

# Simple Classification using Binary Data

Deanna Needell  
Mathematics  
UCLA



Thanks:

- NSF CAREER DMS#1348721
- NSF BIGDATA DMS#1740325
- MSRI NSF DMS#1440140
- Alfred P. Sloan Fdn

# Joint work with



Rayan Saab (UCSD)



Tina Woolf (CGU)

# So much data...



So much data...



# So much data...

**METU *siam*. Student Chapter Presents**  
16-24 May 2012

## WHAT WOULD YOU DO WITH ALL THIS DATA?



**Gerhard Wilhelm Weber**  
*Estimation of Dynamics under Uncertainty*  
May 16, Wednesday, 11:40-12:30 \*

**Aybar Acar**  
*MapReduce and Hadoop: Mining Big Data in the Cloud*  
May 23, Wednesday, 14:00-16:30 \*

**Annette Hohenberger**  
*Fractals in Cognitive Science*  
May 17, Thursday, 11:40-12:30 \*

**Cem İyigün**  
*Introduction to Clustering*  
May 23, Wednesday, 11:40-12:30 \*

**Özlem İlk** *Introduction to R and GGobi*  
May 17, Thursday, 14:00-17:30 \*\*

**Fatma Yerlikaya-Özkurt** *Modeling with MARS and CMARS*  
May 24, Thursday, 14:00-15:30 \*

\* Institute of Applied Math, S 209  
\*\* Department of Mathematics, Computer Lab (M 202)  
For more info visit <http://siam.metu.edu.tr>

**Mathematics, Statistics, and the Data Deluge**  
**MATHEMATICS AWARENESS MONTH**


Institute of Applied Math  
<http://www.metu.edu.tr>

Sponsored by the Joint Policy Board for Mathematics - American Mathematical Society - American Statistical Association - Mathematical Association of America - Society for Industrial and Applied Mathematics

# So much data...

Systems to handle big data might be this generation's moon landing

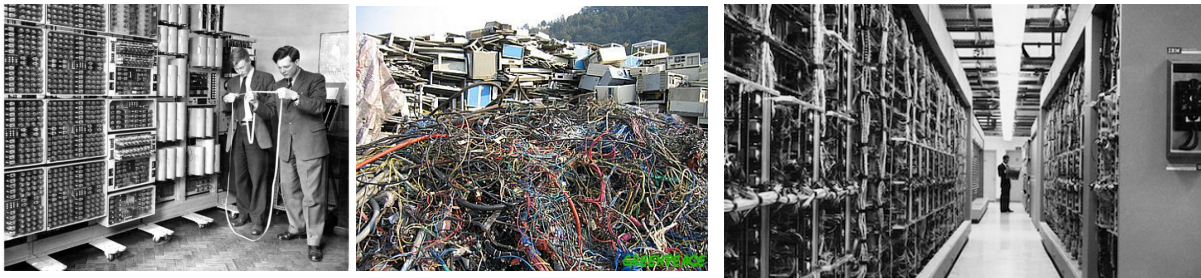
by [Stacey Higginbotham](#)  Apr. 1, 2012 - 9:00 PM PST

 5 Comments

# How can we handle all this data?

Option 1 : Build bigger computing systems

- ❖ We need the resources
- ❖ Fundamental limitations
- ❖ Wasteful (resources, energy, cost, ...)



# How can we handle all this data?



3 MB of internet data transfer = boiling one cup of water

(<https://www.katescomment.com/energy-of-downloads/>)



# How can we handle all this data?

Option 2 : Design more efficient data analysis methods

**Data scientists are the new rock stars of IT**



*“Of course, I don't even get out of bed for less than a petabyte”*

# Compressed sensing



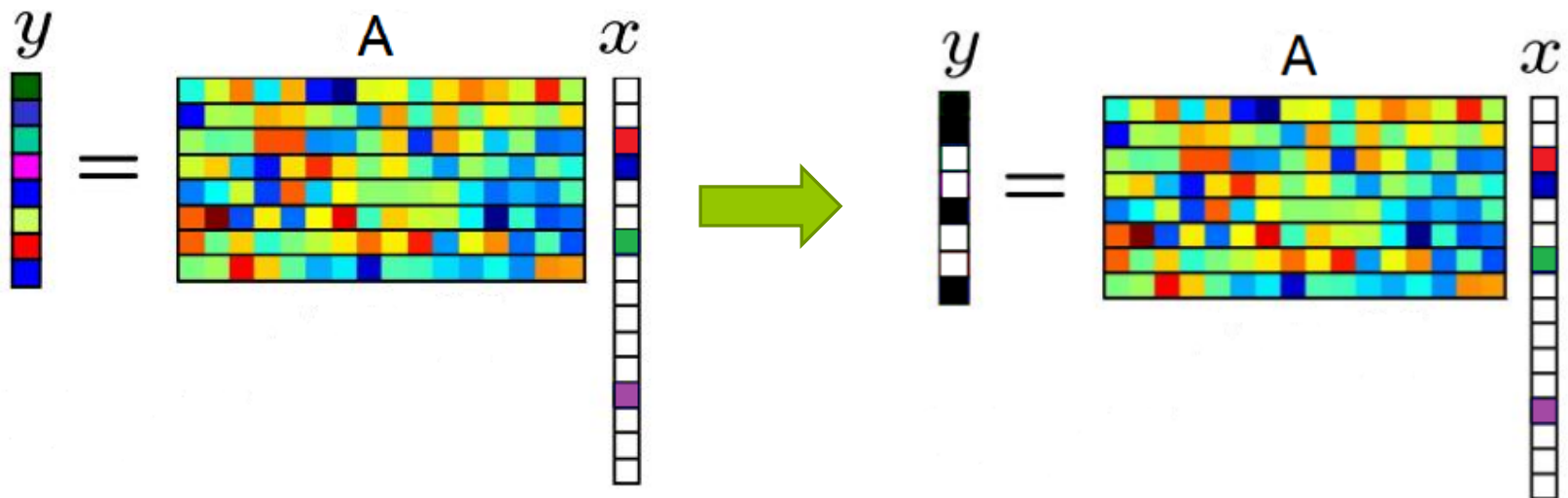
- ❖ Need to solve highly underdetermined linear system

$$y = Ax$$

- ❖ A has a null-space!

# 1-bit compressed sensing

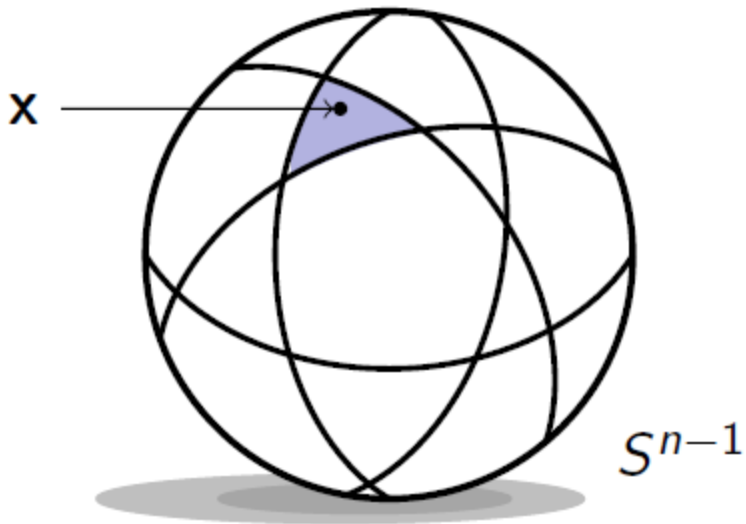
- ❖ Store only the first bit – the SIGN of each measurement



# 1-bit compressed sensing

- ❖ Store only the first bit – the SIGN of each measurement

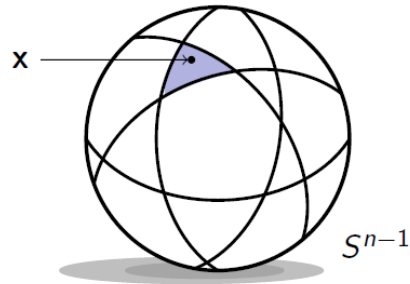
Geometric intuition



# 1-bit compressed sensing

- ❖ Store only the first bit – the SIGN of each measurement

Geometric intuition



- ❖ Remedy: Use “dithers” to estimate the norm of  $x$

$$y_i = \text{sign}(\langle \mathbf{a}_i, \mathbf{x} \rangle - \tau_i)$$

# Moral:

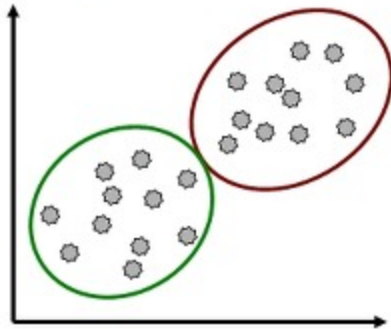
- ❖ One-bit (binary) data is as efficient as it gets
- ❖ It may still contain enough “information” about the signal to perform inference tasks



