# USC Viterbi
School of Engineering

# *Active Methods:*

## *Learning as you go and as fast as you can*

Urbashi Mitra

Ming Hsieh Department of Electrical Engineering Department of Computer Science

*University of Southern California*

# Special thanks

Dhruva Kartik PhD'21
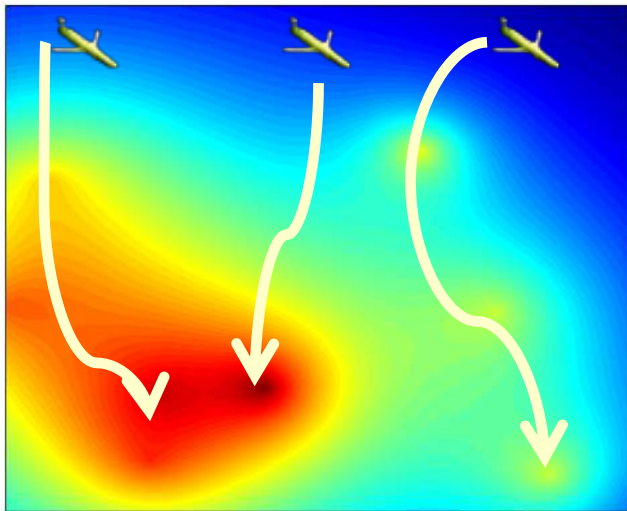
# WHO ARE YOU?

# Background

- Probability/Random Processes

- Detection & Estimation

- Communications

- Information Theory

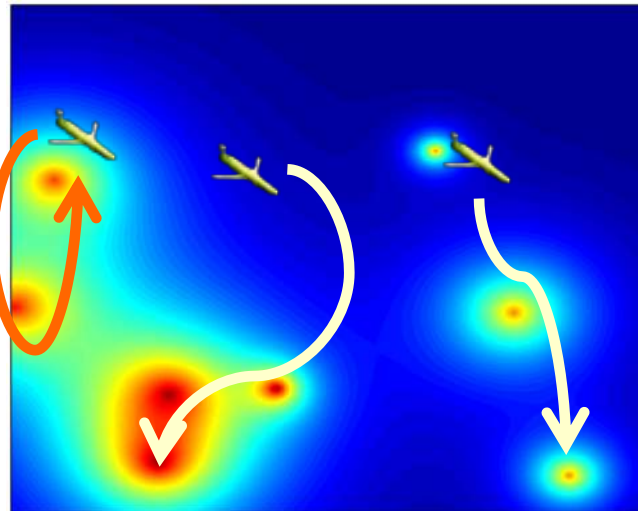- Advanced Information theory

- Machine Learning

# BIG PICTURE

❑ Active hypothesis testing

▪ So many applications!

▪ Information theory in the wild

❑ Important questions

▪ How do you build your tree of actions/observations?

▪ What is the right measure of informativeness that allows you to prune the tree?

❑ Martingales, concentration inequalities

▪ Very useful tools for a wide-range of applications (need more than the CLT)

❑ The classics still matter

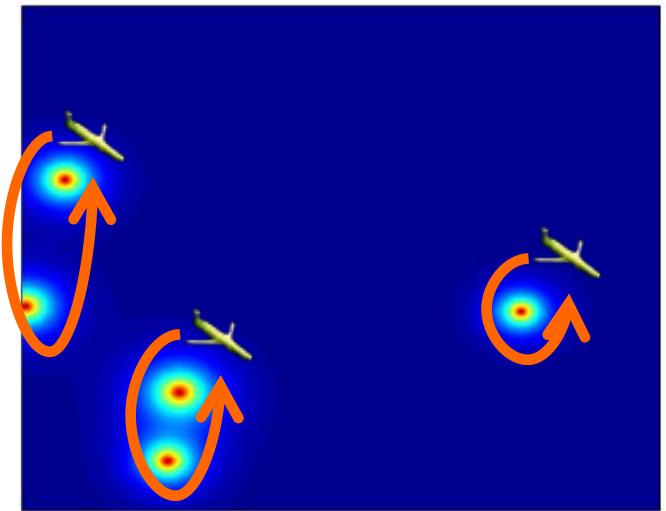▪ Chernoff, Stein, Wald, Blackwell, Fisher, Bayes, Neyman, Pearson

# Exploration-Exploitation
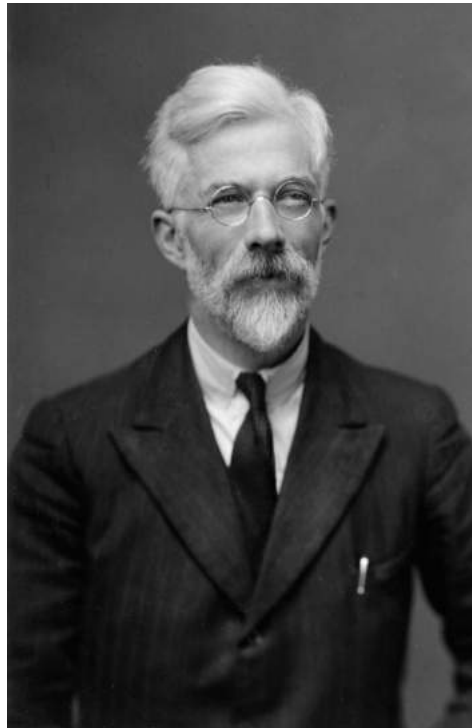


*exploration*
environment unknown

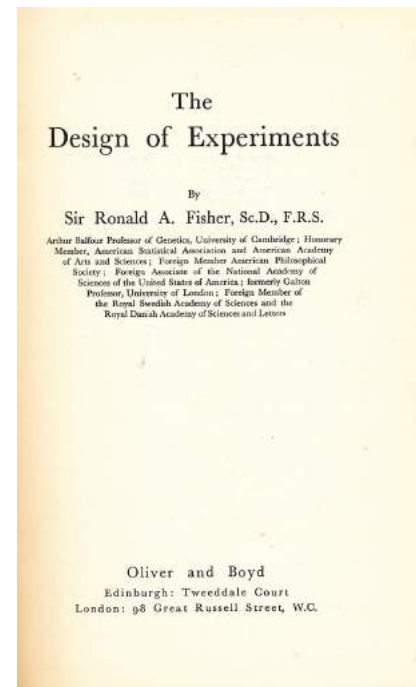*collect observations*
learn

*exploitation*
focus on areas of interest

# Design of Experiments

**Sir Ronald Fisher**
1890-1962



1935

# More broadly

David Blackwell
1919-2010

Herman Chernoff
1923-

Abraham Wald
1902-1950

# MOTIVATING EXAMPLE

# Boundary Detection

- **SENSOR NETWORKS**:

  Actively build boundary

  Data aggregation at

  each layer

- Intrinsic complexity of

  boundary is

$$O(\sqrt{n})$$

Nowak, **M** & Willett, JSAC 2004, IPSN 2003

# Recursive Dyadic Partitions

$\sqrt{n}$ nodes

$\sqrt{n}$ nodes

complete representation

transmit all measurements

# Complete Representation

❏ This is the full tree

# Recursive Dyadic Partitions

$\sqrt{n}$ nodes

$\sqrt{n}$ nodes

pruned representation

transmit averages/some measurements

# Recursive Dyadic Partition

❑ The pruned tree

averaged measurements

from the four

grand- children

averaged measurements

from the two children

# The question

- ❑ What is the optimal grouping?
  - ▪ The cost of keeping fine-grained measurements/size of the tree

$$P = \text{partition}$$
$$|\theta(P)| = \text{size of partition/complexity}$$

  - ▪ The cost of reducing fidelity – squared error

$$R(\theta, x) = \sum_{i,j=1}^{\sqrt{n}} \left(\theta(i,j) - x_{i,j}\right)^2$$

# Connections to Group Testing

❑ Used in WW2 to test soldiers for syphilis

  ▪ R. Dorfman, "The Detection of Defective Members of Large Populations," The Annals of Mathematical Statistics, 1943

  ▪ Binary search

❑ Complexity reduction

$$N \text{ tests} \rightarrow \log(N) \text{ tests}$$

# Estimation Criterion

❑ Penalized empirical risk

▪ Squared error

$$R(\theta, x) = \sum_{i,j=1}^{\sqrt{n}} \left(\theta(i,j) - x_{i,j}\right)^2$$

▪ Complexity of RDP

$$\widehat{\theta}_n = \arg \min_{\theta(P):P \in \mathcal{P}_n} \left\{ R(\theta, x) + 2\sigma^2 f(n)|\theta(P)| \right\}$$

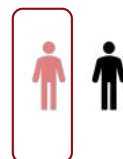$$|\theta(P)| \sim 64 \text{ versus } |\theta(P')| \sim 28$$

# Metric for Pruning

- Over a dyadic partition compare penalized cost of average measurement versus measurements from a finer scale

- Can show

$$\frac{1}{n} \sum_{i,j=1}^{\sqrt{n}} E\left[\left(\hat{\theta}_n(i,j) - \theta^*(i,j)\right)^2\right] \leq O\left(\sqrt{\frac{\log n}{n}}\right)$$

  - Versus minimax lower bound $\mathrm{MSE} \geq O\left(\frac{1}{\sqrt{n}}\right)$

- Optimally pruned partition of order $O\left(\sqrt{n}\right)$

# Numerical Results



65536 Observations — estimates — Partition, $|\theta| = 1111$

1024 Observations — estimates — Partition, $|\theta| = 172$

256 Observations — estimates — Partition, $|\theta| = 70$

# Adaptive Boundary Estimation

- ❑ Actively building up representation, BUT
- ❑ All measurements taken once
  - ▪ Reverse engineering representation
  - ▪ Notion of higher utility/reward
- ❑ Notion of one representation being better than another
- ❑ Not active in measurement collection

# BASICS OF HYPOTHESIS TESTING

# Hypotheses and Likelihoods

❑ Binary Hypotheses:

$$H_0 \quad : \quad \text{null hypothesis}$$
$$H_1 \quad : \quad \text{alternate hypothesis}$$

$X = 0$ : If $H_0$ is true

$X = 1$ : If $H_1$ is true

❑ Model:

$$\mathbb{P}[Y = y \mid X = 0] = p_0(y)$$
$$\mathbb{P}[Y = y \mid X = 1] = p_1(y)$$

$Y \in \mathcal{Y}$
finite alphabet

observation

likelihood functions

# Binary Hypothesis Testing

$H_1$ : alternative hypothesis

$H_0$ : null hypothesis

$f(y) =$ decision rule

$\quad = \{0, 1\}$

$y =$ observation

$p_i(y) =$ pdf of $y$ given $H_i$

partition observation space



$Y_k \sim p_X$

i.i.d. observations

$\hat{X} = f(Y^n, \text{random})$

inference        decision rule

# Good Decision Rules

❏ Log-likelihood Ratio (LLR):

$$L_n = \log \frac{p_0(Y^n)}{p_1(Y^n)} = \sum_{k=1}^{n} \log \frac{p_0(Y_k)}{p_1(Y_k)}$$

❏ Good decision rules

change the metric
change the threshold

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n \geq \tau \\ H_1 & \text{if } L_n < \tau \end{cases}$$

likelihood ratio test

$\tau_{\text{NP}}$

$\tau_{\text{B}}$

# Kullback-Leibler Divergence

❑ DEFINITION:

$$D(p\|q) = \sum_y p(y) \log \frac{p(y)}{q(y)}$$

Expectation of LLR is related to KL-Divergence

❑ Like a ``distance'' between two distributions

❑ BUT, **not** symmetric: $D(p\|q) \neq D(q\|p)$

# Likelihood Ratio Tests

❑ Equivalent representation with respect to the KL divergence

$$L(y^n) > \tau$$

$$D\left(p(y^n)\|p_0(y^n)\right) - D\left(p(y^n)\|p_1(y^n)\right) > \frac{1}{n}\log\tau$$

❑ The empirical distribution is closest to which hypothesis?

