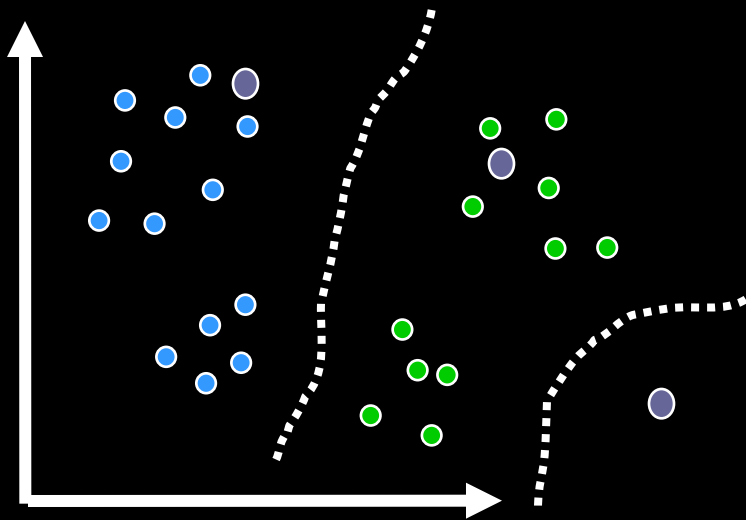


- Start thinking about term project ideas.
- Interesting data set, use of which won the Marr prize in computer vision:
  - <http://vision.cs.stonybrook.edu/~vicente/sbucaptions/>
  - Interesting project?
  - From Larry:
    - [https://lab.nationalmedals.org/img\\_processing](https://lab.nationalmedals.org/img_processing)
- Next 1.5 weeks: Pattern recognition
  - **Concepts, error types (today)**
  - Basic theory and how to use classifiers in MATLAB:
    - Support vector machines (SVM).
    - Neural networks

# Pattern recognition

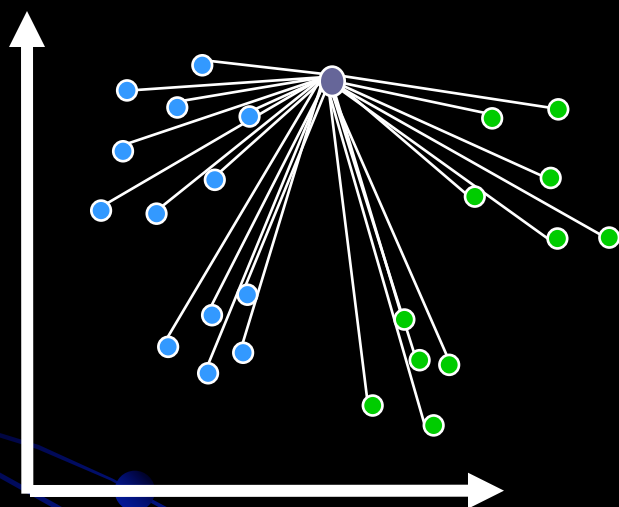


- Making a decision from data
  - A classification problem: *assign a single class label to a datum point*
  - Can include a special class, *reject*,
    - if a sample (a single datum point) appears not to belong to any known class
    - If it is on the boundary between classes
  - Else *forced classification*
- Boundaries between classes-how?
- There's tons of theory, can be applied to many areas. We focus on small subset of those used for vision

# Baseline: Hand-tuned decision boundaries

- You did this based on observations for fruit classification
- You'll do the same thing in Lab 4 for shapes
- But what if the features were much more complex?
  - We now discuss classifiers that learn class boundaries based on *exemplars* (e.g., labeled training examples)

# Ex: Nearest neighbor classifier



- Assumes we have a feature vector for each image
- Calculate distance from *new test sample* to each *labeled training sample*.
- Assign label as closest training sample
- Generalize by assigning same label as the majority of the  $k$  nearest neighbors. No majority?