# Problem: classification

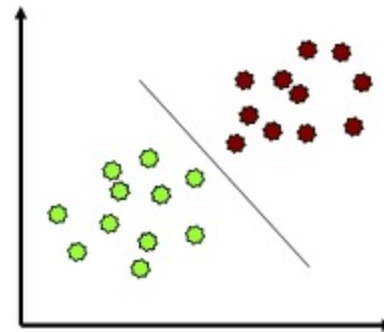
## CLUSTERING

- Data is not labeled
- Group points that are "close" to each other
- Identify structure or patterns in data
- Unsupervised learning

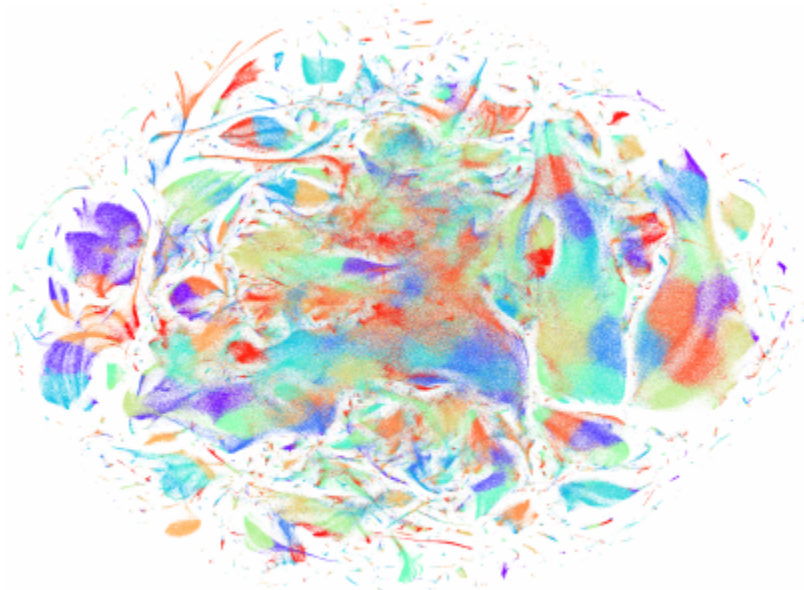


## CLASSIFICATION

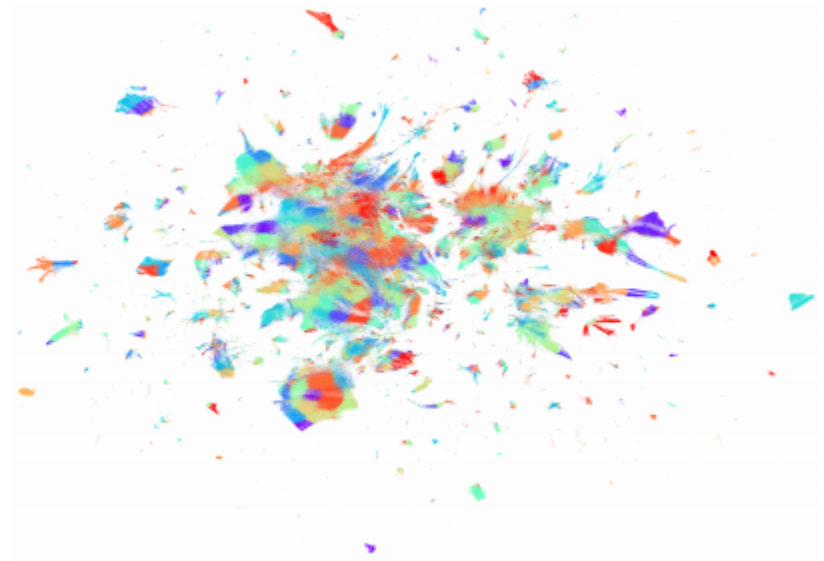
- Labeled data points
- Want a "rule" that assigns labels to new points
- Supervised learning



# Problem: reality



(c) WikiDoc (t-SNE)

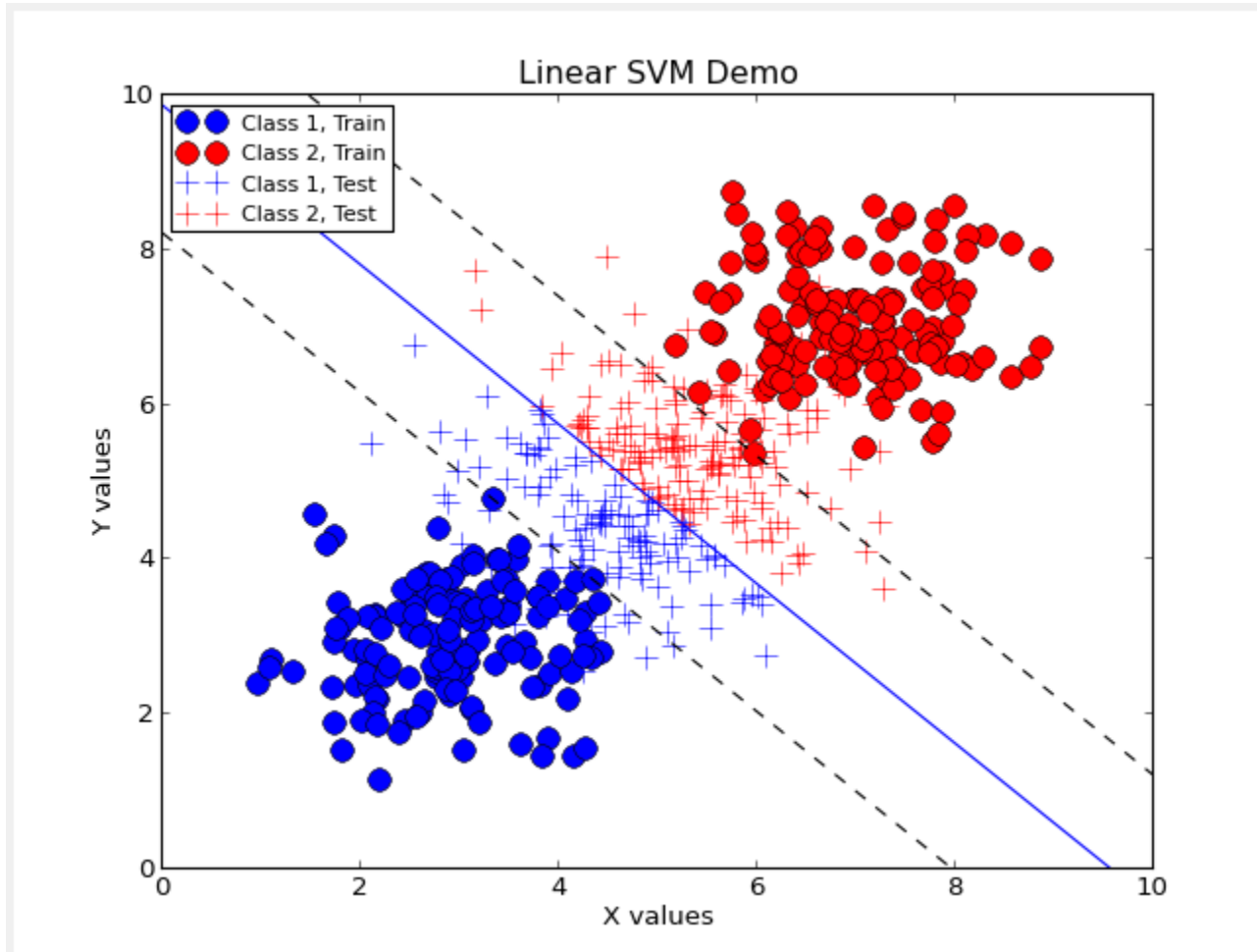


(d) WikiDoc (LargeVis)

