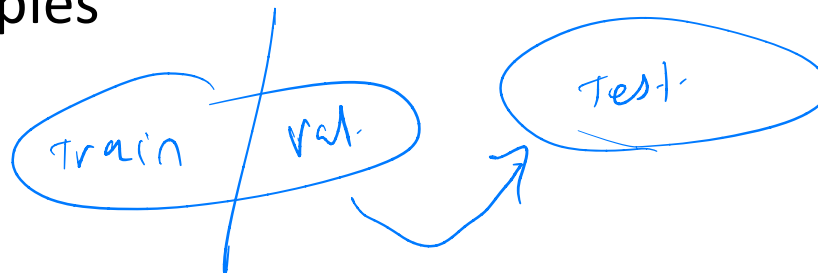| train | test | train | train |
|-------|------|-------|-------|

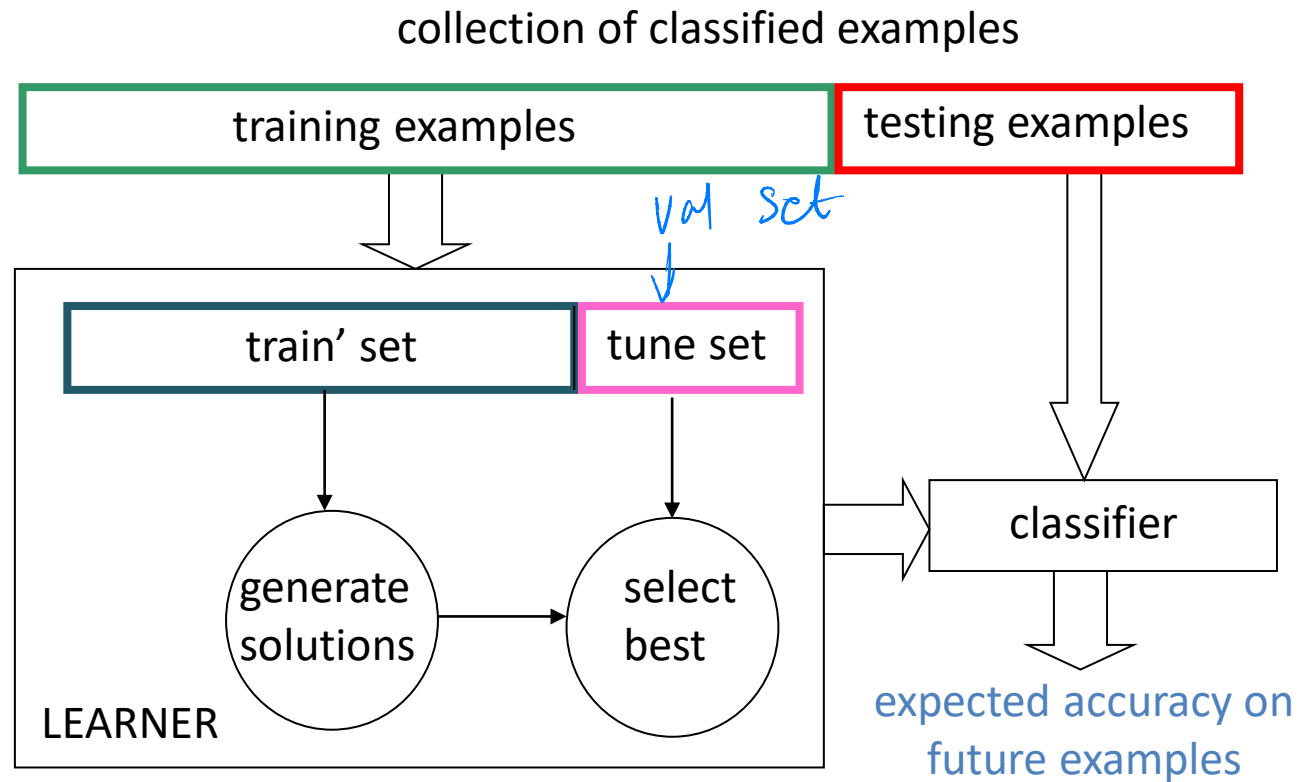Called *N*-fold cross validation

✗ 4 times.

# Validation Sets

- Often, an ML system has to choose when to stop learning, select among alternative answers, etc.