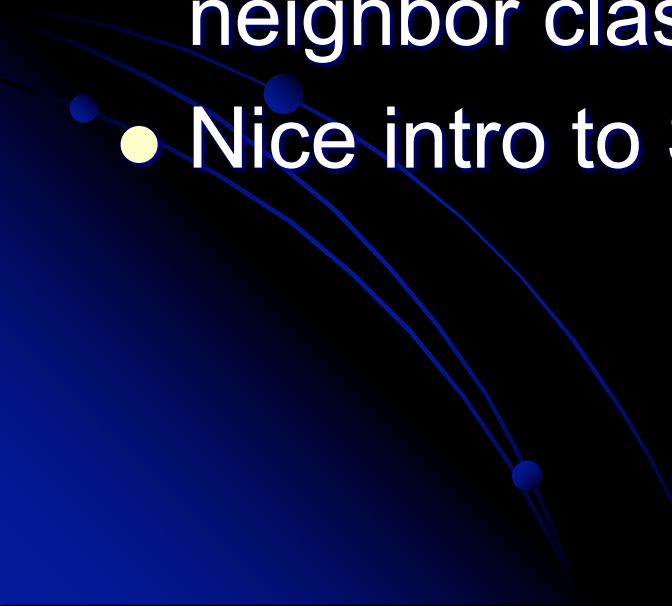
$$\text{In } 2D, \|p_1 - p_2\| = \sqrt{(p_1(x) - p_2(x))^2 + (p_1(y) - p_2(y))^2}$$

$$\text{In } dD, \|p_1 - p_2\| = \sqrt{\sum_{i=1}^d (p_1(i) - p_2(i))^2}$$

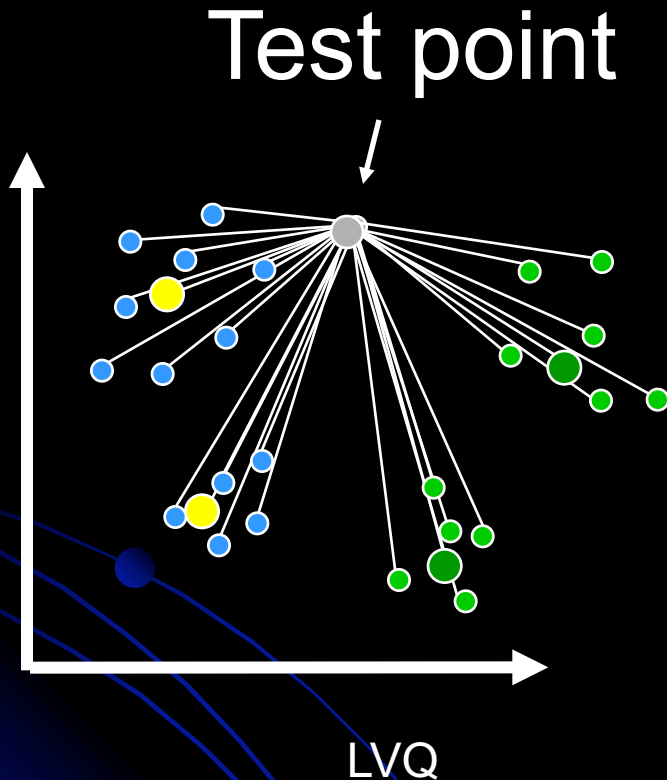
Boundaries:

<http://ai6034.mit.edu/fall12/index.php?title=Demonstrations>

# Look at this to understand nearest neighbor

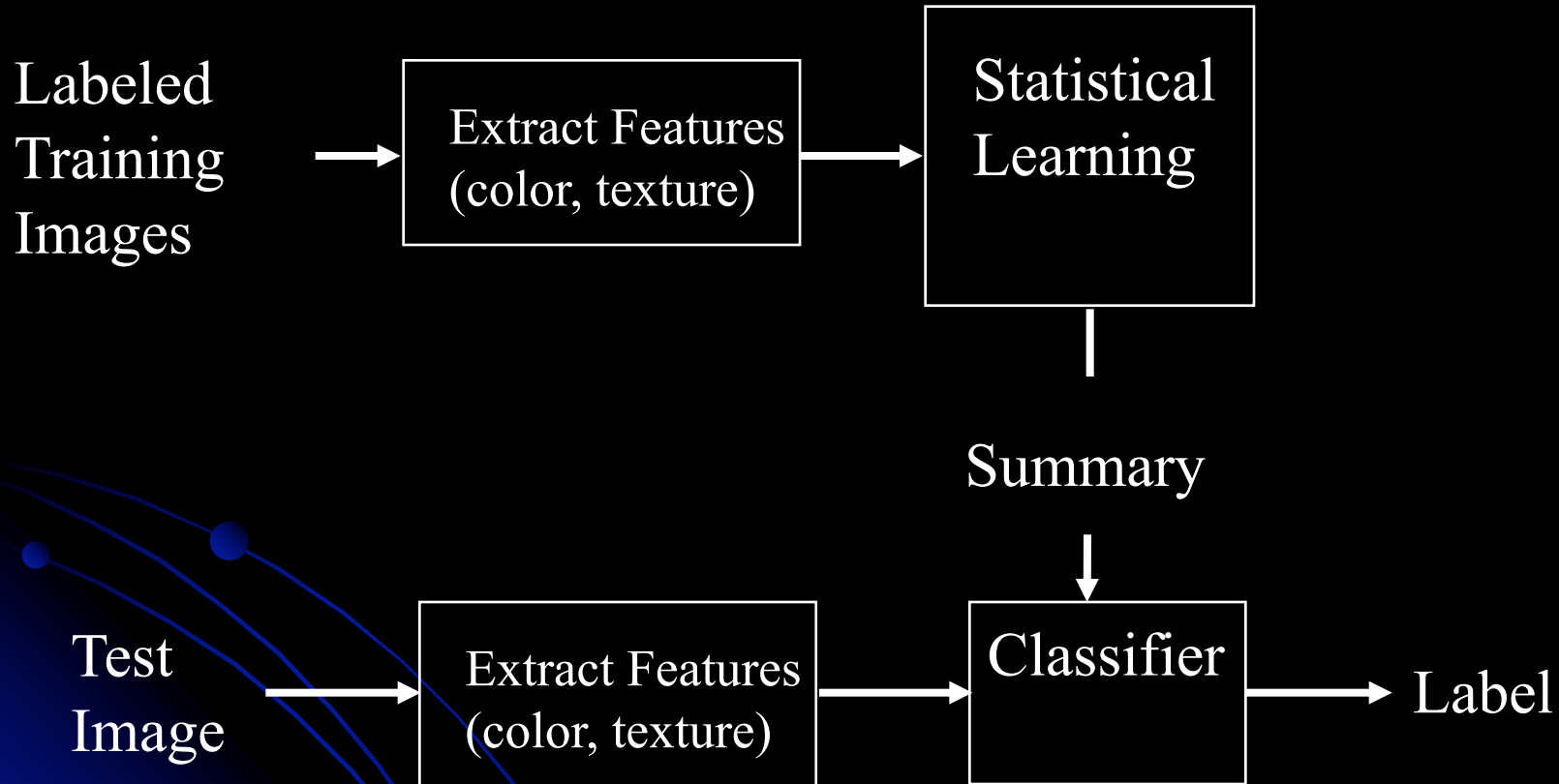
- <http://ai6034.mit.edu/fall12/index.php?title=Demonstrations>
  - Shows Voronoi diagrams for nearest neighbor classifiers
  - Nice intro to SVMs also
- 

# Nearest class mean



- Find class means and calculate distance to each mean
  - Pro?
  - Con?
- Partial solution: clustering
  - Learning vector quantization (LVQ): tries to find optimal clusters

# Common model of learning machines



# How good is your classifier?

Detected \ True	Yes	No
Yes	500 (true pos.)	100 (false neg.)
No	200 (false pos.)	10000 (true neg.)

600  
Total actual positive

10200  
Total actual negative

700  
Total det. as pos.

10100  
Total det. as neg.

- Example from medicine: Disease detection
- Consider costs of false neg. vs. false pos.
- *Lots* of different error measures
  - Accuracy =  $10500/10800 = 97\%$ . Is 97% accuracy OK?
  - Recall (or true positive rate) =  $500/600=83\%$
  - Precision =  $500/700=71\%$
  - False pos rate =  $200/10200 = 2\%$



# How good is your classifier?

- Write out definitions of each measure now

Detected:	Yes	No
Has:		
Yes	500 (true pos.)	100 (false neg.)
No	200 (false pos.)	10000 (true neg.)

- Examples
  - Accuracy =  $10500/10800 = 97\%$ .
  - Recall (or true positive rate) =  $500/600=83\%$
  - Precision =  $500/700=71\%$
  - False pos rate =  $200/10200 = 2\%$

# Thresholding real-valued output allows you to tradeoff TPR and FPR

Simple example: Classes P = positive, N = negative, and single real-valued output.