❑ NOTE:
$$\mathbb{E}_0[L_n] = nD(p_0\|p_1)$$
$$\mathbb{E}_1[L_n] = -nD(p_1\|p_0)$$

❑ Bayes optimal rule versus Neyman-Pearson rule

▪ How to select $\tau$ ?

# Bayes Rule

❑ **Bayesian Risk:**

$$C_{ij} = \text{cost of selecting } i \text{ when } j \text{ is true}$$

$$r(f) = \sum_{j} \pi_j \sum_{i} C_{ij} \mathbb{P}[\hat{X} = i \mid X = j]$$

priors        costs        infer $i$ given truth is $j$

❑ **Bayes rule:**

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n \geq \tau \\ H_1 & \text{if } L_n < \tau \end{cases} \qquad \tau = \log \frac{\pi_1 (C_{01} - C_{11})}{\pi_0 (C_{10} - C_{00})}$$

likelihood ratio test

# Special Cases

- Uniform costs

$$C_{ij} = \delta(i - j)$$

$$\rightarrow \quad \tau = \log \frac{\pi_0}{\pi_1}$$

- *Maximum a posteriori rule*

- Uniform costs and equal priors          all likelihood ratio tests

$$\tau = \log(1) = 0$$

- Maximum likelihood rule

# Gaussian Example

$$\mathcal{H}_i \quad : \quad Y \sim \mathcal{N}(\mu_i, \sigma^2)$$

$$L(y) \quad = \quad \left[\left(\frac{\mu_1 - \mu_0}{\sigma^2}\right)\left(y - \frac{\mu_0 + \mu_1}{2}\right)\right]$$

$$L(y) \quad \underset{<}{\overset{\geq}{\gtrless}} \quad \tau$$

$$y \quad \underset{<}{\overset{\geq}{\gtrless}} \quad \tau' = \frac{\sigma^2}{\mu_1 - \mu_2} \ln \tau + \frac{\mu_0 + \mu_1}{2}$$

$$r(f) = \sum_j \pi_j \sum_i C_{ij} \mathbb{P}[\hat{X} = i \mid X = j]$$

$$\mathbb{P}\left[\hat{X} = 1 | X = j\right] \quad = \quad \mathbb{P}\left[\mathbf{Y} \geq \tau' | X = j\right]$$

$$= \quad Q\left(\frac{\tau' - \mu_j}{\sigma}\right)$$

# How to bound Performance?

❑ Y is a random variable

   ▪ Moment generating function

$$\mu(s) = \mathbb{E}[\exp(-sY)]$$

❑ Chernoff Bound

$$\forall \ s \geq 0 \quad \mathbb{P}\left[Y \geq a\right] \leq e^{-sa}\mu(s) \ \forall \ s$$

$$\rightarrow \mathbb{P}\left[Y \geq a\right] \leq \min_{s} e^{-sa}\mu(s)$$

❑ Proof  - via Markov inequality

$$\mathbb{P}\left[X \leq a\right] \ \leq \ \frac{\mathbb{E}\left[X\right]}{a}$$

$$X \ = \ e^{sY}$$

# Chernoff Bound

USC
Viterbi

31

$$\mathbb{P}\left[Y \geq a\right] \quad \leq \quad \min_{s} e^{-sa}\mathbb{E}\left[e^{sY}\right]$$

Gaussian example

# Error Decay Rates

❑ Often more easily computable than exact probabilities

❑ Enable straightforward comparison across detectors

❑ Provide a measure for how far from asymptotic performance

▪ *When do asymptotics kick in?*

$$\text{Error rate}(\delta) = \lim_{n \to \infty} -\frac{1}{n} \log P_e(f)$$

what about fixed $n$?

$$f(y) = \text{decision rule}$$
$$= \{0, 1\}$$

# Error Decay Rates

## underwater acoustic communication



Simulation performance in UWA channels

Legend:
- 1-hop constructive - MLSD
- 1-hop constructive - DFE
- 1-hop hop destructive - MLSD
- 1-hop destructive - MLSD
- 2-relay cooperative, MLSD
- 2-relay cooperative, DFE
- 4-relay cooperative - MLSD
- 4-relay cooperative - DFE

BER vs total energy/symbol (dB)

TX → relays → RX

diversity

$$\text{decay rate} \quad 2 - 2\frac{\log\log P}{\log P}$$

cost: error propagation at relay

## Distributed Space–Time Cooperative Schemes for Underwater Acoustic Communications

Madhavan Vajapeyam, *Member, IEEE*, Satish Vedantam, Urbashi Mitra, *Fellow, IEEE*, James C. Preisig, *Member, IEEE*, and Milica Stojanovic, *Senior Member, IEEE*

# Error Rate for Bayes Rule

❑ Error rate:

$$\text{Error rate}(f) = \lim_{n \to \infty} -\frac{1}{n} \log r(f)$$

Bayes risk

❑ **Theorem**: error rate for the Bayes optimal rule

$$\text{Error rate(LRT)} = -\min_{0 \le \lambda \le 1} \log \underbrace{\sum_{y} (p_0(y))^\lambda (p_1(y))^{(1-\lambda)}}$$

Chernoff Information

❑ Not a function of the priors! $\pi_i$

# Neyman-Pearson Formulation

❑ Performance Measures:

$$\mathbb{P}[\hat{X} = 0 \mid X = 1] = \mathbb{P}_1[\hat{X} = 0] \text{ (Miss probability)}$$

$$\mathbb{P}[\hat{X} = 1 \mid X = 0] = \mathbb{P}_0[\hat{X} = 1] \text{ (False alarm probability)}$$

❑ Formulation: minimize miss probability while ensuring that false alarm probability is low

$$
\begin{aligned}
&\min_{f} && \mathbb{P}_1[\hat{X} = 0] \\
&\text{subject to} && \mathbb{P}_0[\hat{X} = 1] \leq \epsilon
\end{aligned}
$$

# Neyman Pearson Rule

❑ Optimal Decision Rule is a LRT:

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n > \tau \\ H_0 \text{ w.p. } \gamma & \text{if } L_n = \tau \\ H_1 & \text{if } L_n < \tau \end{cases}$$

❑ How to select parameters:

▪ Challenge when mismatched support and/or discrete RVs

threshold $\tau$ and randomization $\gamma$ unique solutions to
$$\epsilon = \mathbb{P}_0[L_n > \tau] + \gamma \mathbb{P}_0[L_n = \tau]$$

*randomization to achieve $P_F$ exactly*

# Gaussian Example

- ❑ Continuous valued RVs, matching support

- ❑ No randomization necessary

- ❑ False alarm rate determines threshold

$$\alpha = \mathbb{P}\left[\hat{X} = 1 | X = 0\right]$$

$$= Q\left(\frac{\tau' - \mu_0}{\sigma}\right)$$

$$\rightarrow \tau' = \mu_0 + \sigma Q^{-1}(\alpha)$$



likelihood function

| | pdf of +1 received |
| | pdf of  1 received |
| | zero threshold |
| | optimal threshold |

❑ NP best tradeoff between $P_F(\delta)$ and $P_D(\delta)$

# Chernoff-Stein Lemma

❑ Kullback-Leibler Divergence:

$$D(p||q) = \sum_y p(y) \log \frac{p(y)}{q(y)}$$

$$\mathbb{E}_0[L_n] = nD(p_0||p_1)$$
$$\mathbb{E}_1[L_n] = -nD(p_1||p_0)$$

Expectation of LLR is related to KL-Divergence

❑ Chernoff-Stein Lemma: Miss rate of NP rule is

$$\lim_{n \to \infty} -\frac{1}{n} \log \mathbb{P}_1[\hat{X} = 0] = D(p_0||p_1)$$

# Bayes Rule versus NP Rule

❏ Bayes rule

$$\text{Error rate(Bayes)} = -\min_{0 \leq \lambda \leq 1} \log \underbrace{\sum_y (p_0(y))^\lambda (p_1(y))^{(1-\lambda)}}_{\text{Chernoff Information}}$$

❏ Neyman Pearson rule

$$\text{Error rate(NP)} = \underbrace{D(p_0 \| p_1)}_{\text{Chernoff-Stein exponent}}$$

# SEQUENTIAL OBSERVATIONS

# Sequential Probability Ratio Tests

❑ Should you always use all of the data?

  ▪ Stop when confident!

❑ A Wald, *The Annals of Mathematical Statistics*, 1945

❑ Problem set up

  ▪ Samples $\mathbf{y}_m = [y_1, y_2, \cdots, y_m]$

$$L_m = \log \frac{p_0(\mathbf{y}_m)}{p_1(\mathbf{y}_m)}$$

$$\alpha = \text{false alarm rate}$$
$$\beta = \text{miss probability}$$

# SPRT solution

$$f(\mathbf{y}_m) = \begin{cases} H_0 & L_m \geq \frac{1}{B} \\ H_1 & L_m < \frac{1}{A} \\ \text{keep sampling} & \text{else} \end{cases}$$



time, number of samples

$$A \approx \log \frac{1-\beta}{\alpha}$$

$$B \approx \log \frac{\beta}{1-\alpha}$$

**same experiment**

# Now….

# Now….

# Now….

- ❑ Allow myself to take more observations and change experiment

  - ▪ Different experiments:  different sensors, different groupings

- ❑ Now, how to develop algorithms and analyze?



which action is more informative?

wumbo.net

# Now....

USC
Viterbi

- ❑ Allow myself to take more observations and change experiment
  - ▪ Different experiments: different sensors, different groupings
- ❑ Now, how to develop algorithms and analyze?



$$u_n = \text{experiment/observation mode}$$

$$y_n = \text{observation}$$

$$H = \text{true hypothesis}$$

$$\text{cost} = c\left(\{y_1, \cdots, y_{n-1}\}, \{u_1, \cdots, u_{n-1}\} \mid H\right)$$

which action is more informative?

wumbo.net

# YOU KNOW ``ACTIVE'' TESTING ALREADY

# Bayes Rule

$$\mathbb{P}\left[A\mid B\right] \;=\; \frac{\mathbb{P}\left[A,B\right]}{\mathbb{P}\left[B\right]}$$

$$=\; \frac{\mathbb{P}\left[B\mid A\right]\mathbb{P}\left[A\right]}{\mathbb{P}\left[B\right]}$$

# Monty Hall Problem

❑ Three doors: one car and two goats

❑ Pick a door!

# Monty Hall Problem

- ❑ Three doors: one car and two goats

- ❑ Pick a door!



- ❑ You select **A**

$$P(A|A) = \frac{1}{3}$$

car is behind door A

door A is chosen

# Now reveal a door

❑ Your door is still closed

❑ Do you change doors?