- One wants the model that produces the highest accuracy on **future** examples ("overfitting avoidance")

- It is a **"cheat"** to look at the **test** set while still learning

- Better method
  - Set aside part of the training set
  - Measure performance on this validation data to estimate future performance for a given set of hyperparameters
  - Use best hyperparameter settings, train with **all** training data (except **test** set) to estimate future performance on **new** examples

# A typical Learning system

collection of classified examples

Statistical techniques such as 10-fold cross validation and *t*-tests are used to get meaningful results

training examples | testing examples

*val set*

train' set | tune set

generate solutions → select best

LEARNER

classifier

expected accuracy on future examples

# Multiple Tuning sets

- Using a **single** tuning set can be unreliable predictor, plus some data "wasted"

  1) For each possible set of hyperparameters

     a) Divide <u>training</u> data into **train** and **valid.** sets, using **N-fold cross validation**

     b) Score this set of hyperparameter values:  average **valid.** set accuracy over the $N$ folds

  2) Use **best** set of hyperparameter settings and **all** (train + valid.) examples

  3) Apply resulting model to **test** set

# EVALUATING ML MODELS

# Contingency Tables

(special case of 'confusion matrices')

$$Acc = \frac{n(1,1) + n(0,0)}{n}$$

**True Answer**

|  | + | - |
|---|---|---|
| **+** | n(1,1) [true pos] | n(1,0) [false pos] |
| **-** | n(0,1) [false neg] | n(0,0) [true neg] |

**Algorithm Answer**

$$n = n(1,1) + n(1,0) + n(0,1) + n(0,0)$$

Counts of occurrences

# TPR and FPR

True

$$\text{Pred} \quad + \quad \begin{array}{|c|c|} \hline n(1,1) & n(1,0) \\ \hline n(0,1) & n(0,0) \\ \hline \end{array} \begin{array}{c} + \\ \\ - \end{array}$$

**True Positive Rate** $\quad = \quad n(1,1) \,/\, (\,n(1,1) \,+\, n(0,1)\,)$

(TPR) $\qquad\qquad\qquad = \;$ correctly categorized +'s / total positives

$\qquad\qquad\qquad\quad \sim \;$ P(algo outputs + | + is correct)

**False Positive Rate** $\quad = \quad n(1,0) \,/\, (\,n(1,0) \,+\, n(0,0)\,)$

(FPR) $\qquad\qquad\qquad = \;$ incorrectly categorized −'s / total neg's

$\qquad\qquad\qquad\quad \sim \;$ P(algo outputs + | - is correct)