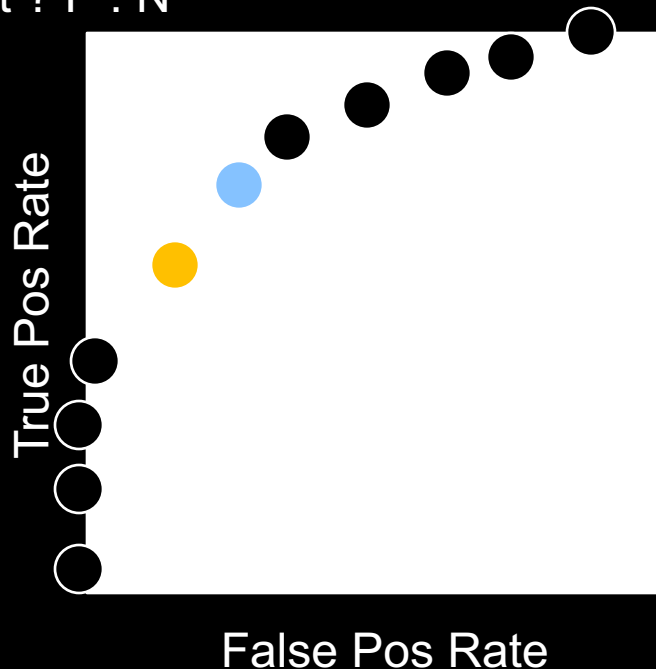


Threshold output to get class. label = output > t ? P : N  
Choice of threshold a

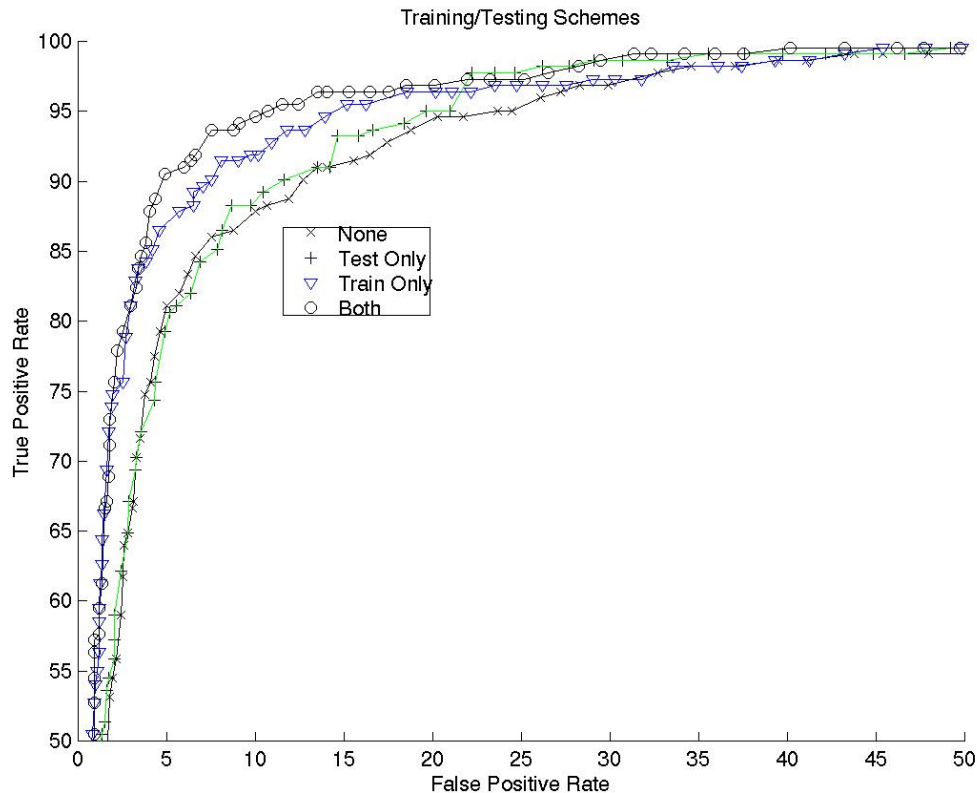
If  $t == 0$ : TPR =  $\frac{9}{12}$ , FPR =  $\frac{2}{8}$