# Key Assumptions

❑ The car is equally likely to be behind all three doors

❑ The player is equally likely to pick one of the doors (independent of car's location_

❑ After player picks a door, the <span style="color:darkred">host must open a different door with a goat</span> and let player switch if they wish

❑ If selected door has car, host is equally likely to pick one of the goat doors

❑ KEY – non-uniform sample space/probabilities

# Decision Tree

MIT 6.042/18.062J  Leighton & Rubinfeld

# Decision Tree

# Now reveal a door

❑ Your door is still closed

❑ Do you change doors?



car is behind door A    door A is chosen

revealed    stay

$$P(A|C, A) = \frac{P(C|A, A)P(A|A)}{P(C|A)}$$

car is behind door A

door C is revealed

door A is chosen

# Now reveal a door

- ❑ Your door is still closed
- ❑ Do you change doors?



$$P(A|C, A) = \frac{\frac{1}{2}\frac{1}{3}}{1\frac{1}{3} + \frac{1}{2}\frac{1}{3}} = \frac{1}{3}$$

car is behind door A

door C is revealed

door A is chosen

# Now reveal a door

- Your door is still closed
- Do you change doors?



$$P(B|C, A) = \frac{1\frac{1}{3}}{1\frac{1}{3} + \frac{1}{2}\frac{1}{3}} = \frac{2}{3}$$

car is behind door B

door C is revealed

door A is chosen

# Monty Hall Problem

- ❑ This is a sequential decision-making problem

- ❑ The decision tree

$$\boxed{A} \xrightarrow{\text{revealed}} \boxed{C} \begin{array}{c} \xrightarrow{\text{stay}} \boxed{A} \quad \mathbb{P}\left[\text{win}\right] = \frac{1}{3} \\ \searrow_{\text{switch}} \boxed{B} \quad \mathbb{P}\left[\text{win}\right] = \frac{2}{3} \end{array}$$

- ❑ Action: switch, since the odds of winning are higher

- ❑ labels are arbitrary → optimal strategy: always switch

# Getting closer

- ❑ Some elements of our desired framework
  - ▪ Sequential decisions/observations
    - • A tree, but we are not pruning yet
  - ▪ Adversarial "game"

- ❑ Still one kind of experiment
  - ▪ One type of observation is not more informative than another

- ❑ Can we quantify informativeness?
  - ▪ How do we prune the tree?

# Wireless Body Area Sensing Network

ECG & Tilt sensor

SpO2 & Motion sensor

**Body Area Network**

Motion sensors

COMMUNICATIONS IN UBIQUITOUS HEALTHCARE

**KNOWME: A Case Study in Wireless Body Area Sensor Network Design**

Urbashi Mitra, B. Adar Emken, Sangwon Lee, and Ming Li, University of Southern California
Viktor Rozgic, Raytheon BBN Technologies
Thatte, TrellisWare Technologies, Inc.
dhan Vathsangam, Daphney-Stavroula Zois, Murali Annavaram, and Shrikanth Narayanan, Uni-
f Southern California
vorato, Stanford University and University of Southern California
pruijt-Metz and Gaurav Sukhatme, University of Southern California

May 2012, Vol. 50, No. 5

**IEEE Communications MAGAZINE**
www.comsoc.org

Free ComSoc Tutorial
Next-Gen Multi Gbps WLANs
See Page 11

- Ubiquitous Healthcare: Wireless Sensors, Devices and Solutions
- Optical Networking Advances
- Automotive Networking
- Smart Grid

IEEE COMMUNICATIONS SOCIETY
A Publication of the IEEE Communications Society

*Jovanov et al. Journal of NeuroEngineering and Rehabilitation* 2005

# What is my problem?

$$\mathbf{x} =$$

$$\mathbf{y} = f(\mathbf{x}, \mathbf{u}) \rightarrow \hat{\mathbf{x}}$$

observation    state, control

this is **active hypothesis testing**
for a **time-varying process**

# Problem Framework

- Sensor time-series (ECG, accelerometer, etc.) converted to features

- Each state indicated by a standard basis vector

- 
$$\mathbf{e}_i \quad = \quad [0, \cdots, 0, 1, 0 \cdots 0]$$

*i'th* component

Zois & **M**, TSP'17, ICASSP'14, ISIT'14, Globecom'14, Asilomar'13, GlobalSIP'13

Zois, Levorato &**M**, TSP'14, TSP'13

# Heterogeneity

❑ Different sensors are good at discriminating different states

❑ Chicken and egg problem…

# What is my problem?

❑ Goal: track temporal evolution of a discrete–time, finite–state Markov chain

❑ Design control (sensor allocation problem)

- Heterogeneous fidelity across sensors
- Heterogeneous costs across sensors
- **Optimize performance, minimize cost**

❑ Contrast to standard control problems:

- **control influences observations (not state)**

# POMDP System

partially observable Markov decision process (POMDP)

# Signal Model

❑ **System state**

▪ First order Markov process

$$\mathcal{X} = \{\mathbf{e}_1, \mathbf{e}_2, \ldots, \mathbf{e}_n\}$$
$$\mathbf{e}_i = [0, \ldots, 0, 1, 0, \ldots, 0]^T$$

❑ **Sensor features**

$$\mathbf{y}_k \big| \mathbf{e}_i, \mathbf{u}_{k-1} \sim \mathcal{N}(\mathbf{m}_i^{\mathbf{u}_{k-1}}, \mathbf{Q}_i^{\mathbf{u}_{k-1}})$$

control input *(can affect size, form, etc)*

▪ control is which sensor to listen to and for how long

▪ Validated by real world experiments

# Non-linear Decision Regions



Decision regions for bivariate Gaussians for six activities

❑ **Distinct** means and covariance matrices for **each** subject | personalized training

# Differential Entropy

- Definition

$$h(X) = -\int_{\mathcal{X}} f(x) \log f(x) dx \quad X \sim f(x)$$

- Properties

1. $h(X + c) = h(X)$ $c$ is a constant

2. $h(cX) = h(X) + \log |c|$ $c \neq 0, c$ is a constant

3. $X \sim \mathcal{N}(0, \sigma^2)$
   $\rightarrow h(X) = \frac{1}{2} \log \left(2\pi e \sigma^2\right)$ maximal differential entropy

4. $X$ is a mixed random variable $\rightarrow h(X) = -\infty$

# Bounds on estimation error

$$\mathbb{E}\left[\left(X - \hat{X}\right)^2\right] \geq \frac{1}{2\pi e} e^{2h(X)}$$

$$\hat{X} = \mathbb{E}[X] \quad \textit{MSE optimizing estimator}$$

$$\mathbb{E}\left[\left(X - \hat{X}|Y\right)^2\right] \geq \frac{1}{2\pi e} e^{2h(X|Y)}$$

$$\hat{X} = \mathbb{E}[X|Y] \quad \textit{MSE optimizing estimator}$$

these are the variances
differential entropy bounded by that of a Gaussian

# State Estimator

- <u>Minimize</u>: Mean-Square Error (MSE)

MMSE estimator $\quad \mathbf{x}_{k|k} \doteq \mathbb{E}\{\mathbf{x}_k | \mathcal{F}_k\}$

history of observations and control inputs

- MMSE estimator equals conditional belief (probability)



$$\mathbf{x}_{k|k} = \mathbf{p}_{k|k}$$

$$\mathbf{e}_i = [0, \ldots, 0, 1, 0, \ldots, 0]^T$$

$$\mathbf{p}_{k|k} = [p^1_{k|k}, p^2_{k|k}, \ldots, p^n_{k|k}]^T$$

with

$$p^i_{k|k} = P(\mathbf{x}_k = \mathbf{e}_i | \mathcal{F}_k)$$

- Designed a Kalman-like estimator (recursive/discrete states)

# Optimal Control Policy

❑ Control inputs sequence to optimize filter performance (MSE performance)

Cost function

$$J_\gamma = \mathbb{E}\left\{ \sum_{k=1}^{L} \mathrm{tr}\left( \boldsymbol{\Sigma}_{k|k}(\mathbf{y}_k, \mathbf{u}_{k-1}) \right) \right\}$$

filtering error covariance matrix

❑ Optimal solution via dynamic programming (DP)

optimal cost to go $= \displaystyle\min_{\mathbf{u}_{k-1} \in \mathcal{U}} \Big[$ current cost

$+$ expected future cost

Zois, Levorato, **M**, ICASSP 2013

# Include energy cost

Cost function

$$J = \mathbb{E}\left\{ \sum_{k=1}^{L} (1-\lambda)\mathsf{MSE}(\mathbf{y}_k, \mathbf{u}_{k-1}) + \lambda\mathcal{E}(\mathbf{u}_{k-1}) \right\}$$

trade–off
parameter

❑ **Partially observable** stochastic control problem: *determine control sequence to optimize trade–off between MSE performance and energy cost*

$$\min_{\mathbf{u}_0,\mathbf{u}_1,\ldots,\mathbf{u}_{L-1}} J$$

# Challenges of DP

- ❑ Curse of dimensionality
    - ▪ Predicted belief state drawn from uncountably infinite set
    - ▪ Control space can be exponentially large in N, K

- ❑ Non-linear POMDP

- ❑ expected future cost requires N–dimensional integration,   N = number of measurements

DP impractical for large-scale applications

# Goal & Approach

❑ **Goal**: determine

- Structural properties of the cost $-$ to $-$ go function

- Sufficient conditions to characterize optimal control

❑ **Assumptions**:

- discriminate between *two states*, $\mathbf{e}_1$ and $\mathbf{e}_2$

- Select 1 out of *N* available sensors (scalar measurements)

❑ Two hypotheses

$$\mathbf{p}_{k|k} = [p, 1 - p]$$

# Cost – to – go function properties

- ❑ Current cost

$$\ell(\mathbf{p}_{k|k-1}, \mathbf{u}_{k-1}) \doteq \mathbf{p}_{k|k-1}^T \mathbf{h}(\mathbf{p}_{k|k-1}, \mathbf{u}_{k-1})$$

- ❑ **Lemma**: current cost is concave function of $\mathbf{p}_{k|k-1}$

- ❑ **Theorem**: The cost – to – go function $\overline{J}_k(\mathbf{p}_{k|k-1})$ is a concave function of $\mathbf{p}_{k|k-1}$

$$k = L, L - 1, \ldots, 1$$

Zois, Levorato, **M**, GlobalSIP 2013

# Graphical interpretation

□ What does the **Theorem** really mean?



cost versus belief for different
controls/observation modes

# Graphical interpretation

❑ What does the **Theorem** really mean?



$\overline{J}_k(p)$

0      $p^*$      1

— $\mathbf{u}^\alpha$
— $\mathbf{u}^\beta$
— $\mathbf{u}^\gamma$

❑ Optimal policy has **threshold structure**

$$\mathbf{u}^{opt} = \begin{cases} \mathbf{u}^\gamma, & p \leqslant p^* \\ \mathbf{u}^\alpha, & p > p^* \end{cases}$$

well − known for
**linear POMDPs**
our system is **non-linear**

# Informativeness

❑ **Definition**: Given two conditional pdfs $f_\alpha$ and $f_\beta$ from $\mathcal{X}$ to $\mathcal{Y}$,

$f_\beta$ is *less informative than* $f_\alpha$ ( $f_\beta \leqslant_B f_\alpha$ ) if $\exists$

stochastic transformation $W : \mathcal{Y} \to \mathcal{Y}$

Blackwell
Ordering

$$f_\beta(\mathbf{y}|\mathbf{x}) = \int f_\alpha(\mathbf{z}|\mathbf{x}) W(\mathbf{z}; \mathbf{y}) d\mathbf{z}, \ \forall \mathbf{x} \in \mathcal{X}$$

# Informativeness

- **Fact**: Consider observation kernels $f(y|\mathbf{x}, \mathbf{u}^\alpha)$ and $f(y|\mathbf{x}, \mathbf{u}^\beta)$. If $f(y|\mathbf{x}, \mathbf{u}^\beta) \leqslant_B f(y|\mathbf{x}, \mathbf{u}^\alpha)$, then $\mathbf{u}^\alpha$ better than $\mathbf{u}^\beta$

  - Why? *Lower future cost* $V(p, \mathbf{u}^\alpha) \leqslant V(p, \mathbf{u}^\beta)$
  - Directly exploits the concavity of the cost-to-go function