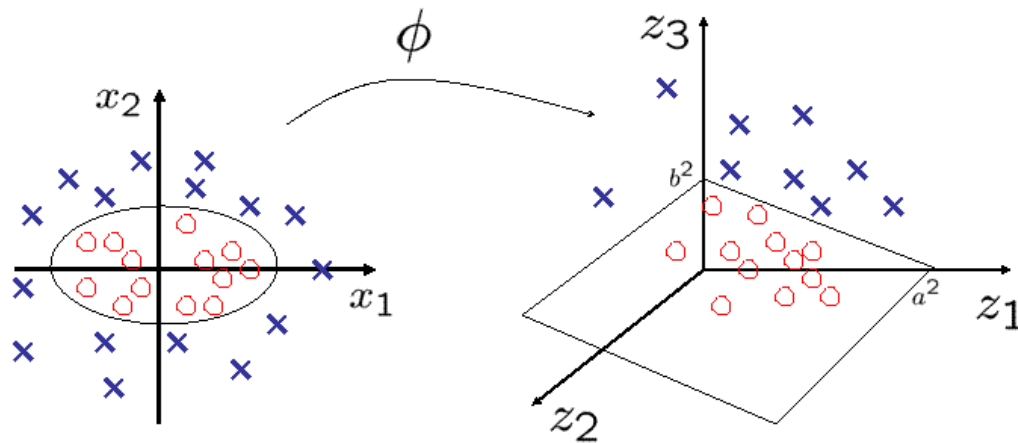
(Tang et al. 2016)



# Background: SVM



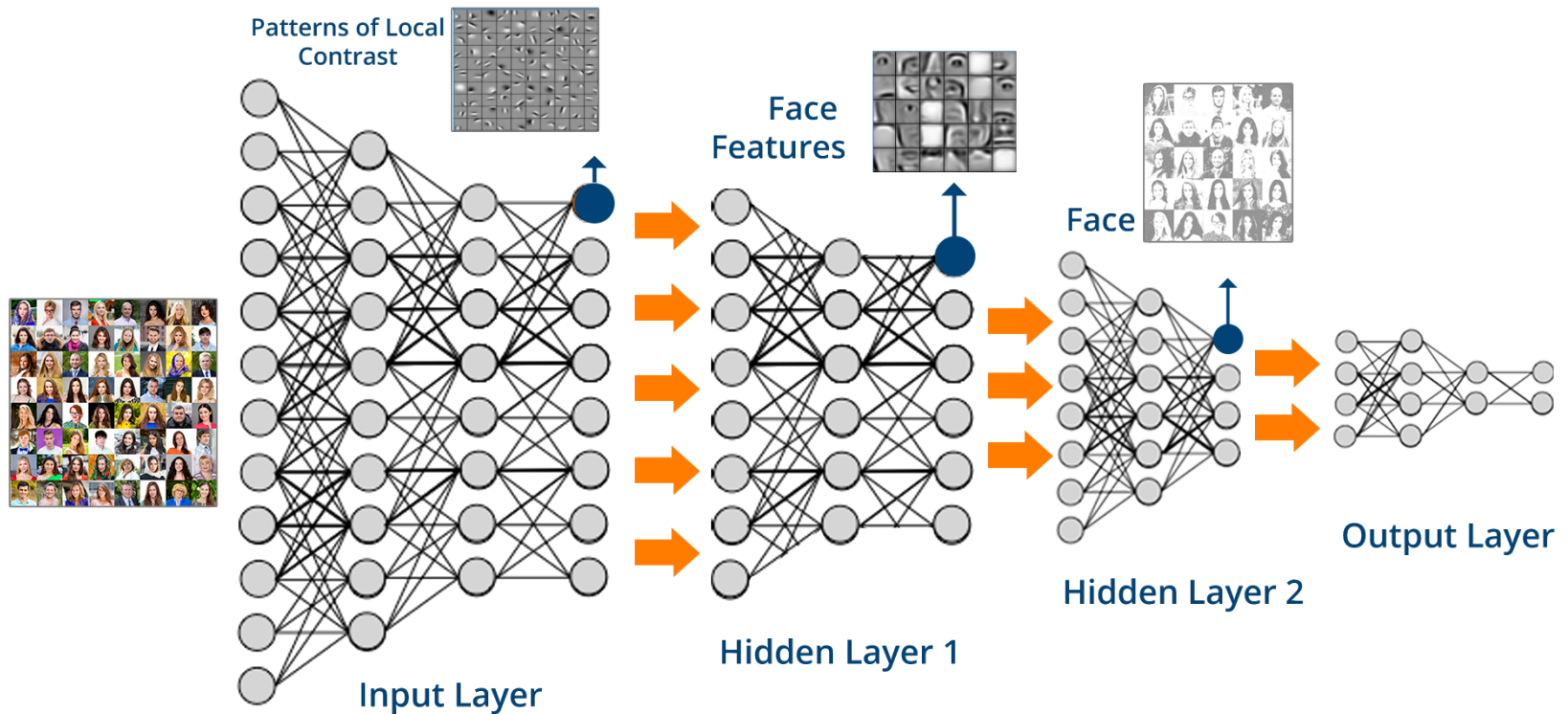
# Background: SVM



$$\phi : (x_1, x_2) \longrightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\left(\frac{x_1}{a}\right)^2 + \left(\frac{x_2}{b}\right)^2 = 1 \longrightarrow \frac{z_1}{a^2} + \frac{z_3}{b^2} = 1$$

# Background: Deep Learning



# Our goals

Design classification scheme that:

- ❖ Uses binary data
- ❖ Is simple and efficient
- ❖ Uses layers in an interpretable way



# Notation

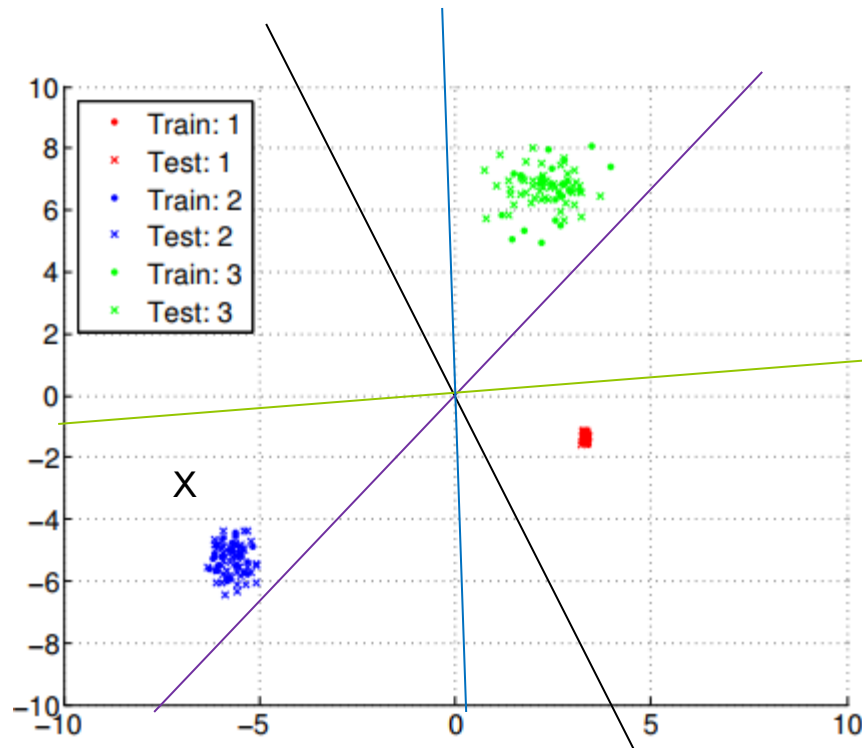
- ❖  $X$  :  $n \times p$  training data matrix (data in columns)
- ❖  $A$  :  $m \times n$  (random) matrix
- ❖  $Q = \text{sign}(AX)$  :  $m \times p$  binary training data
- ❖  $G$  : # of classes
- ❖  $b$  :  $p$  training labels ( $1-G$ )
- ❖  $L$  : # of layers in our design

# Main idea

- ❖ Each row of  $A$  corresponds to a hyperplane
- ❖ Each binary measurement  $Q_{ij}$  indicates on which side of the  $i^{\text{th}}$  hyperplane data point  $X_j$  lies
- ❖ If we gather enough of this info for all the training data, we can use it to predict the class label for a new test point  $x$