Can similarly define False Negative Rate and True Negative Rate

$$\downarrow \qquad\qquad\qquad\qquad\qquad \downarrow$$

$$\frac{n(0,1)}{n(1,1) + n(0,1)} \qquad\qquad \frac{n(0,0)}{n(1,0) + n(0,0)}$$
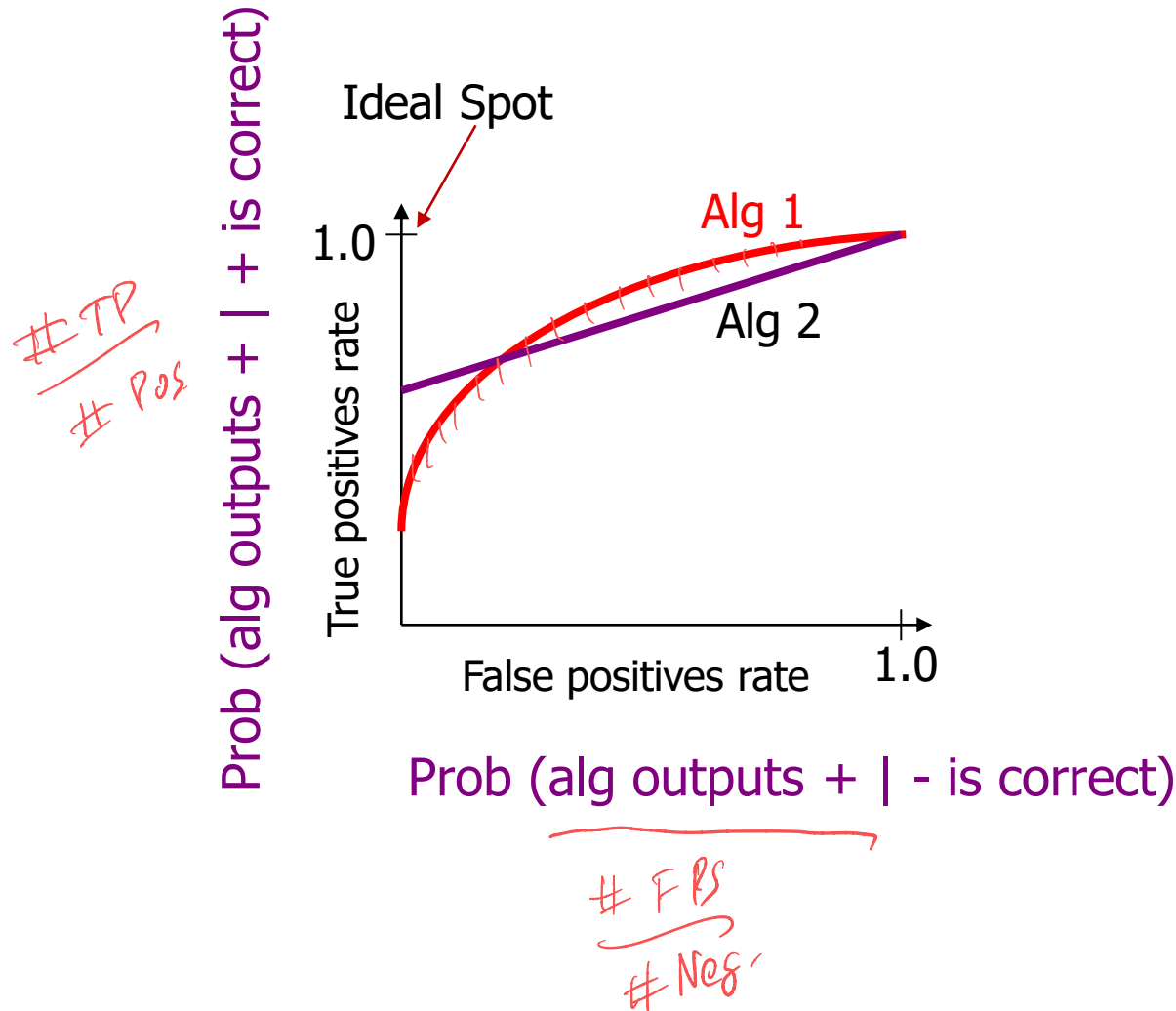
# ROC Curves

*10% Positives     90% Negatives*

- ROC: *Receiver Operating Characteristics*

- Started for radar research during WWII

- Judging algorithms on accuracy alone may not be good enough when **getting a positive wrong** costs more than **getting a negative wrong** (or vice versa)
  - e.g., medical tests for serious diseases
  - e.g., a movie-recommender system

*search lAds*

*Algo 1:   100% Negatives (Acc: 90%)*

*Algo 2:   5% pos, 95% Neg.*

# ROC Curves Graphically

Ideal Spot

**Prob (alg outputs + | + is correct)**

True positives rate

1.0

Alg 1

Alg 2

False positives rate

1.0

**Prob (alg outputs + | - is correct)**

$\frac{\# TP}{\# Pos}$

$\frac{\# FPs}{\# Neg}$

Different algorithms can work better in different parts of ROC space. This depends on cost of false + vs false -

**The Standard Approach:**

- You need an ML algorithm that outputs NUMERIC results such as prob(example is +)   *Alg Score ∈ [0,1]*

- You can use ensemble methods to get this from a model that only provides Boolean outputs
  - e.g., have 100 models vote & count votes

**Step 1**: Sort predictions on test set

**Step 2**: Locate a *threshold* between examples with opposite categories

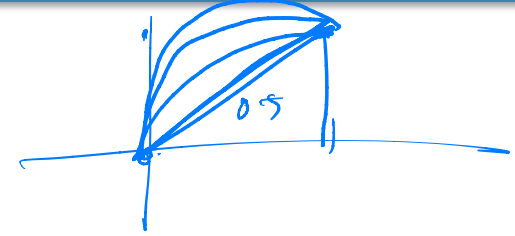**Step 3**: Compute TPR & FPR for each threshold of Step 2