- *Like a data processing inequality*
  - The stochastic transformation $W : \mathcal{Y} \to \mathcal{Y}$
    is processing the kernel $f(y|\mathbf{x}, \mathbf{u}^\alpha)$

# Data Processing Inequality

❑ Markov chains

$$p(x, y, z) \;=\; p(x)p(y|x)p(z|y)$$
$$p(x, z|y) \;=\; p(x|y)p(z|y)$$

❑ The inequality

$$X - Y - Z \;\rightarrow\; I(X;Y) \geq I(X;Z)$$
$$\rightarrow\; I(Y;Z) \geq I(X;Z)$$

▪ *processing Y cannot increase the information about X*

# Determining optimal control

Case I: *same mean, same variance*

Case II: *same mean, different variance*

Case III: *different mean, same variance*

- ❑ **Case II**: Blackwell ordering of observation kernels determines optimal control

- ❑ **Case III**: ordering of current cost is achieved by ordering of function of means $(m_1^{\mathbf{u}} - m_2^{\mathbf{u}})^2$

# Myopic Solution

❑ Optimal solution: expensive to determine over finite horizon

- Classical engineering fix: don't look too far into the future

❑ **Basic idea:** minimize one – step ahead cost

$$\mathbf{u}_{k-1}^{myopic} = \arg\min \ell(\mathbf{p}_{k|k-1}, \mathbf{u}_{k-1})$$

# Myopic Solution

❑ Current cost is concave with respect to $\mathbf{p}_{k|k-1}$ for 2 activity states and 1 measurement

- Policy has a **threshold structure also!**

$$\mathbf{u}^{myopic} = \begin{cases} \mathbf{u}^{\gamma}, & p \leqslant p^* \\ \mathbf{u}^{\alpha}, & p > p^* \end{cases}$$

$\ell(p, \mathbf{u}_{k-1})$



— $\mathbf{u}^{\alpha}$
— $\mathbf{u}^{\beta}$
— $\mathbf{u}^{\gamma}$

0     $p^*$     1

❑ This seems to be true for > 2 activity states and multi-dimensional measurement vectors (via *numerical validation*)

# Trade–off Curves

- ❑ Equal allocation: request same number of samples from each sensor

- ❑ Compared to equal allocation, energy gains as high as 60% for the same estimation/detection performance

# Summary

❑ Active hypothesis testing problem

▪ Individual's state is time-varying across time

▪ Allocate # measurements/which sensor (observation mode)

❑ Notion of informative observation modes

▪ (Blackwell ordering)

❑ Given belief for each state, we know which sensor to select

$$\mathbf{p}_{k|k-1} \rightarrow \mathbf{u}^{\alpha}$$

❑ How do we analyze performance?

# OPTIMAL DECAY RATE?

# Analysis of Interest

❑ Determining closed form probability of error intractable for WBAN case

- How to analyze so that we can determine design strategies/resource choices?

❑ How well does the approach work as the number of observations get large?

- Still interested in non-asymptotic/finite horizon performance

$$\lim_{N \to \infty} -\frac{1}{N}\mathbb{P}[\hat{X} = j | X \neq j] \quad \text{probability of error}$$

$$\text{subject to } \mathbb{P}[\hat{X} = j | X = j] \geq 1 - \epsilon$$

correct detection

# Let's go back to basics

❑ To find desired results, need to go simpler/abstract

❑ **Fixed** true hypothesis (not time-varying)

candidate hypotheses

$$
\begin{array}{llllllll}
h_1 & h_1 & h_1 & h_1 & h_1 \\
h_2 & h_2 & h_2 & h_2 & h_2 & h_2 & h_2 & & h_2 & h_2 & h_2 & h_2 & h_2 & h_2 \\
h_3 & h_3 & h_3 & h_3 & h_3 & h_3 & h_3 & & h_3 & h_3 & h_3 \\
h_4 & h_4 & h_4 & h_4 & h_4 & h_4 & h_4 \\
h_5 & h_5 & h_3 & h_5
\end{array}
$$

policies/experiments

$$
\begin{array}{llllllll}
u_1 & u_2 & u_2 & u_3 & u_2 & u_1 & u_3 & & u_2 & u_3 & u_3 & u_2 & u_2 & u_2 \\
\downarrow & \downarrow & \downarrow \\
y_1 & y_2 & y_3
\end{array}
$$

observations



Kartik, Nayyar & **M**, TAC'22, ISIT'20, ISIT'19, Asilomar'18
Kartik & **M**, TSP'22

# Recall: Neyman Pearson Rule

❑ Optimal Decision Rule is a LRT:

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n > \tau \\ H_0 \text{ w.p. } \gamma & \text{if } L_n = \tau \\ H_1 & \text{if } L_n < \tau \end{cases}$$

❑ How to select parameters:

▪ Challenge when mismatched support and/or discrete RVs

threshold $\tau$ and randomization $\gamma$ unique solutions to
$$\epsilon = \mathbb{P}_0[L_n > \tau] + \gamma \mathbb{P}_0[L_n = \tau]$$

*threshold choice determines NP rule*

# Near-Optimal Decision Rule

❑ Simpler Near-optimal Decision Rule:

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n \geq \tau \\ H_1 & \text{if } L_n < \tau \end{cases} \qquad \tau \approx nD(p_0||p_1)$$

a threshold test like optimal likelihood ratio test

❑ **Lemma**: miss probability probability for this decision rule

$$\mathbb{P}_1[\hat{X} = 0] \leq \exp(-\tau)$$
$$\approx \exp(-nD(p_0||p_1))$$

Large $\tau$ leads to high false-alarm probability
need to balance miss and false-alarm probabilities

# Moment Generating Function of LLR

- ❑ MGF of LLR:

$$\mu(s) = \mathbb{E}[\exp(-sL) \mid H_0]$$
$$= \sum_{y \in \mathcal{Y}} (p_0(y))^{1-s} (p_1(y))^s \qquad L = \log \frac{p_0(Y)}{p_1(Y)}$$

- ❑ Recall Chernoff Information

$$- \min_{0 \le \lambda \le 1} \log \sum_y (p_0(y))^\lambda (p_1(y))^{(1-\lambda)}$$

  - ▪ (and recall Chernoff bound)

- ❑ The idea: use new measures to drive hypothesis testing

# MGF of LLR – connections

- MGF of LLR:

$$\mu(s) = \mathbb{E}[\exp(-sL) \mid H_0]$$

$$= \sum_{y \in \mathcal{Y}} (p_0(y))^{1-s}(p_1(y))^s \qquad L = \log \frac{p_0(Y)}{p_1(Y)}$$

- Chernoff Information:

$$C(p_0 \| p_1) = - \min_{0 \le s \le 1} \log \mu(s)$$

- Kullback-Leibler Divergence:

$$D(p_0 \| p_1) = \lim_{s \to 0} -\frac{1}{s} \log(\mu(s))$$

# MGF of LLR – connections

- MGF of LLR:

$$\mu(s) = \mathbb{E}[\exp(-sL) \mid H_0]$$
$$= \sum_{y \in \mathcal{Y}} (p_0(y))^{1-s}(p_1(y))^s$$

$$L = \log \frac{p_0(Y)}{p_1(Y)}$$

- Chernoff Information:

$$C(p_0||p_1) = - \min_{0 \le s \le 1} \log \mu(s) \qquad \text{Bayes rate}$$

- Kullback-Leibler Divergence:

$$D(p_0||p_1) = \lim_{s \to 0} -\frac{1}{s} \log(\mu(s)) \qquad \text{NP rate}$$

# Renyi Entropy & Divergence

❑ Renyi Entropy

- generalizes entropy $H_\alpha(p) = \frac{1}{1-\alpha} \log \sum_{i=1}^{n} p_i^\alpha$

❑ Renyi Divergence

$$D_\alpha(p_0\|p_1) = \frac{1}{\alpha - 1} \log \left( \sum_{y \in \mathcal{Y}} (p_0(y))^\alpha (p_1(y))^{1-\alpha} \right)$$

❑ Renyi Divergence and MGF of LLR

$$D_{(1-s)}(p_0\|p_1) = -\frac{1}{s} \log \mu(s)$$

# Chernoff Bound and Renyi Divergence

- False-alarm probability bound using Chernoff bound:

$$\mathbb{P}_0[\hat{X} = 1] = \mathbb{P}_0[L_n < \tau] \leq e^{s\tau}\mathbb{E}_0[\exp(-sL_n)]$$
$$= e^{s\tau}(\mu(s))^n$$

- False-alarm decay rate:

$$-\frac{\log(\mathbb{P}_0[L_n < \tau])}{n} \geq \sup_{s \geq 0}\left(sD_{1-s}(p_0\|p_1) - s\tau/n\right)$$
$$= -\inf_{s \geq 0}\left(\log(\mu(s)) + s\tau/n\right)$$

# Example: Gaussian Likelihoods

- Null and Alternate Hypotheses:

$$H_0 : Y_n \sim \mathcal{N}(0, \sigma^2)$$
$$H_1 : Y_n \sim \mathcal{N}(\mu, \sigma^2)$$

$$p_0(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{y^2}{2\sigma^2}}$$

$$p_1(y) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

- Log-likelihood ratio is also Gaussian:

$$L_n = \sum_{k=1}^{n} \frac{\mu^2 - 2\mu Y_n}{2\sigma^2}$$

Mean: $\dfrac{\mu^2}{2\sigma^2}$ under $H_0$
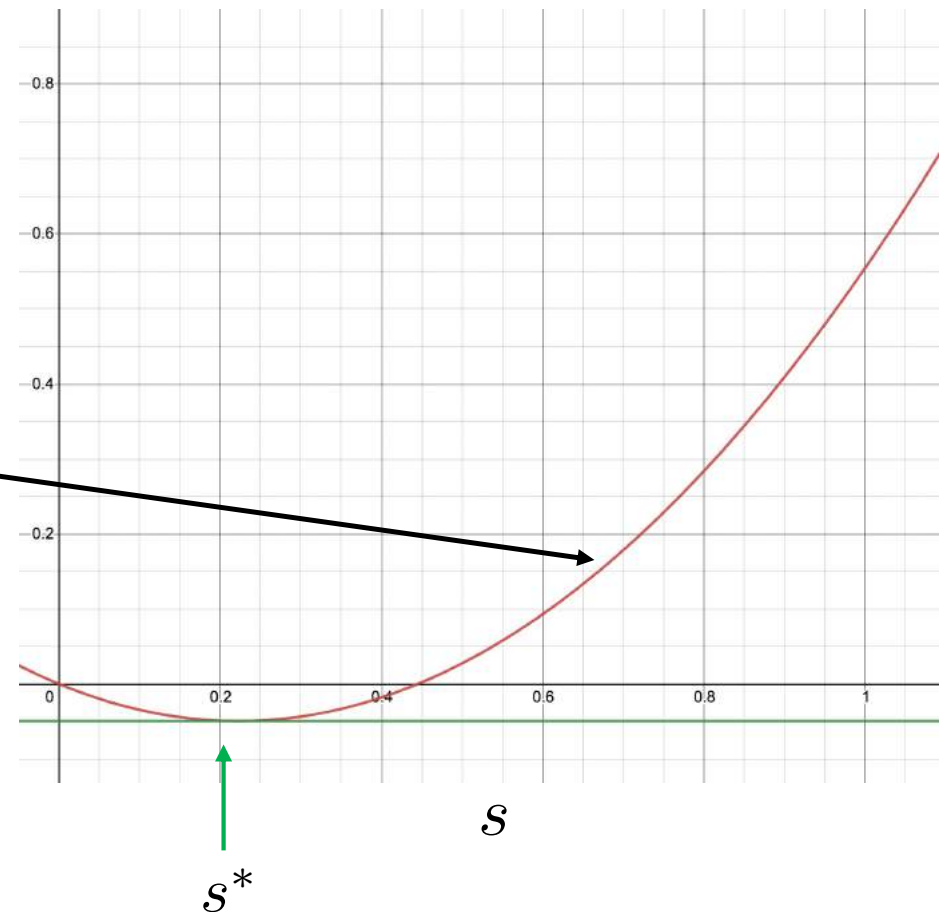
Variance: $\dfrac{\mu^2}{\sigma^2}$ under $H_0$