# Single layer

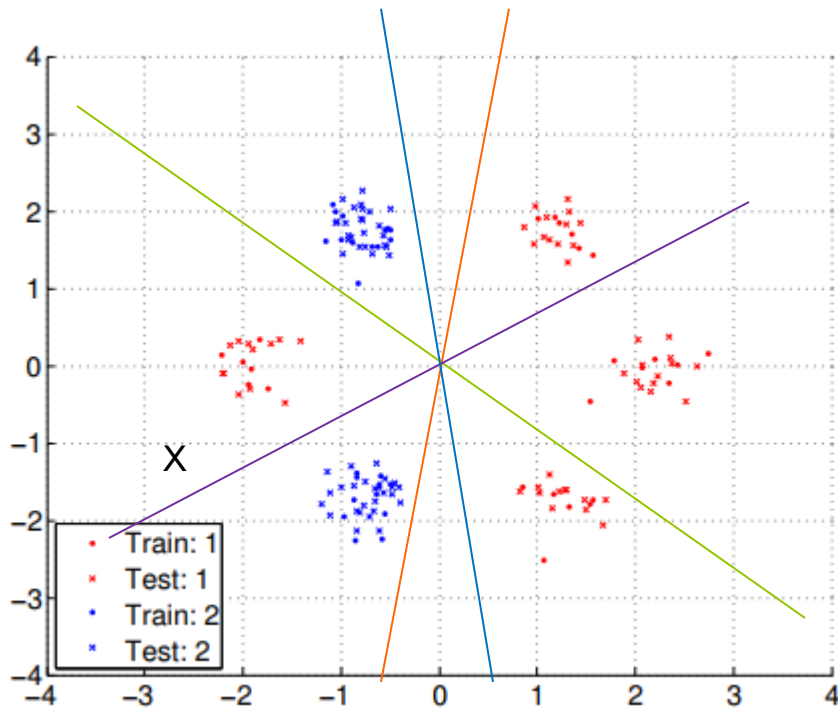
- ❖ All the hyperplane information in  $Q$  is enough



For a new point  $x$ , simply compare its sign pattern with those of the training points and choose the label it matches the most often

# Multiple layers

❖ What about?



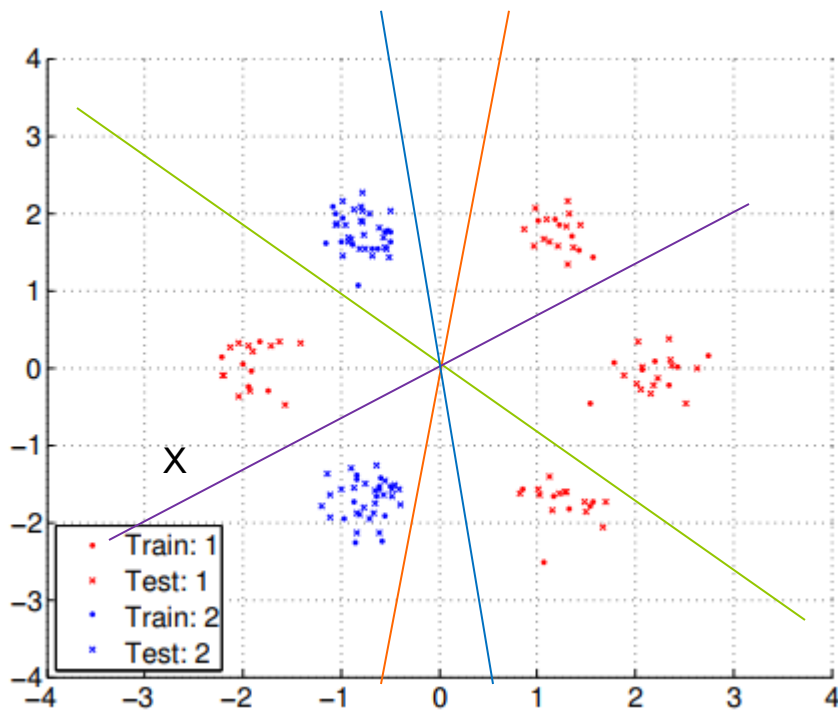
For ANY hyperplane, there are both red and blue on at least one side of it





# Multiple layers

❖ Now PAIRS of hyperplanes do the trick



For a new point  $x$ , simply compare its sign pattern for hyperplane PAIRS with those of the training points and choose the label it matches the most often

# Multiple layers

What about higher dimensions?

- ❖ We continue this strategy to build layers, the 1<sup>th</sup> layer corresponding to 1-tuples of hyperplanes
- ❖ For simplicity (and computation), we consider  $m$  1-tuples at each layer, selected randomly from all  $\binom{m}{\ell}$  possible
- ❖ For a new test point  $x$ , we use the sign patterns across all layers for classification

# Using all layers

How to integrate the info from all layers?

- ❖  $l$  : layer
- ❖  $i$  : index from 1 to  $m$
- ❖  $\Lambda_{\ell,i}$  : the set of 1 indices indicating which hyperplanes are selected in the  $i^{\text{th}}$  1-tuple
- ❖  $t$  : a possible sign pattern of 1 +/- 1s
- ❖  $g$  : a class index (from 1 to  $G$ )
- ❖  $P_{g|t}$  : the # of training points from the  $g^{\text{th}}$  class having sign pattern  $t$  from the hyperplanes in  $\Lambda_{\ell,i}$

# Using all layers

- ❖  $P_{g|t}$ : the # of training points from the  $g^{\text{th}}$  class having sign pattern  $t$  from the hyperplanes in  $\Lambda_{\ell,i}$

$$r(\ell, i, t, g) = \frac{P_{g|t}}{\sum_{j=1}^G P_{j|t}} \frac{\sum_{j=1}^G |P_{g|t} - P_{j|t}|}{\sum_{j=1}^G P_{j|t}}$$

# Using all layers

- ❖  $P_{g|t}$ : the # of training points from the  $g^{\text{th}}$  class having sign pattern  $t$  from the hyperplanes in  $\Lambda_{\ell,i}$

$$r(\ell, i, t, g) = \frac{P_{g|t}}{\sum_{j=1}^G P_{j|t}} \frac{\sum_{j=1}^G |P_{g|t} - P_{j|t}|}{\sum_{j=1}^G P_{j|t}}$$



fraction of training points in class  $g$  out of all points with pattern  $t$

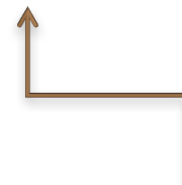
# Using all layers

- ❖  $P_{g|t}$ : the # of training points from the  $g^{\text{th}}$  class having sign pattern  $t$  from the hyperplanes in  $\Lambda_{\ell,i}$

$$r(\ell, i, t, g) = \frac{P_{g|t}}{\sum_{j=1}^G P_{j|t}} \frac{\sum_{j=1}^G |P_{g|t} - P_{j|t}|}{\sum_{j=1}^G P_{j|t}}$$



fraction of training points in class  $g$  out of all points with pattern  $t$



gives more weight to group  $g$  when its size is much different than others with same sign pattern

# Our method: training

---

## Algorithm 1 Training

---

**input:** binary training data  $Q$ , training labels  $b$ , number of classes  $G$ , number of layers  $L$

**for**  $\ell$  from 1 to  $L$ ,  $i$  from 1 to  $m$  **do**

**select:** Randomly select  $\Lambda_{\ell,i} \subset [m]$ ,  $|\Lambda_{\ell,i}| = \ell$

**determine:** Determine the  $T_{\ell,i} \in \mathbb{N}$  unique column patterns in  $Q^{\Lambda_{\ell,i}}$

**for**  $t$  from 1 to  $T_{\ell,i}$ ,  $g$  from 1 to  $G$  **do**

**compute:** Compute  $r(\ell, i, t, g)$  by (1)

**end for**

**end for**

---

- ❖ “For each layer, pick the 1-tuples and then compute all values of  $r(1,i,t,g)$ ”



# Our method: testing

---

## Algorithm 2 Classification

---

**input:** binary data  $q$ , number of classes  $G$ , number of layers  $L$ , learned parameters  $r(\ell, i, t, g)$ ,  $T_{\ell, i}$ , and  $\Lambda_{\ell, i}$  from Algorithm 1