If  $t == 1$ : TPR =  $\frac{6}{12}$ , FPR =  $\frac{4}{8}$

Repeat for many values of t



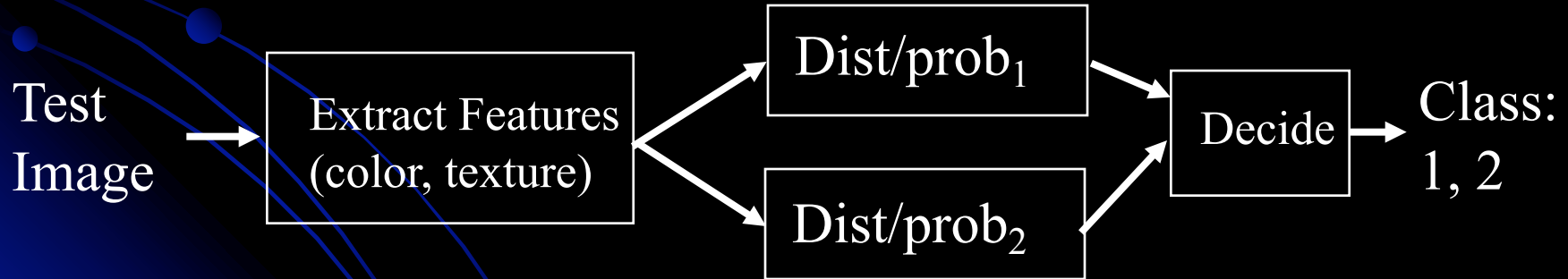
# ROC curve



- Receiver-operating characteristic
- Useful when you can change a threshold to get different true and false positive rates
- Consider extremes
- Much more information recorded here!

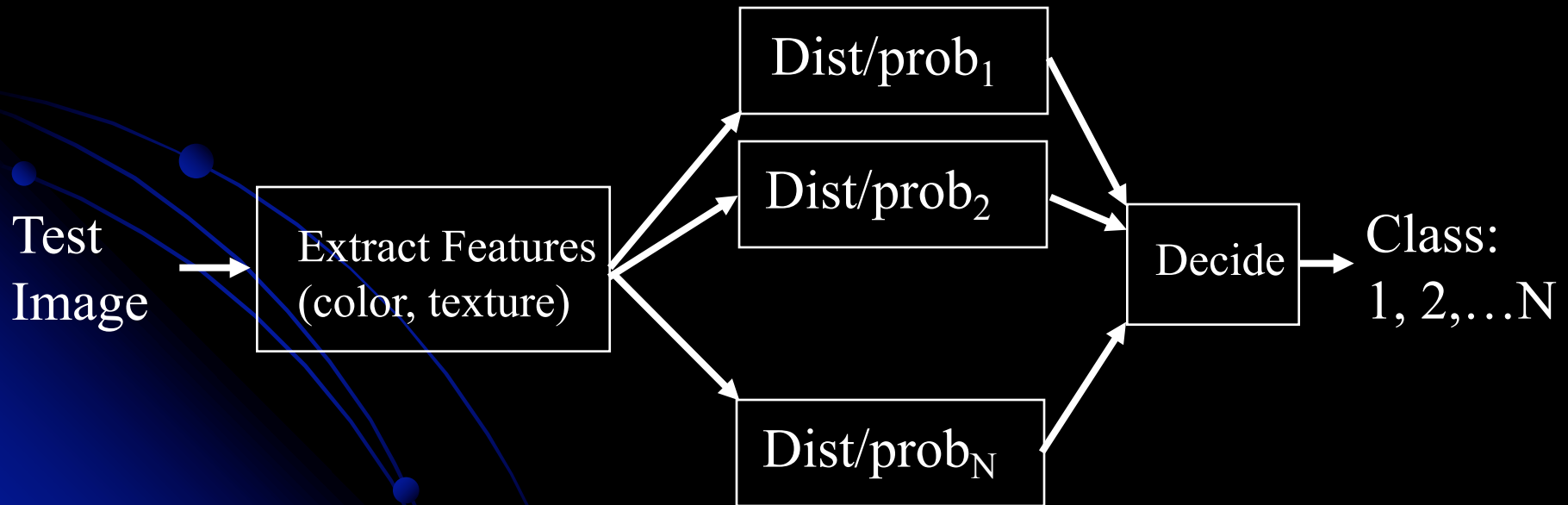
# Focus on testing

- Let  $m$  = the number of possible class labels
- Consider  $m=2$ .
- Example: Calculate distance to cluster means for 2 classes.



# Multiclass problems

- Consider  $m > 2$ .
- Example: Calculate distance to cluster means for 10 classes.



# Confusion matrices for $m > 2$ (outdoor image example)

Detected

	Bch	Sun	FF	Fld	Mtn	Urb
Bch	169	0	2	3	12	14
Sun	2	183	5	0	5	5
FF	3	6	176	6	4	5
Fld	15	0	1	173	11	0
Mtn	11	0	2	21	142	24
Urb	16	4	8	5	27	140

- Beach recall:  $169 / (169 + 0 + 2 + 3 + 12 + 14) = 84.5\%$
- Note confusion between mountain and urban classes due to features
  - Similar colors and spatial layout

# Why do we need separate training and test sets?

## Exam analogy

But working on practice questions is helpful...get the analogy? *We hope our ability to do well on practice questions helps us on the actual exam*

Application to nearest-neighbor classifiers

Often reserve a 3<sup>rd</sup> set for validation as well (to tune parameters of training set)