# Example: Gaussian Likelihoods

- MGF of negative LLR:

$$\mu(s) = \exp\left(\frac{-\mu^2 s}{2\sigma^2} + \frac{\mu^2 s^2}{2\sigma^2}\right)$$

$\inf_{s \geq 0}\left(\log(\mu(s)) + s\tau/n\right)$ can be obtained in closed form



$s$

$s^*$

# Example: Gaussian Likelihoods

❑ False-alarm decay rate:

$$-\frac{\log(\epsilon)}{n} \geq -\frac{\log(\mathbb{P}_0[L_n < \tau])}{n}$$

$$\geq -\inf_{s \geq 0} \left( \log(\mu(s)) + s\tau/n \right)$$

$$= -\inf_{s \geq 0} \left( \left( \frac{-\mu^2}{2\sigma^2} + \frac{\tau}{n} \right) s + \frac{\mu^2 s^2}{2\sigma^2} \right)$$

$$= \frac{\left( \frac{-\mu^2}{2\sigma^2} + \frac{\tau}{n} \right)^2}{\frac{4\mu^2}{2\sigma^2}} \qquad \text{if } \frac{\tau}{n} \leq \frac{\mu^2}{2\sigma^2}$$

$$\therefore \tau \leq \frac{\mu^2 n}{2\sigma^2} - \sqrt{\frac{2\mu^2 n \log(\frac{1}{\epsilon})}{\sigma^2}}$$

# Summary: Gaussian Likelihoods

- Decision-rule:

$$\hat{X} = \begin{cases} H_0 & \text{if } L_n \geq \tau \\ H_1 & \text{if } L_n < \tau \end{cases}$$

likelihood ratio test

- Miss probability lemma: large threshold desirable

$$\mathbb{P}_1[\hat{X} = 0] \leq \exp(-\tau)$$

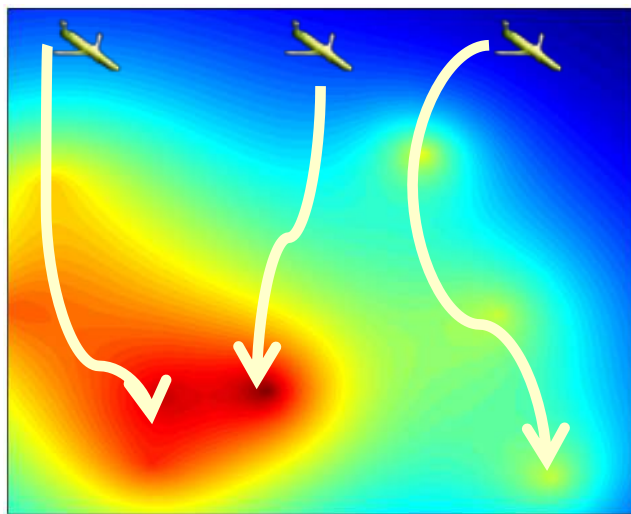- False-alarm probability: cannot have very large threshold

Sufficient to satisfy constraint

$$\tau \leq \frac{\mu^2 n}{2\sigma^2} - \sqrt{\frac{2\mu^2 n \log(\frac{1}{\epsilon})}{\sigma^2}}$$

asymptotically optimal error rate

non-asymptotic term

# Exploration-Exploitation



*exploration*
environment unknown

*collect observations*
learn

*exploitation*
focus on areas of interest

# Active Hypothesis Testing

EXPLORATION

candidate hypotheses

$h_1$ $h_1$ $h_1$ $h_1$ $h_1$

$h_2$ $h_2$ $h_2$ $h_2$ $h_2$ $h_2$ $h_2$

$h_3$ $h_3$ $h_3$ $h_3$ $h_3$ $h_3$ $h_3$

$h_4$ $h_4$ $h_4$ $h_4$ $h_4$ $h_4$ $h_4$

$h_5$ $h_5$ $h_3$ $h_5$

$h_2$ $h_2$ $h_2$ $h_2$ $h_2$ $h_2$

$h_3$ $h_3$ $h_3$

EXPLOITATION

$u_1$ $u_2$ $u_2$ $u_3$ $u_2$ $u_1$ $u_3$     $u_2$ $u_3$ $u_3$ $u_2$ $u_2$ $u_2$

policies/experiments

# focus on exploitation

# Active Hypothesis Testing – Prior Work

❑ Chernoff, H., 1959. Sequential design of experiments. *The Annals of Mathematical Statistics*

❑ Nitinawarat, S., Atia, G.K. and Veeravalli, V.V., 2013. Controlled Sensing for Multihypothesis Testing. *IEEE Transactions on Automatic Control*

  ▪ Considers decay rate of maximal error probability with fixed sample size

  ▪ Asymptotic optimality of stopping time formulation

❑ Naghshvar, M. and Javidi, T., 2013. Active sequential hypothesis testing. *The Annals of Statistics*

  ▪ POMDP formulation - Bounds on value function and asymptotic optimality

❑ Huang, B., Cohen, K. and Zhao, Q., 2019. Active Anomaly Detection in Heterogeneous Processes. *IEEE Transactions on Information Theory*

  ▪ Group testing-type approach and asymptotic optimality

We focus on **non-asymptotics**:
performance analysis and policy design

# Stopping Time Formulation

- ❑ **Classical approach**
- ❑ Perform experiments until confident – inconclusive declaration not allowed
- ❑ Stochastic time-horizon

- ❑ Minimize:

$$\mathbb{E}[N] + L \times \mathbb{P}[\hat{X} \neq X]$$
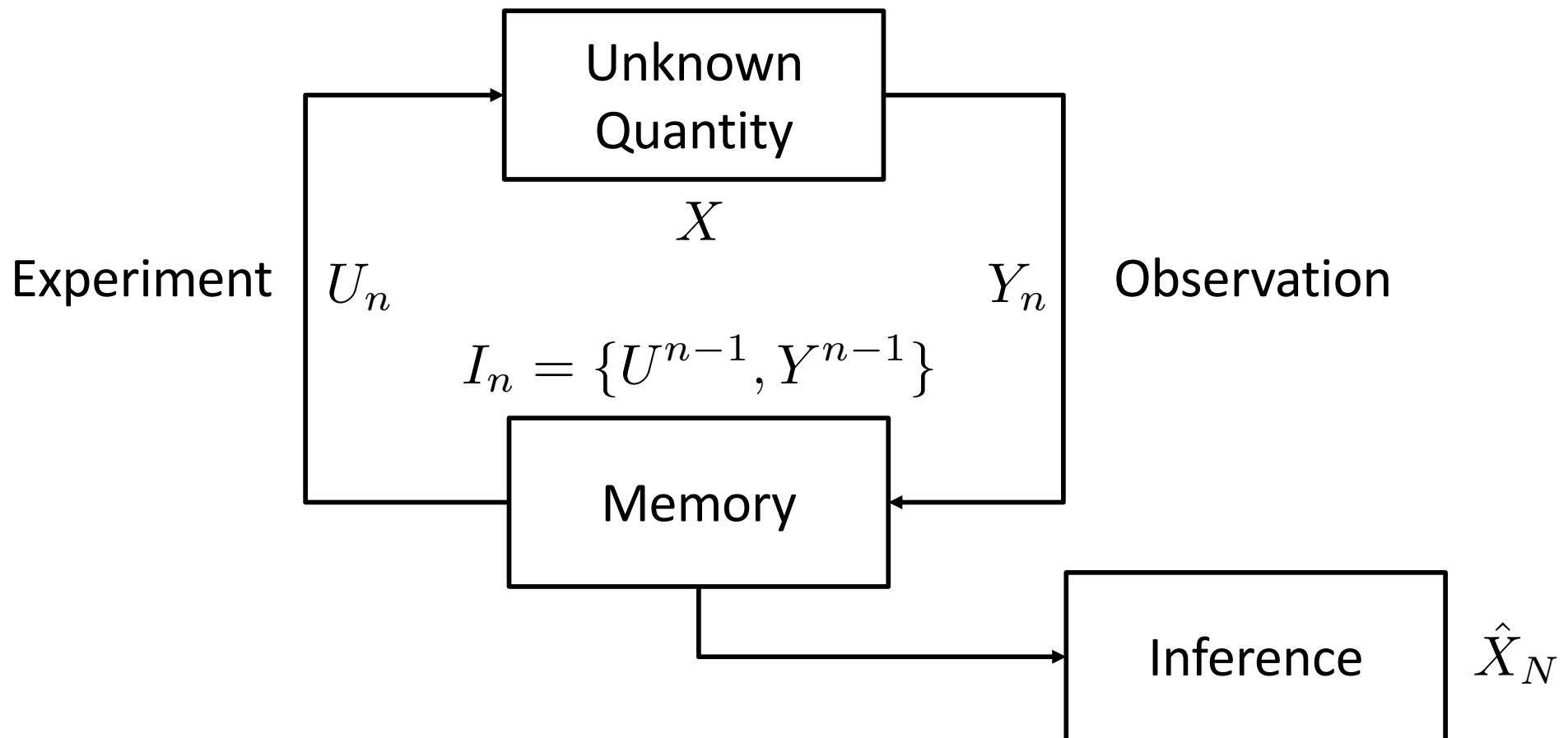
expected stopping time

fixed, usually very large

Bayesian error probability

**room for improvement in the *non-asymptotic* regime**

# Active Hypothesis Testing

❑ Access to multiple **experiments** and can select them in a data-driven fashion

# System Model

❑ Experiment Selection Strategy:



$$(U_1, Y_1) \quad (U_2, Y_2) \quad \ldots \quad (U_n, Y_n) \quad \ldots \quad (U_N, Y_N) \quad \hat{X}_N$$

Past information $I_n$ $\qquad\qquad$ $U_n = g_n(I_n, \text{random}) \in \mathcal{U}$

Observation $Y_n$ independent of past given $U_n$ and $X$

❑ Inference Strategy: infer after gathering all data – may declare inconclusive if necessary

$$\hat{X}_N = f(I_{N+1}, \text{random}) \in \mathcal{X} \cup \{\varnothing\}$$

# System Model

❑ Observations:

$$\mathbb{P}[Y = y \mid X = i, U_n = u] = p_i^u(y)$$

$$Y \in \mathcal{Y}$$
Finite alphabet

↑      ↑      ↑

Observation    Experiment    Likelihood functions

Observation $Y_n$ independent of past given $U_n$ and $X$

# Neyman-Pearson Formulation (P1)

❑ Incorrect conclusion: very expensive – must be avoided

$$\gamma_N = \mathbb{P}^{f,g}\big[\cup_{i\in\mathcal{X}}\{\hat{X}_N = i, X \neq i\}\big]$$

**Misclassification probability**: Probability of making an incorrect conclusion

Misclassification probability 0 if always declare inconclusive

❑ **Correct inference**: need to make correct inference with sufficiently large probability

$$\psi_N(i) \doteq \mathbb{P}^{f,g}\big[\hat{X}_N = i \mid X = i\big]$$

Correct inference probability of type-$i$

# Neyman-Pearson Formulation (P1)

❑ Optimization Problem:

$$\min_{f \in \mathcal{F}, g \in \mathcal{G}} \gamma_N$$

$$\text{subject to} \quad \psi_N(i) \geq 1 - \epsilon_N, \ \forall i \in \mathcal{X}$$