**initialize:**  $\tilde{r}(g) = 0$  for  $g = 1, \dots, G$ .

**for**  $\ell$  from 1 to  $L$ ,  $i$  from 1 to  $m$  **do**

**identify:** Identify the pattern  $t^* \in [T_{\ell, i}]$  to which  $q^{\Lambda_{\ell, i}}$  corresponds

**for**  $g$  from 1 to  $G$  **do**

**update:**  $\tilde{r}(g) = \tilde{r}(g) + r(\ell, i, t^*, g)$

**end for**

**end for**

**scale:** Set  $\tilde{r}(g) = \frac{\tilde{r}(g)}{Lm}$  for  $g = 1, \dots, G$

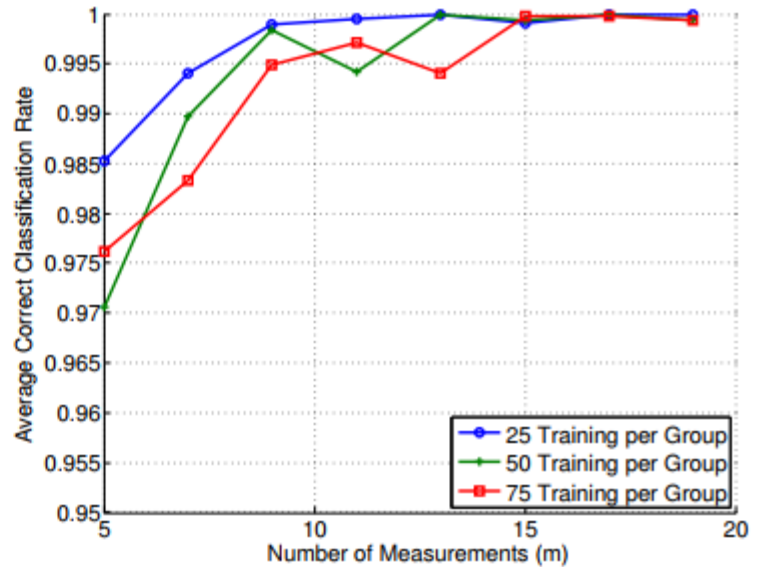
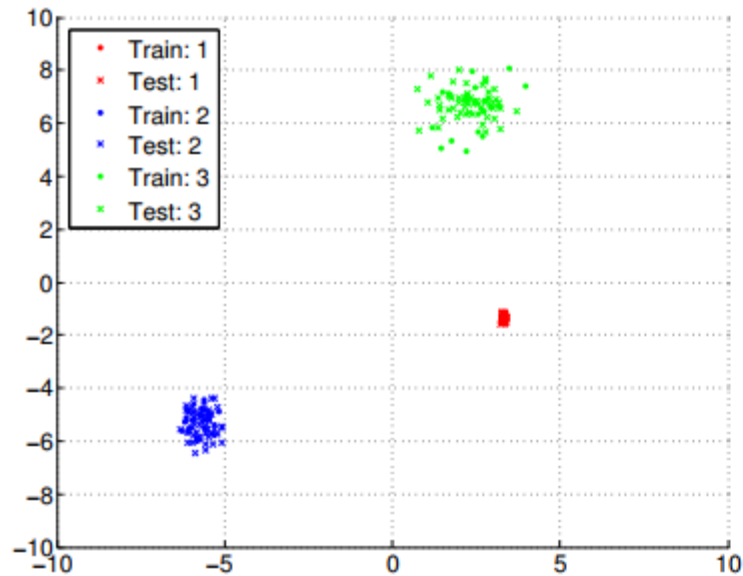
**classify:**  $\hat{b}_x = \operatorname{argmax}_{g \in \{1, \dots, G\}} \tilde{r}(g)$

---

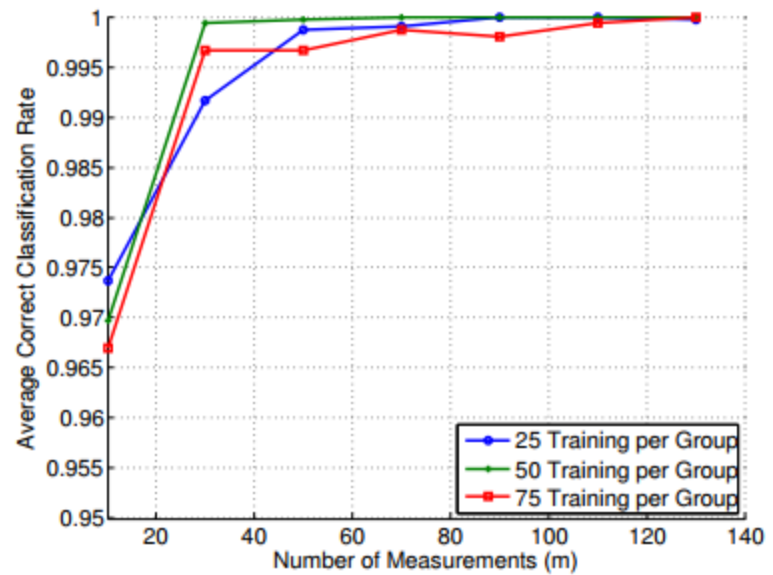
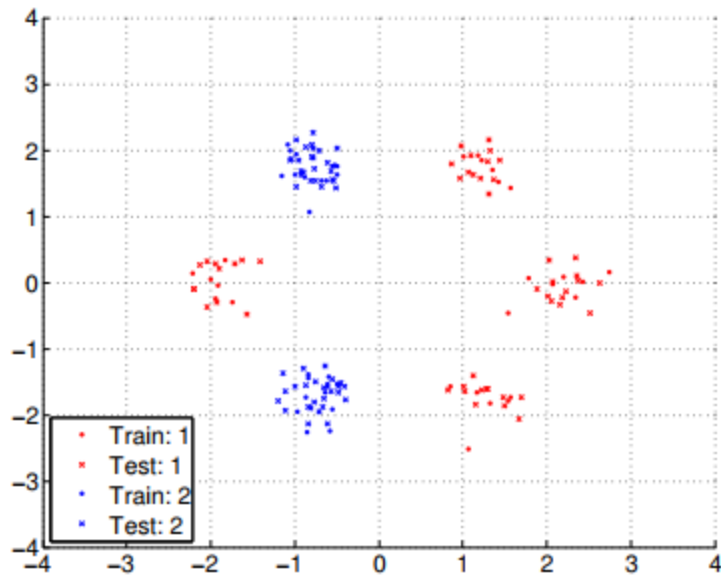
- ❖ “For a new point  $x$  with sign pattern  $t^*$ , compute the sum of all  $r(1, i, t^*, g)$  for each class  $g$  and then assign the label  $g$  which has the largest sum.”



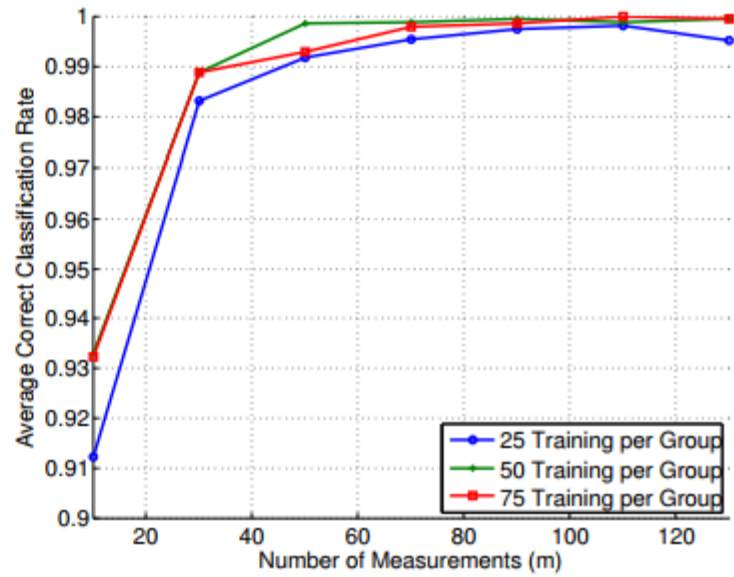
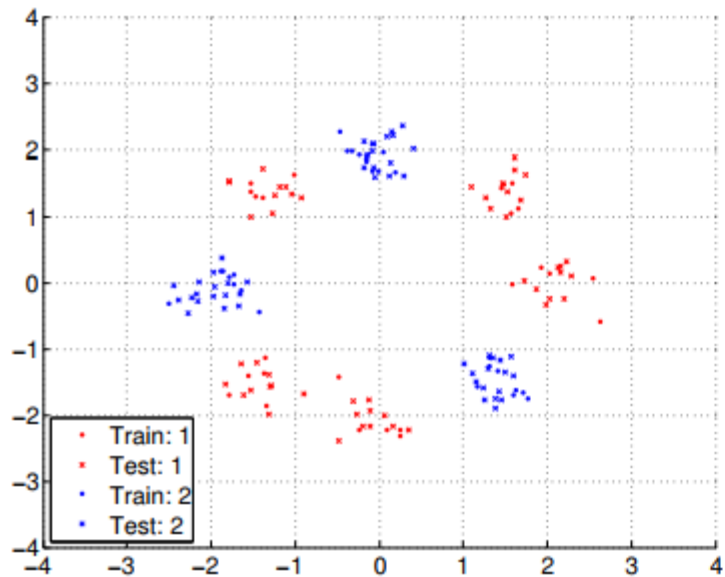
# Results (L=1)