**Step 4**: Connect the dots



Example ROC Curve
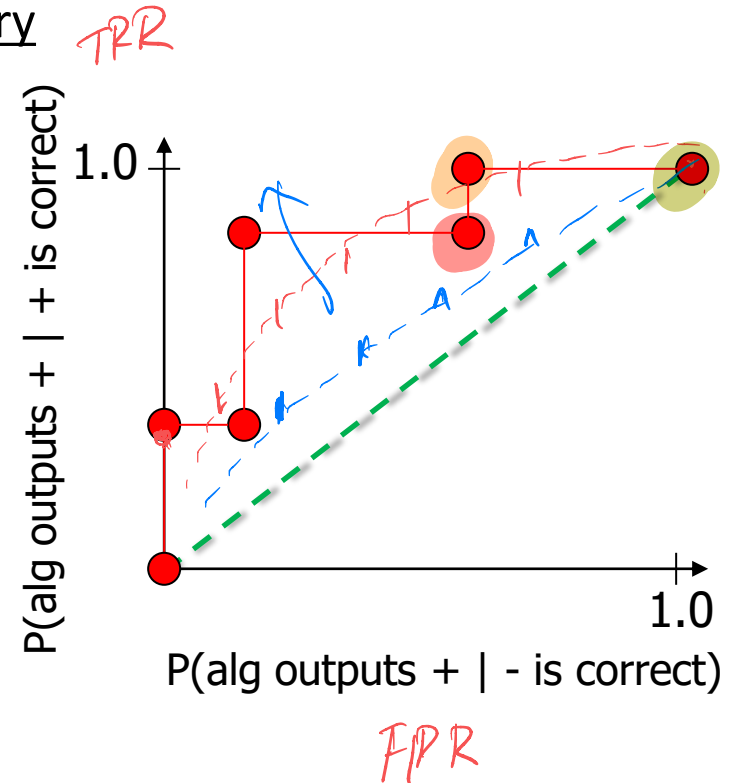
# Plotting ROC Curves - Example

## ML Algo Output (Sorted)          Correct Category

| | | | |
|---|---|---|---|
| Ex 9 | .99 | | + |
| Ex 7 | .98 | TPR=(2/5), FPR=(0/5) | + |
| Ex 1 | .72 | TPR=(2/5), FPR=(1/5) | - |
| Ex 2 | .70 | | + |
| Ex 6 | .65 | TPR=(4/5), FPR=(1/5) | + |
| Ex 10 | .51 | | - |
| Ex 3 | .39 | TPR=(4/5), FPR=(3/5) | - |
| Ex 5 | .24 | TPR=(5/5), FPR=(3/5) | + |
| Ex 4 | .11 | | - |
| Ex 8 | .01 | TPR=(5/5), FPR=(5/5) | - |

0.25

0.12

0

Algorithm predicts + if its output is ≥ thresh.

TPR

P(alg outputs + | + is correct)

1.0

P(alg outputs + | - is correct)

1.0

FPR

# Area Under ROC Curve

- A common metric for experiments is to numerically integrate the ROC Curve

  - Usually called AUC

  - Probability that ML alg. will "rank" a randomly chosen positive instance higher than a randomly chosen negative one

  - Can summarize the curve **too much** in practice

True positives (y-axis), False positives (x-axis), axis labeled 1.0 on both

$$AUC = Prob\left(Score_A(Rand\ +) \geq Score_A(Rand\ -)\right)$$

# Asymmetric Error Costs

- Assume that cost(FP) ≠ cost(FN)

- You would like to pick a threshold that minimizes

$$E(total\ cost)$$
$$= cost(FP) \times \text{pr}(FP) \times (\#\ of\ neg\ ex's) +$$
$$cost(FN) \times \text{pr}(FN) \times (\#\ of\ pos\ ex's)$$

- You could also have (maybe negative) costs for TP and TN (assumed zero in above)

# ROC's & Skewed Data

- One strength of ROC curves is that they are a good way to deal with skewed data (|+| >> |-|) since the axes are fractions (rates) independent of the # of examples

- You must be careful though!

  - Low FPR * (many negative ex) = sizable number of FP

  - Possibly more than # of TP

$$TPR = \frac{TP}{\#POS}$$

$$FPR = \frac{FP}{\#Neg.}$$

# Precision vs. Recall (PR Curve)

- Think about search engines...