Infimum value:

$$\gamma_N^*$$

among all strategies that make correct inference with high probability, pick those that misclassify least

symmetric formulation

# Symmetric Cases

misclassification probability

$$\gamma_N = \mathbb{P}^{f,g}\left[\cup_{i \in \mathcal{X}}\{\hat{X}_N = i, X \neq i\}\right]$$
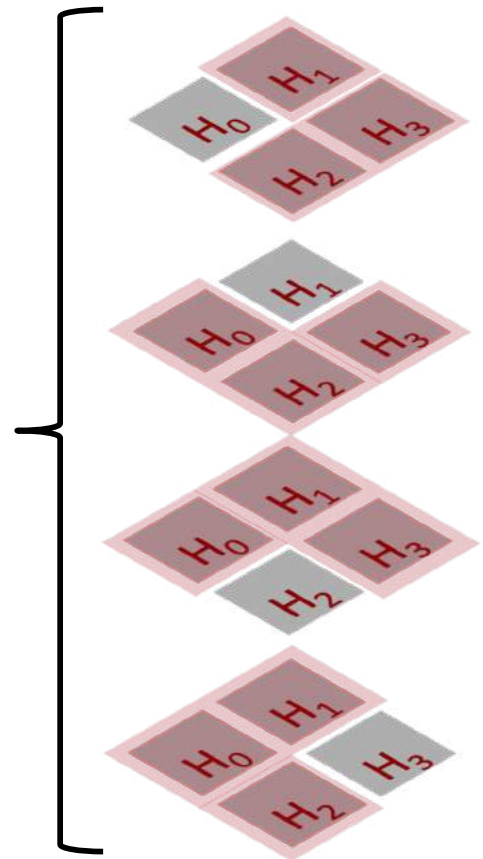
**P1**

symmetric
formulation

$$\min_{f \in \mathcal{F}, g \in \mathcal{G}} \quad \gamma_N$$

subject to $\quad \psi_N(i) \geq 1 - \epsilon_N, \; \forall i \in \mathcal{X}$

correct inference probability

❑ Incorrect conclusion: focus on a particular hypothesis

$$\phi_N(i) \doteq \mathbb{P}^{f,g}[\hat{X}_N = i \mid X \neq i]$$

Incorrect inference probability of type-$i$
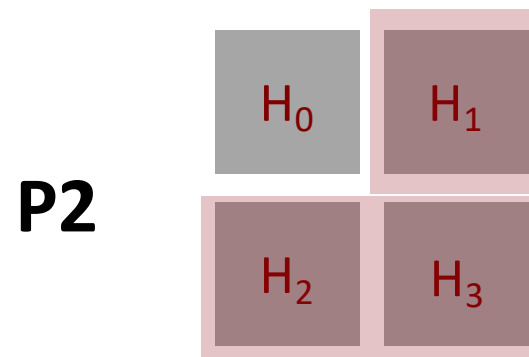
Probability of incorrectly inferring hypothesis $i$

❑ Correct inference: need to make correct inference with sufficiently large probability

$$\psi_N(i) \doteq \mathbb{P}^{f,g}[\hat{X}_N = i \mid X = i]$$

Correct inference probability of type-$i$

# Asymmetric Case

**P2**

| | |
|---|---|
| $H_0$ | $H_1$ |
| $H_2$ | $H_3$ |

asymmetric

$$H_0 \text{ versus } \{H_1, H_2, H_3\}$$

# Composite Test

- $H_i$ is a single hypothesis

- $H_i^c$ is all other hypotheses

$$H_i \text{ versus } H_i^c$$

$$H_i^c = \{H_0, H_1 \cdots H_{i-1}, H_{i+1}, \cdots H_M\}$$

$$\phi_N(i) \doteq \mathbb{P}^{f,g}[\hat{X}_N = i \mid X \neq i]$$

$$= \mathbb{P}[\hat{X}_N = i | H_i^c] \quad \text{incorrect inference}$$

$$\psi_N(i) \doteq \mathbb{P}^{f,g}[\hat{X}_N = i \mid X = i]$$

$$= \mathbb{P}[\hat{X}_N = i | H_i] \quad \text{correct inference}$$

# Neyman-Pearson Formulation (P2)

❑ Optimization Problem:

$$\min_{f \in \mathcal{F}, g \in \mathcal{G}} \quad \phi_N(i)$$

$$\text{subject to} \quad \psi_N(i) \geq 1 - \epsilon_N$$

Infimum value:
$$\phi_N^*(i)$$

Simple Null $\{X = i\}$ vs Composite Alternate $\{X \neq i\}$

Problem (P2) is easier to analyze
P2 will get us to solving P1

# Neyman-Pearson Formulation (P2)

❑ Asymmetric Hypothesis Test:

Fix experiment selection strategy $g$
and view as single-shot hypothesis testing problem

$$(U_1, Y_1) \quad (U_2, Y_2) \quad \ldots \quad (U_n, Y_n) \quad \ldots \quad (U_N, Y_N)$$

$$I_{N+1}$$

$$P_{N,i}^g(\mathcal{I}_{N+1}) \qquad\qquad\qquad\qquad Q_{N,i}^g(\mathcal{I}_{N+1})$$
$$\mathbb{P}^g[I_{N+1} = \mathcal{I}_{N+1} \mid X = i] \qquad\qquad \mathbb{P}^g[I_{N+1} = \mathcal{I}_{N+1} \mid X \neq i]$$

test if $I_{N+1}$ comes from $P$ or $Q$

**asymmetric** formulation

# Asymmetric vs Symmetric Cases

**P2**

| $H_0$ | $H_1$ |
|-------|-------|
| $H_2$ | $H_3$ |

asymmetric
$H_0$ versus $\{H_1, H_2, H_3\}$

**P1**

symmetric
formulation



$$\gamma_N = \mathbb{P}^{f,g}[\cup_{i \in \mathcal{X}}\{\hat{X}_N = i, X \neq i\}]$$

$$\min_{f \in \mathcal{F}, g \in \mathcal{G}} \gamma_N$$

subject to $\quad \psi_N(i) \geq 1 - \epsilon_N, \ \forall i \in \mathcal{X}$

❑ Confidence Level:

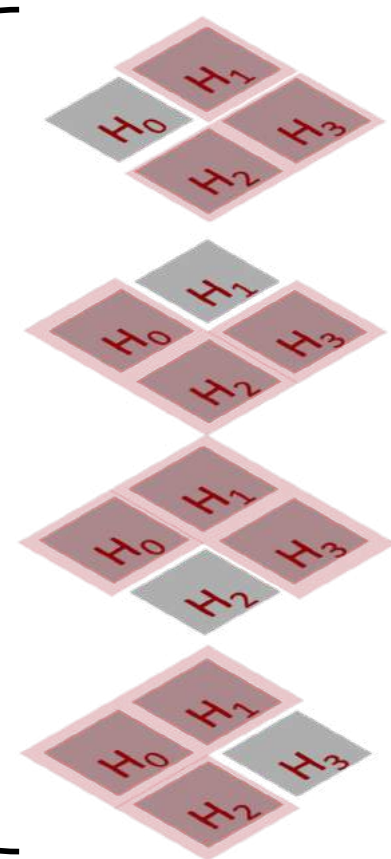$$\mathcal{C}_i(\rho) \doteq \log \frac{\rho(i)}{1 - \rho(i)}$$

$$\rho_n(i) = \mathbb{P}[X = i \mid U_{1:n-1}, Y_{1:n-1}]$$

Posterior belief

$i$ versus **not** $i$

❑ Expected Confidence Rate: Average Kullback-Leibler Divergence of the Asymmetric Hypothesis Test

$$J_N^g(i) \doteq \frac{1}{N}\mathbb{E}_i^g\left[\mathcal{C}_i(\rho_{N+1}) - \mathcal{C}_i(\rho_1)\right] = \frac{1}{N}\mathbb{E}_i^g\left[\log \frac{P^g(I_{N+1})}{Q^g(I_{N+1})}\right]$$

# Useful Information-theoretic Quantities

❑ **Max-min KL-Divergence**

$$D^*(i) \doteq \max_{\alpha \in \Delta \mathcal{U}} \min_{j \neq i} \sum_u \alpha(u) D(p_i^u \| p_j^u)$$

$$= \min_{\beta \in \Delta \tilde{\mathcal{X}}_i} \max_{u \in \mathcal{U}} \sum_{j \neq i} \beta(j) D(p_i^u \| p_j^u)$$

- Distributions over set of experiments: $\Delta \mathcal{U}$

- Max-minimizer: $\alpha^{i^*}$

- Distributions over set of alternate hypotheses: $\Delta \tilde{\mathcal{X}}_i$

- Min-maximizer: $\beta^{i^*}$

# Max-min Divergence

❑ Max-min KL-Divergence

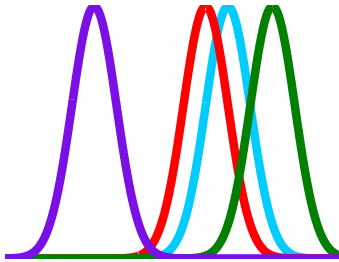$$D^*(i) \doteq \max_{\alpha \in \Delta\mathcal{U}} \boxed{\min_{j \neq i} \sum_u \alpha(u) D(p_i^u \| p_j^u)}$$

$$\alpha(u) = \mathbb{P}\left[\text{select experiment } u\right]$$

- Distributions over set of experiments: $\Delta\mathcal{U}$

- Max-minimizer: $\alpha^{i*}$

best probability distribution for hypothesis $i$

given $\alpha$, averaging over all experiments which two hypotheses yield the smallest divergence?
→ hardest to distinguish

# Max-min optimization

- Distributions over set of experiments: $\Delta \mathcal{U}$

- Max-minimizer: $\alpha^{i^*}$

we want to select the experiment that maximally separates the distributions for each hypothesis

# Min-Max optimization

❑ Equivalent optimization

$$D^*(i) \doteq \max_{\alpha \in \Delta\mathcal{U}} \min_{j \neq i} \sum_u \alpha(u) D(p_i^u || p_j^u)$$

$$= \min_{\beta \in \Delta\tilde{\mathcal{X}}_i} \boxed{\max_{u \in \mathcal{U}} \sum_{j \neq i} \beta(j) D(p_i^u || p_j^u)}$$
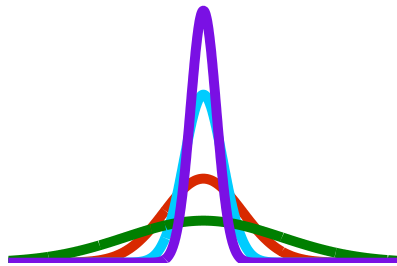
- Distributions over set of alternate hypotheses: $\Delta\tilde{\mathcal{X}}_i$

- Min-maximizer: $\beta^{i^*}$

given the best u , which priors make the two easiest hypothesis hard to distinguish

worst prior probability distribution for the other null hypotheses
*prior on hypotheses*

# Min-max optimization

- Distributions over set of alternate hypotheses: $\Delta \tilde{\mathcal{X}}_i$

- Min-maximizer: $\beta^{i*}$

the adversary wants to maximize the ``prior'' of the wrong hypothesis closest to the true hypothesis
P[purple] <<< P[blue]

# Data Processing Inequality

❑ Markov chains

$$p(x, y, z) = p(x)p(y|x)p(z|y)$$
$$p(x, z|y) = p(x|y)p(z|y)$$

❑ The inequality

$$X - Y - Z \rightarrow I(X; Y) \geq I(X; Z)$$
$$\rightarrow I(Y; Z) \geq I(X; Z)$$