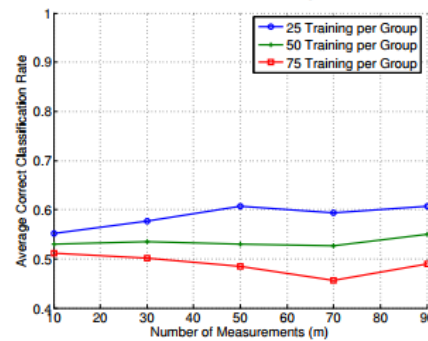
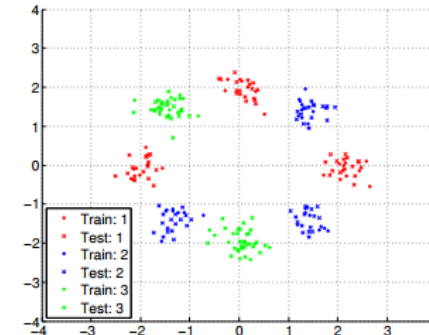
# Results (L=4)



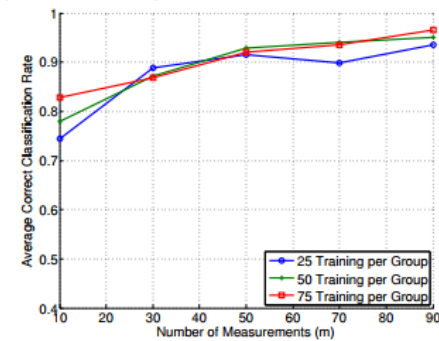
# Results (L=5)



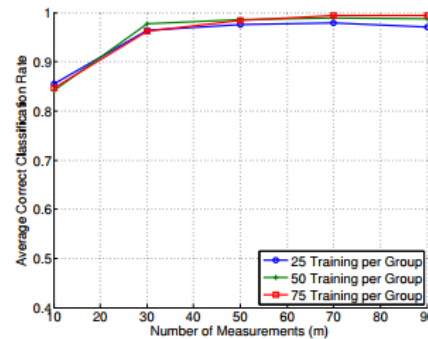
# Results



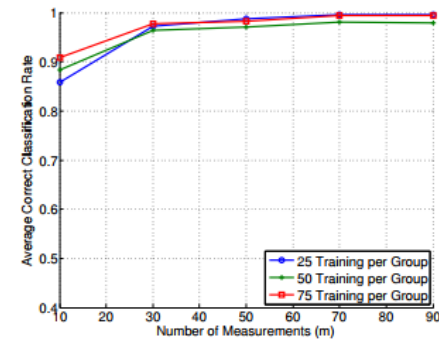
(a)  $L = 1$



(b)  $L = 2$

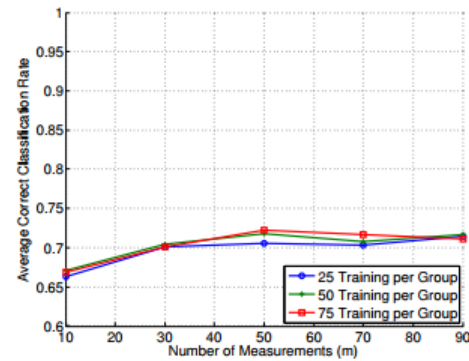
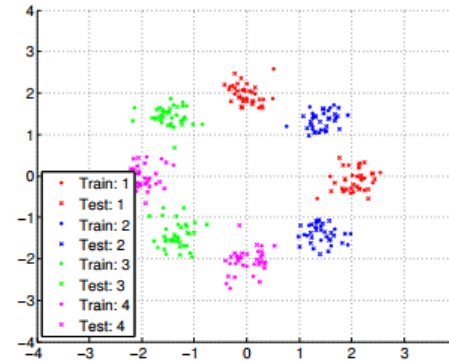


(c)  $L = 3$

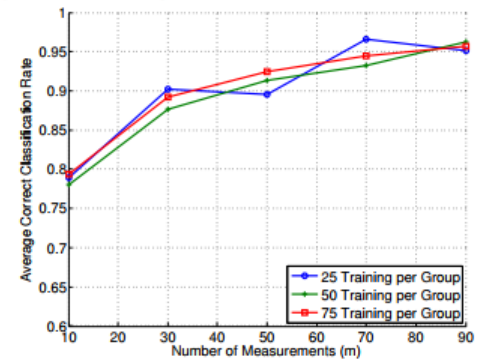


(d)  $L = 4$

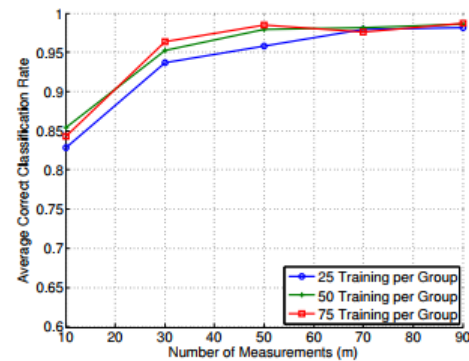
# Results



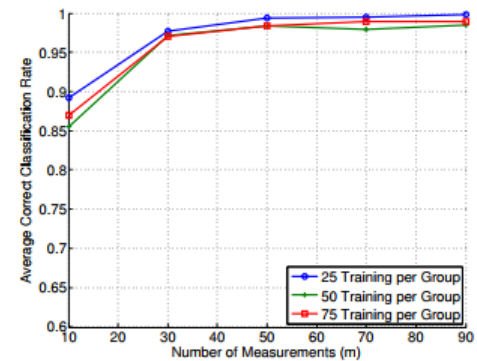
(a)  $L = 1$



(b)  $L = 2$

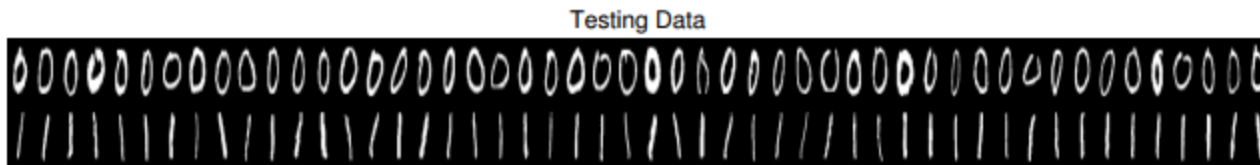
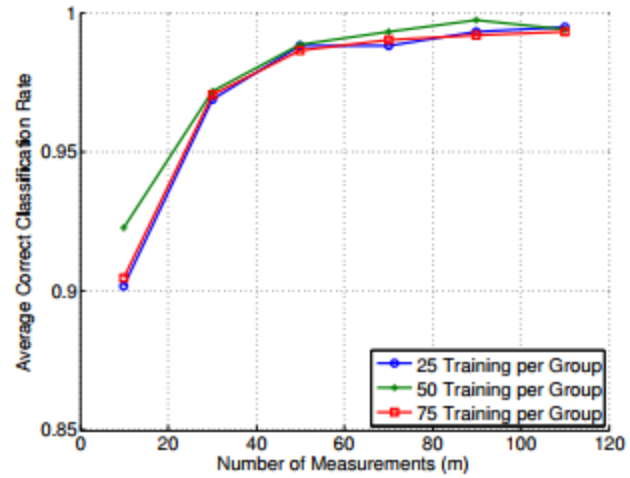
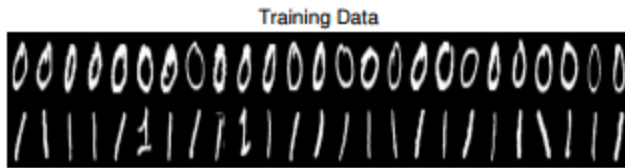