- **Precision** = (# of relevant items retrieved) / (total # of items retrieved)

  = $n(1,1)$ / ( $n(1,1)$ + $n(1,0)$ )

  $\cong$ P(is pos | called pos)

- **Recall** = (# of relevant items retrieved) / (# of relevant items that exist)

  = $n(1,1)$ / ( $n(1,1)$ + $n(0,1)$ ) = TPR

  $\cong$ P(called pos | is pos)

- Notice that n(0,0) is not used in either formula
  Therefore you get <u>no</u> credit for filtering out <u>ir</u>relevant items
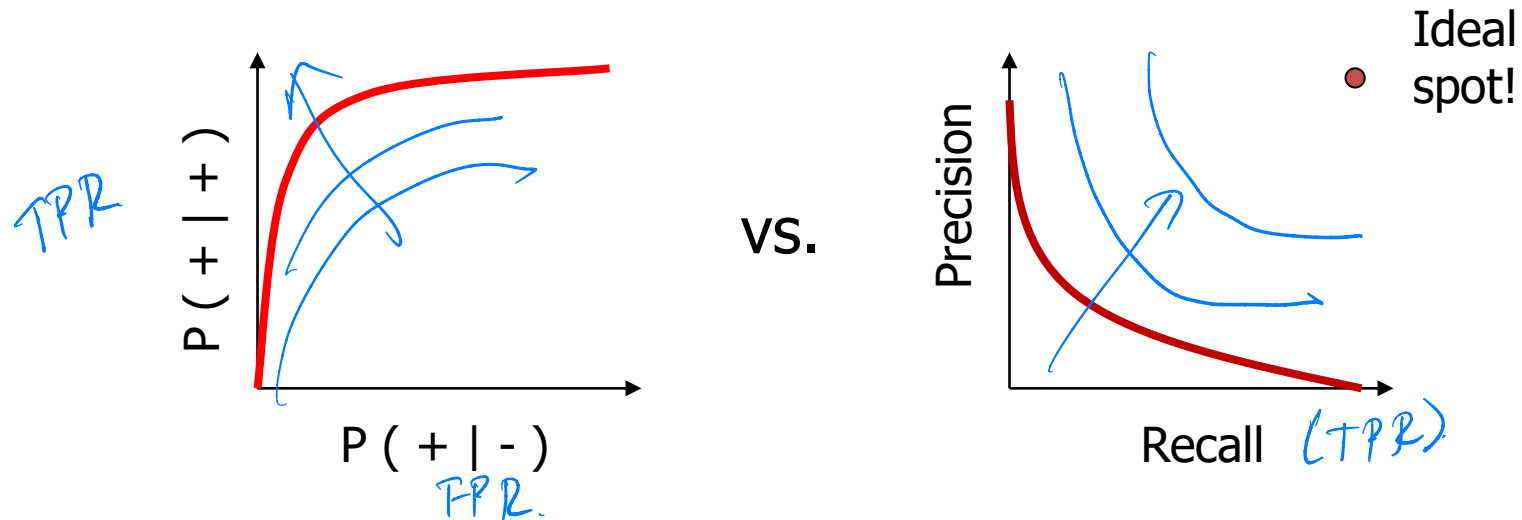
$$\frac{TP}{TP + FP} = Precision$$