▪ *processing Y cannot increase the information about X*

# DPI for Divergence

❑ Channel $\qquad X \to \boxed{p_{y|x}} \to Y$

❑ Two input distributions:
$$\text{if } X \sim p_X \text{ then } Y \sim p_Y$$
$$\text{if } X \sim q_X \text{ then } Y \sim q_Y$$

❑ DPI: $D(p_x\|q_x) \geq D(p_y\|q_y)$

- Processing the observation makes it more challenging to determine whether it came from p or q

❑ $p_{y|x}$ can be deterministic $\quad Y = \mathbf{1}_{\mathcal{A}}(X)$ for event $\mathcal{A}$

$$Y \sim \text{Ber with probability } \mathbb{P}(\mathcal{A}) \text{ or } \mathbb{Q}(\mathcal{A})$$

$$D(p_x\|q_x) \geq D\left(\text{Ber}(\mathbb{P}(\mathcal{A})\|\text{Ber}(\mathbb{Q}(\mathcal{A}))\right)$$

# Asymmetric Converse (P2)

❑ Weak converse: using DPI for binary hypothesis testing

$$-\frac{1}{N}\log\phi_N(i) \leq J_N^g(i) + \Theta(1/N) \leq D^*(i) + \Theta(1/N)$$

❑ Asymptotically optimal strategies: Using Chernoff bound

▪ achievability

$$-\frac{1}{N}\log\phi_N^*(i) > D^*(i) - \Theta(1/\sqrt{N})$$

# Chernoff's Strategy

❑ For asymmetric formulation, *i* specified

❑ Randomly select experiment, open-loop from distribution

Set of all distributions on experiments

$$\boldsymbol{\alpha}_i^* := \arg \max_{\boldsymbol{\alpha} \in \Delta\mathcal{U}} \min_{j \neq i} \sum_u \alpha_u D(p_i^u \| p_j^u)$$

Kullback-Leibler divergence

❑ For symmetric formulation,

- select most likely *i* based on data
- For most likely I, use $\alpha_i^*$ above

❑ Other works use a similar approach

# Asymmetric Achievable Strategy (P2)

$$f(\rho_{N+1}) =$$

$$\begin{cases} i & \text{if } \mathcal{C}_i(\rho_{N+1}) - \mathcal{C}_i(\rho_1) \geq \theta \\ \varnothing & \text{otherwise.} \end{cases}$$

Threshold based inference strategy

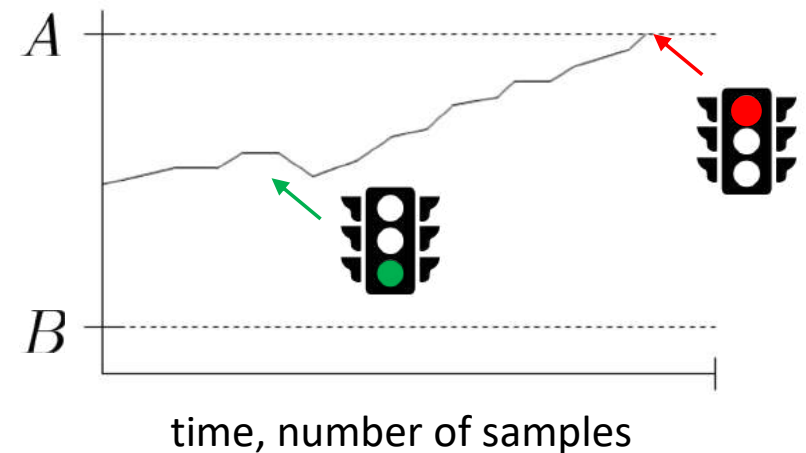Randomly select experiment with distribution $\alpha^{i*}$

Experiment selection strategy

# Asymmetric Achievable Strategy (P2)

$$f(\rho_{N+1}) =$$

$$\begin{cases} i & \text{if } \mathcal{C}_i(\rho_{N+1}) - \mathcal{C}_i(\rho_1) \geq \theta \\ \varnothing & \text{otherwise.} \end{cases}$$

Threshold based inference strategy



time, number of samples

recall SPRT:
stop if confident enough

Randomly select experiment with distribution $\alpha^{i*}$

Experiment selection strategy

# Achievable Strategy

❑ Note that achievable strategy is

- Data driven in inference (hypothesis selection)
    - Confidence function is a function of the data
- Randomized in experiment selection
- Devised to prove asymptotic results of best possible strategy

# Symmetric Converse (P1)

❑ Converse: use total probability theorem and the converse for (P2)

$$-\frac{1}{N}\log\gamma_N = -\frac{1}{N}\log\left(\sum_i \mathbb{P}[X \neq i]\phi_N(i)\right)$$

$$\leq \min_i D^*(i) + \Theta(1/N)$$

# Symmetric Achievability (P1)

❑   Achievability: a variant of the previous strategy

$$f(\rho_{N+1}) =$$

$$\begin{cases} i & \text{if } \mathcal{C}_i(\rho_{N+1}) - \mathcal{C}_i(\rho_1) \geq \theta \\ \varnothing & \text{otherwise.} \end{cases}$$

Threshold based inference strategy

Current most-likely hypothesis: $\hat{i}$

Randomly select experiment with distribution $\alpha^{\hat{i}*}$

Experiment selection strategy

# Optimal Error Rates

❑ **Theorem:** Chernoff-Stein Exponent for Asymmetric case (P2):

$$\lim_{N \to \infty} -\frac{1}{N} \log \phi_N^*(i) = D^*(i)$$

❑ **Theorem:** Chernoff-Stein Exponent for Symmetric case (P1):

$$\lim_{N \to \infty} -\frac{1}{N} \log \gamma_N^* = \min_{i \in \mathcal{X}} D^*(i)$$

# Active Experiment Selection Strategy

❑ MGF of LLR: now depends on the experiment

$$\mu_j^i(u, s) \doteq \mathbb{E}_i \exp\left(-s \log \frac{p_i^u(Y)}{p_j^u(Y)}\right)$$

❑ MGF based metric for experiment selection:

$$\mathcal{M}_i(u, \rho, s) \doteq \frac{\sum_{j \neq i}(\rho(j))^s \mu_j^i(u, s)}{\sum_{j \neq i}(\rho(j))^s} \qquad s_N \doteq \sqrt{\frac{2 \log \frac{M}{\epsilon_N}}{NB^2}}$$

Select the experiment $u \in \mathcal{U}$ that minimizes $\mathcal{M}_i(u, \rho_n, s_N)$

# Performance Guarantees

□ **Theorem**: the experiment selection strategy is asymptotically optimal and achieves significantly better performance in the non-asymptotic regime

- $s = s_N$ chosen ``just right'' so the right sums converge
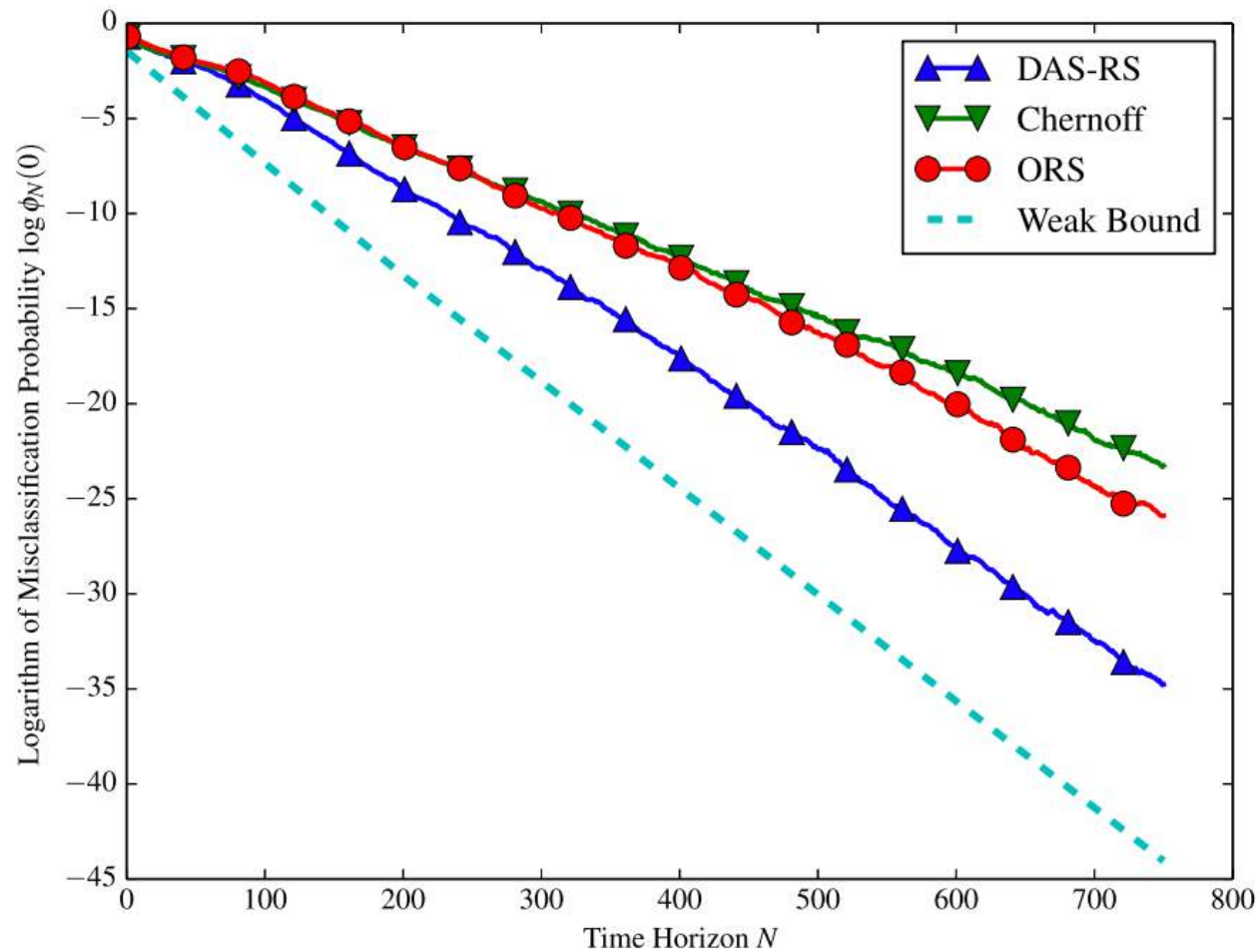
# Some Finite Horizon results

- ❑ For the general case (we will specialize to anomaly detection)

- ❑ Determine a Chernoff bound for active experiment selection

- ❑ Key:  bounded LLRs

$$\left\| \log \frac{p_0^u(Y)}{p_1^u(Y)} \right\| < B$$

bounded variables are sub-Gaussian

- ❑ Can determine bound and optimized threshold

# Numerical Results

# Some Takeaways

❑ For binary hypothesis testing, select the experiment with largest KL-Divergence

- Exploitation does not need to be active

- NOT always true for M-ary testing (multiple alternatives)