(c)  $L = 3$

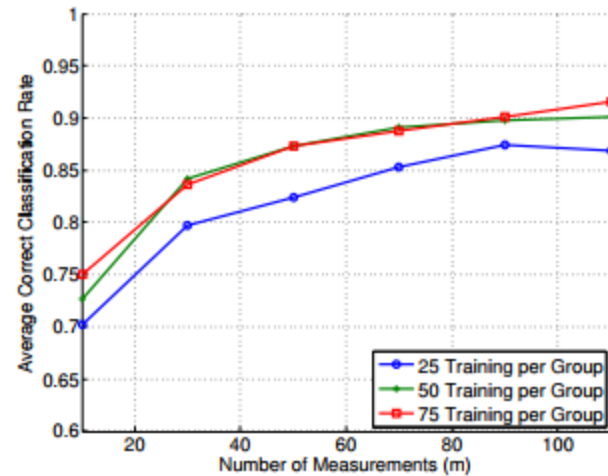


(d)  $L = 4$

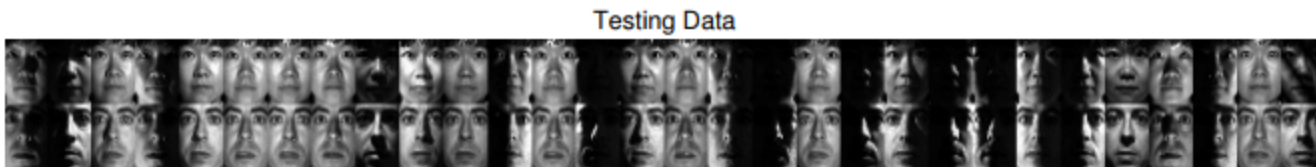
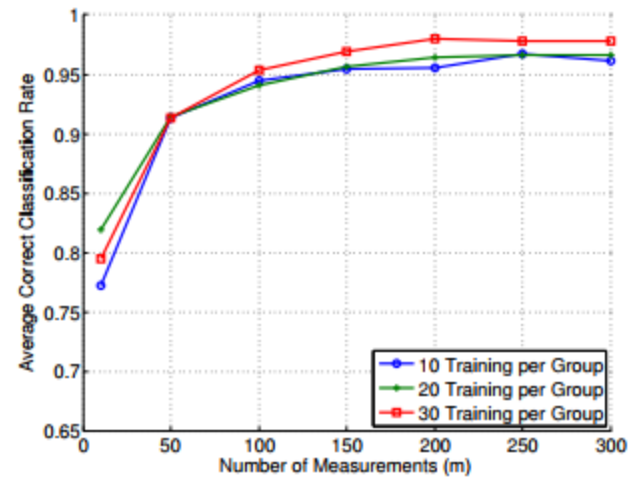
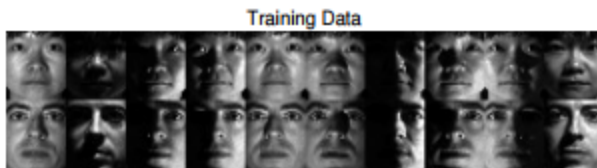
# Results (L=1)



# Results (L=4)

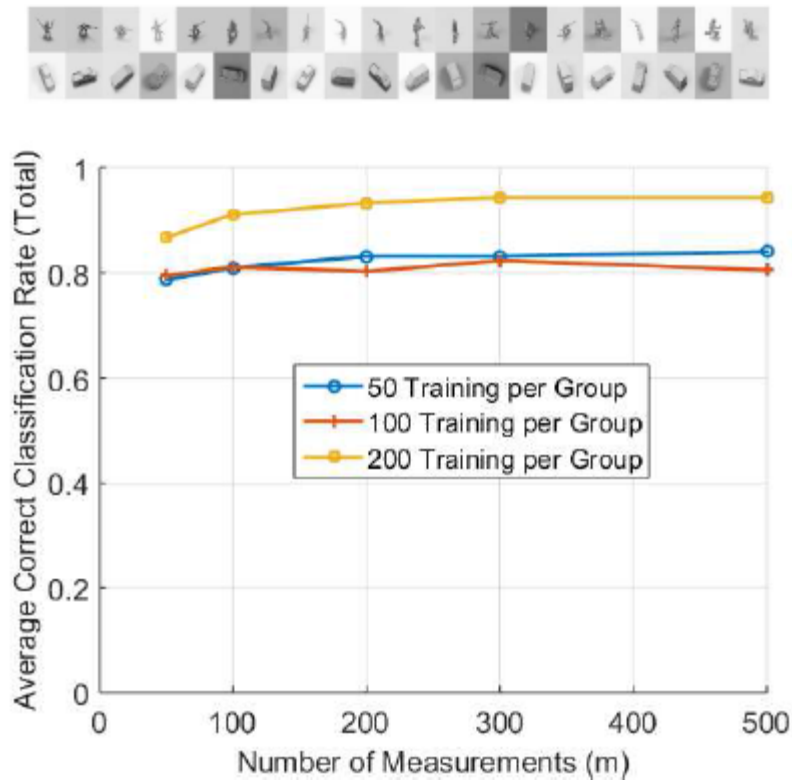


# Results (L=5)



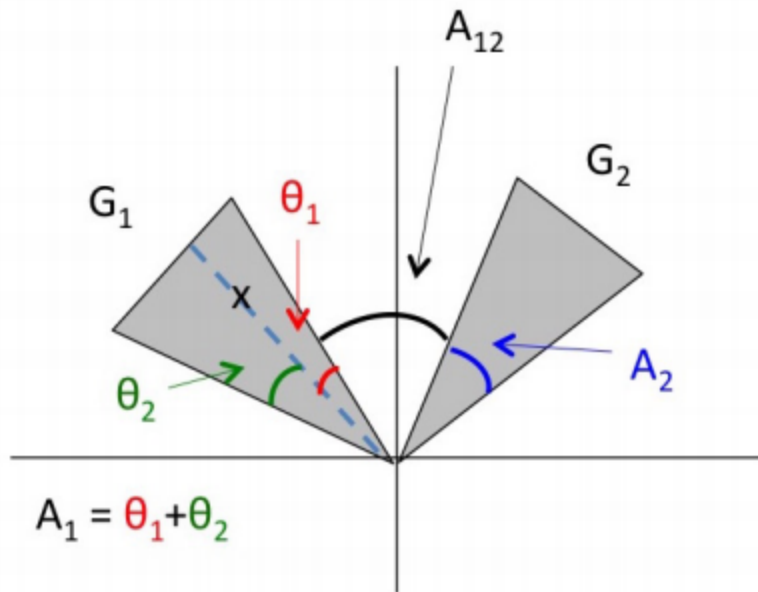


# Results (L=5)

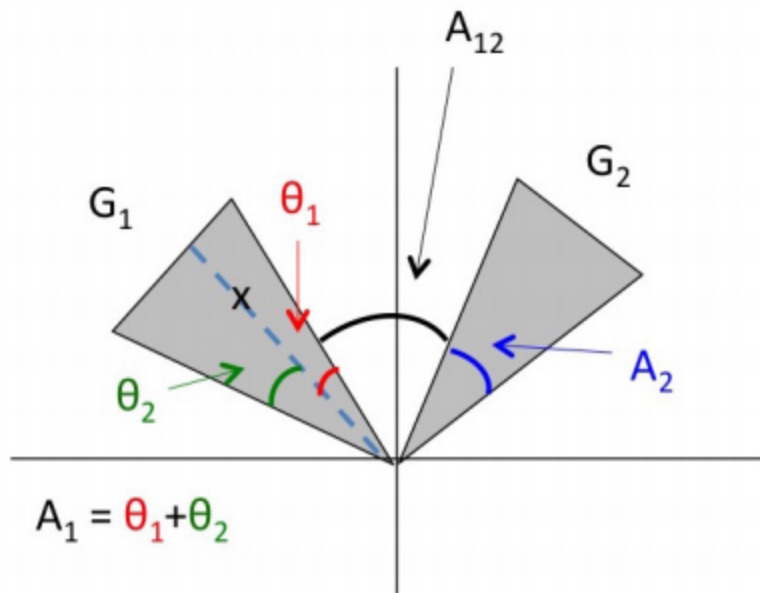


(b)  $L = 5$

# Theory



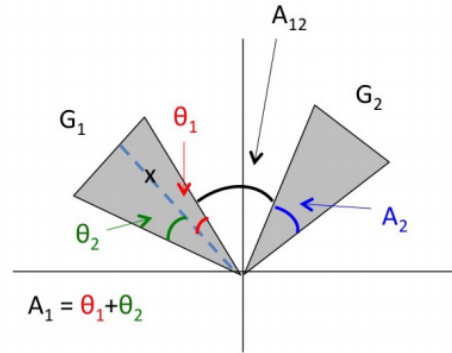
# Theory



$A_{g|t}$  : the angle of class  $g$  with sign pattern  $t$  for the  $i^{\text{th}}$  1-tuple in layer 1

$$r(\ell, i, t, g) = \frac{A_{g|t}}{\sum_{j=1}^G A_{j|t}} \frac{\sum_{j=1}^G |A_{g|t} - A_{j|t}|}{\sum_{j=1}^G A_{j|t}}$$