$$\frac{TP}{TP + FN} = Recall = TPR$$

$$TN$$

$$FPR = \frac{FP}{FP + TN}$$

$$(TN)$$

# ROC vs. Precision-Recall

You can get very different visual results
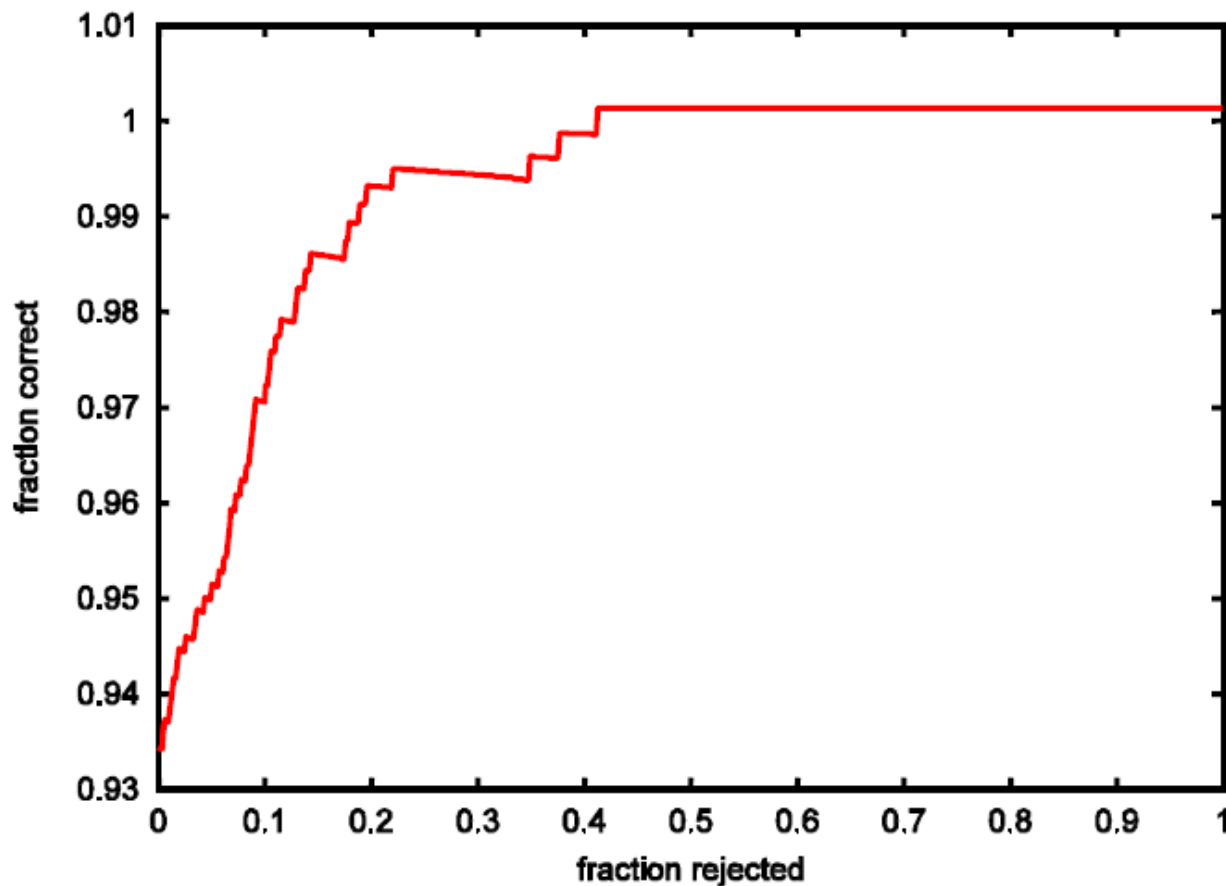on the same data!



vs.

Ideal
spot!

The reason for this is that there may be lots of – ex's
(e.g., might need to include 100 neg's to get 1 more pos)

# Rejection Curves

- In most learning algorithms, we can specify a threshold for making a rejection decision

  - Probabilistic classifiers: adjust cost of rejecting versus cost of FP and FN

  - Decision-boundary method: if a test point $x$ is within $\theta$ of the decision boundary, then reject

    - Equivalent to requiring that the "activation" of the best class is larger than the second-best class by at least $\theta$

# Rejection Curves

- Vary θ and plot fraction correct versus fraction rejected

# The F1 Measure

- Figure of merit that combines precision and recall

(Assumes you have threshold)

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

where $P$ = precision; $R$ = recall. This is twice the harmonic mean of $P$ and $R$.

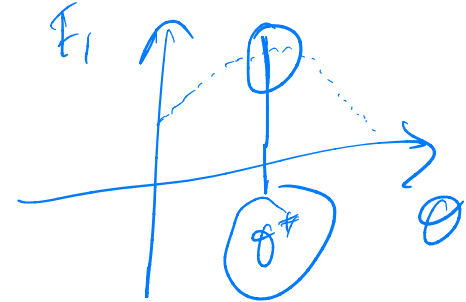- We can plot $F1$ as a function of the classification threshold $\theta$

- Figure of merit that combines precision and recall

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

where $P$ = precision; $R$ = recall. This is twice the harmonic mean of $P$ and $R$.

- We can plot $F1$ as a function of the classification threshold $\theta$

# Multi-Class

$F_1(y)$, $\forall y$.

$Avg-F = Avg\left( F_1(y), \forall y \right)$

$1, 2, 3$

$F_1(1) : F_1( 1 \text{ v/s } 2-3 )$

$F_1(2) : F_1( 2 \text{ v/s } 1-3 )$

$F_1(3) = F_1( 3 \text{ v/s } 1-2 )$

# The F1 Measure

- Figure of merit that combines precision and recall

$$F_1 = 2 \cdot \frac{P \cdot R}{P + R}$$

  where $P$ = precision; $R$ = recall. This is twice the harmonic mean of $P$ and $R$.

- We can plot $F1$ as a function of the classification threshold $\theta$