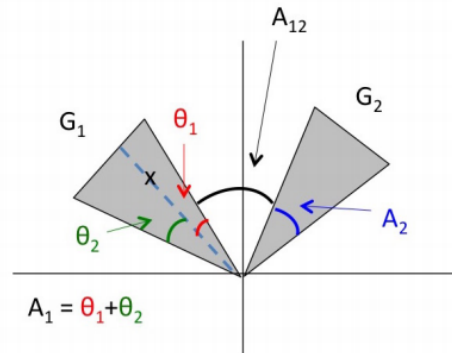
# Theory



**Theorem 1.** *Let the classes  $G_1$  and  $G_2$  be two cones in  $\mathbb{R}^2$  defined by angular measures  $A_1$  and  $A_2$ , respectively, and suppose regions of the same angular measure have the same density of training points. Suppose  $A_1 = A_2$ ,  $\theta_1 = \theta_2$ , and  $A_{12} + A_1 + A_2 \leq \pi$ . Then, the probability that a data point  $x \in G_1$  gets classified in class  $G_1$  by Algorithms 1 and 2 using a single layer and a measurement matrix  $A \in \mathbb{R}^{m \times 2}$  with independent standard Gaussian entries is bounded as follows,*

$$\begin{aligned} \mathbb{P}[\hat{b}_x = 1] \geq 1 - \sum_{j=0}^m \sum_{k_1, \theta_1=0}^m \sum_{k_1, \theta_2=0}^m \sum_{k_2=0}^m \sum_{k=0}^m & \binom{m}{j, k_1, \theta_1, k_1, \theta_2, k_2, k} \left(\frac{A_{12}}{\pi}\right)^j \left(\frac{A_1}{2\pi}\right)^{k_1, \theta_1 + k_1, \theta_2} \\ & j + k_1, \theta_1 + k_1, \theta_2 + k_2 + k = m, k_1, \theta_2 \geq 9(j + k_1, \theta_1) \\ & \times \left(\frac{A_1}{\pi}\right)^{k_2} \left(\frac{\pi - 2A_1 - A_{12}}{\pi}\right)^k. \end{aligned} \quad (3)$$

# Theory

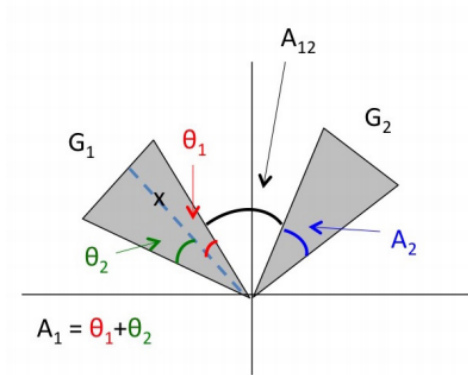
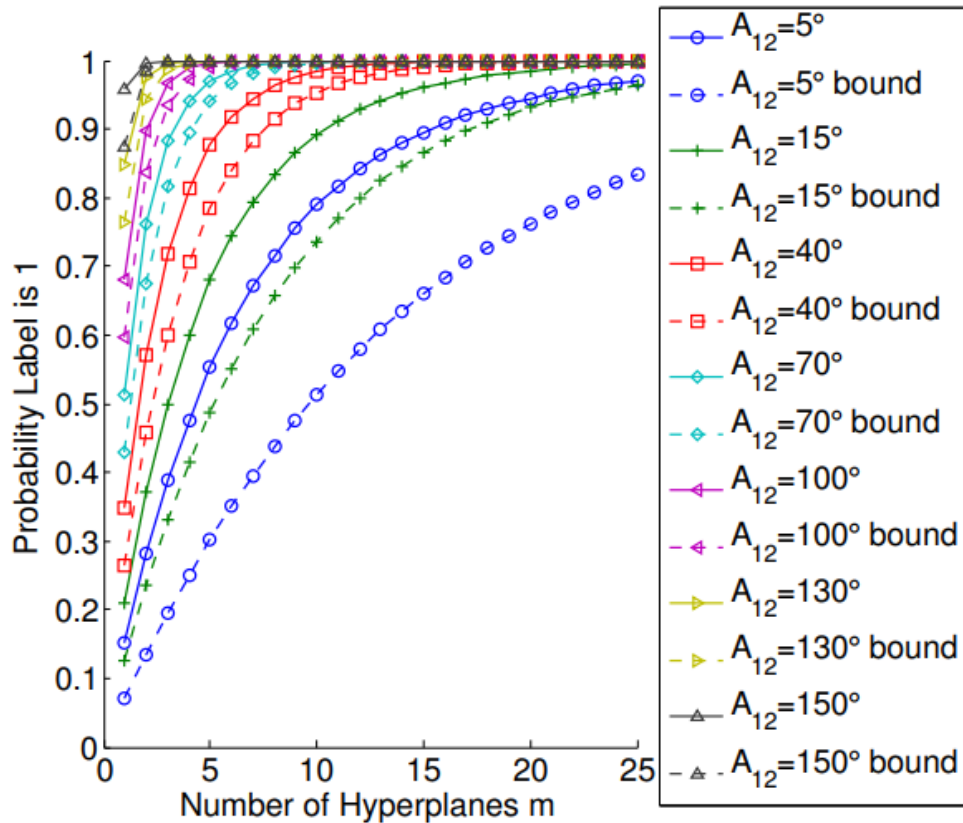


**Theorem 1.** Let the classes  $G_1$  and  $G_2$  be two cones in  $\mathbb{R}^2$  defined by angular measures  $A_1$  and  $A_2$ , respectively, and suppose regions of the same angular measure have the same density of training points. Suppose  $A_1 = A_2$ ,  $\theta_1 = \theta_2$ , and  $A_{12} + A_1 + A_2 \leq \pi$ . Then, the probability that a data point  $x \in G_1$  gets classified in class  $G_1$  by Algorithms 1 and 2 using a single layer and a measurement matrix  $A \in \mathbb{R}^{m \times 2}$  with independent standard Gaussian entries is bounded as follows,

$$\mathbb{P}[\hat{b}_x = 1] \geq 1 - \sum_{j=0}^m \sum_{k_1, \theta_1=0}^m \sum_{k_1, \theta_2=0}^m \sum_{k_2=0}^m \sum_{k=0}^m \binom{m}{j, k_1, \theta_1, k_1, \theta_2, k_2, k} \left(\frac{A_{12}}{\pi}\right)^j \left(\frac{A_1}{2\pi}\right)^{k_1, \theta_1 + k_1, \theta_2} \\ j + k_1, \theta_1 + k_1, \theta_2 + k_2 + k = m, k_1, \theta_2 \geq 9(j + k_1, \theta_1) \\ \times \left(\frac{A_1}{\pi}\right)^{k_2} \left(\frac{\pi - 2A_1 - A_{12}}{\pi}\right)^k. \quad (3)$$

❖ As  $m \rightarrow \infty$ ,  $P(\text{correct label}) \rightarrow 1$

# Theory



# Take-away



- ❖ Simple classification from binary data
  - ❖ Efficient storage of the data
  - ❖ Efficient and simple algorithm
  - ❖ Theoretical analysis possible
  - ❖ Already competes with state of the art
- ❖ Future work
  - ❖ Dithers to allow for more complicated geometries?
  - ❖ Theoretical analysis of the discrete case?

Thanks for  
**listening**  
to me...



- [deanna@math.ucla.edu](mailto:deanna@math.ucla.edu)
- [math.ucla.edu/~deanna](http://math.ucla.edu/~deanna)
- “Simple Classification using Binary Data”
  - Needell, Saab, Woolf