❑ For M-ary case, we care about the event
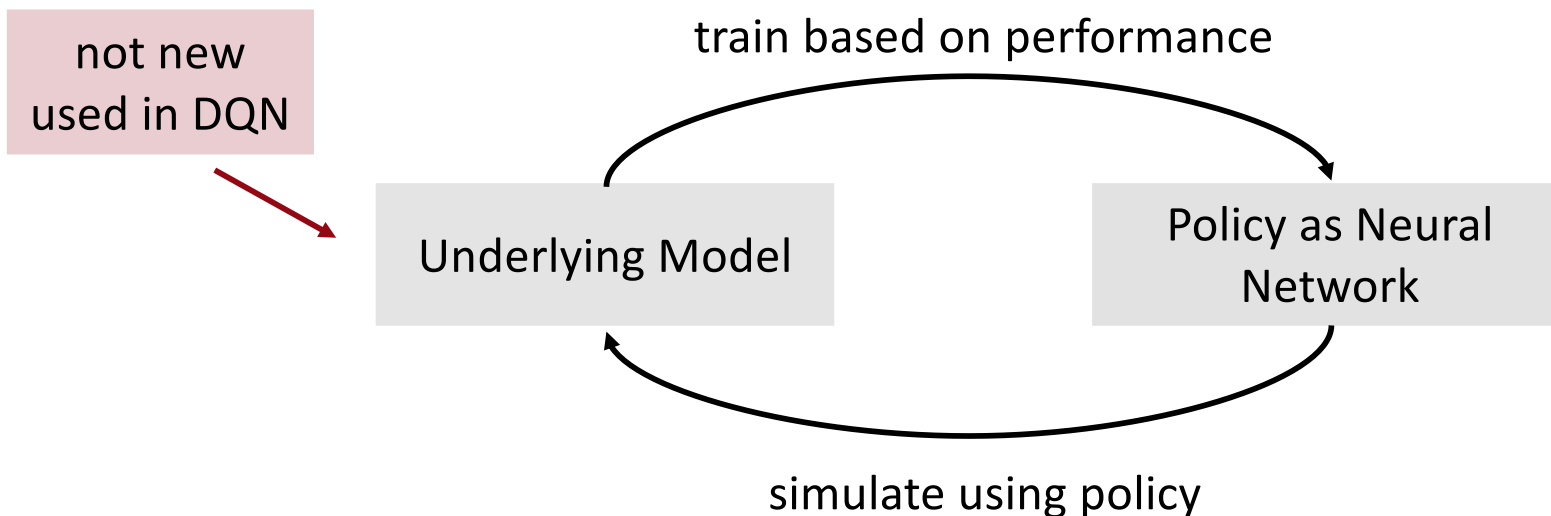$$\min_{j \in \text{alt. hyp}} L_n^j \geq \tau$$

Pairwise LLR for each alternate must exceed the threshold

❑ Similar achievability bounds can be derived in this case – these achievability bounds lead to our MGF based scheme

# Neural Networks as Policy Optimizers

❑ Consider the following framework

- DNNs as policy optimizers

- Simulate underlying model, generate data, evaluate performance

- With simulated data, train DNN via gradient descent

train based on performance

not new
used in DQN

Underlying Model

Policy as Neural
Network

simulate using policy

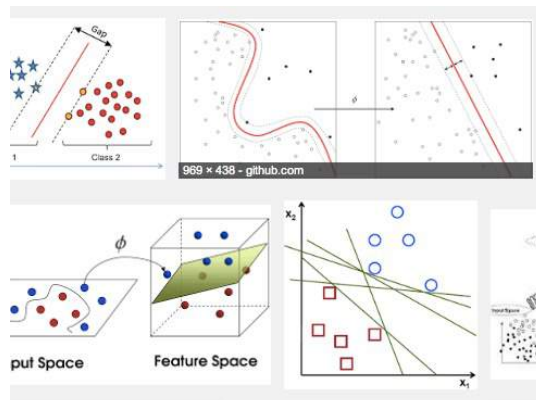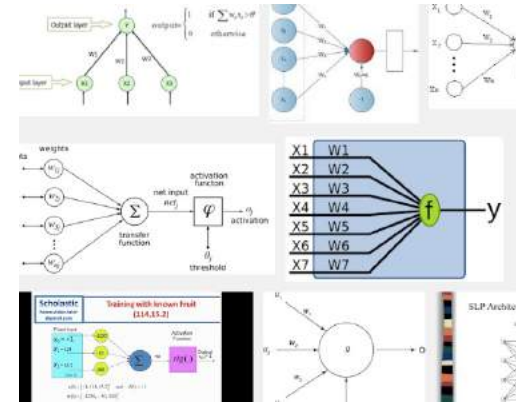Q: How to properly design NNs for experiment selection and classification?

# Third Wave of NN

ANSACTIONS ON COMMUNICATIONS, VOL. 43, NO. 2/3/4, FEBRUARY/MARCH/APRIL 1995

## Adaptive Receiver Algorithms for Near-Far Resistant CDMA

Urbashi Mitra, *Member, IEEE* and H. Vincent Poor, *Fellow, IEEE*

single layer perceptron

IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS, VOL. 12, NO. 9, DECEMBER

## Neural Network Techniques for Adaptive Multiuser Demodulation
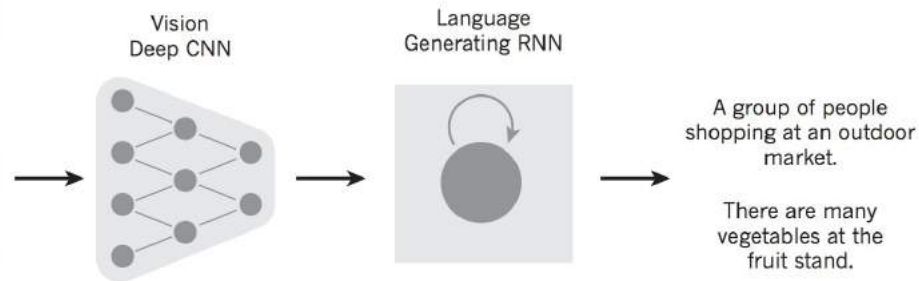
U. Mitra and H. Vincent Poor, *Fellow, IEEE*

support vector machine

**NOW: COMPUTATIONAL HORSEPOWER & NEW ANALYSIS TOOLS**

# Architecture Challenge
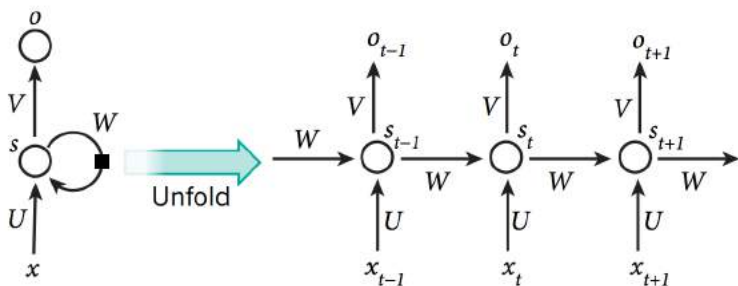
❑ Theoretically, neural networks are universal approximators

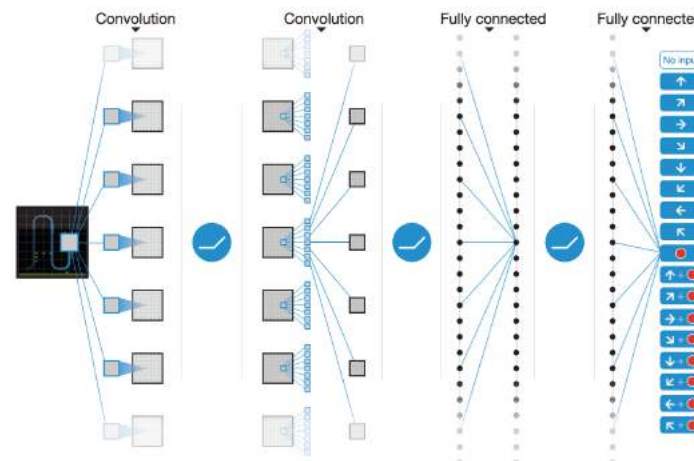❑ Challenge is finding the right architecture



Source: *Deep learning, Nature*

Convolutional Neural Network and Recurrent Neural Network for caption generation



Recurrent Neural Network
Source: *Deep learning, Nature*

Deep Q Network for reinforcement learning

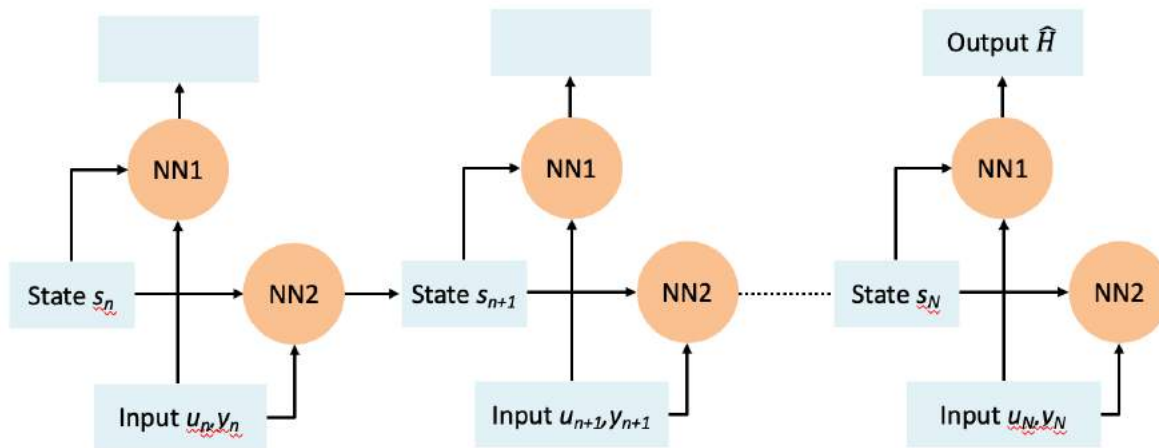Source: *Human-level control through deep reinforcement learning, Nature*

# Design Goals

use insights from information and control theory
to design architecture and features

❑ Deep reinforcement learning is an adaptation of Q-learning

- examined Recurrent Neural Networks and Q-Networks

- Q-Networks learn efficient query selection policies

Kartik,  Sabir, **M** & Natarajan, *Allerton'18*

# Recurrent Neural Network



Learns to classify

- Sequentially provide query-observation pairs

- After $N$ time steps, guess hypothesis

- If correct, 0 loss and 1 otherwise

- BUT

  • Fails to learn policy

  • Backpropagation has numerical stability issues

# Solution: Deep Q-Network

Represent Q values as a neural network vs a matrix

❑ Cannot simply assign Q-value updates

Fit Q-value update to network with MSE loss using gradient descent

❑ Optimize loss using gradient descent

$$\mathrm{MSE} = ||\mathrm{DQN}(\rho) - Q'(\rho)||^2$$

❑ Issues

- Belief space infinite – ε exploration
- Numerical stability issues/normalization



$\rho(1)$, $\rho(2)$, $\rho(3)$ Belief vector

Fully connected hidden layers

$Q(u_1)$, $Q(u_2)$ Q-values

Agent

$R_n$ Simulate belief update $\rho_n$

$\rho_{n+1}$ $U_n$

# Numerical Comparison

❑ Extrinsic  Jensen-Shannon Divergence (EJS):

$$EJS(\boldsymbol{\rho}, u) = \mathbb{E}[\mathcal{C}(F(\boldsymbol{\rho}, u, \mathbf{Y})) - \mathcal{C}(\boldsymbol{\rho})]$$

- Greedy: select experiment that maximizes EJS

- Naghshvar & Javidi, *Extrinsic Jensen-Shannon divergence with application in active hypothesis testing,* ISIT, 2012

❑ Open loop verification (OPE):

- Explore using EJS

- If $\rho_i > 0.7$ (confidence) select experiment with distribution

  – Recall Chernoff approach $\qquad\qquad \boldsymbol{\alpha}_i > 0.7$

- Naghshvar and Javidi, *Active Sequential Hypothesis Testing,* The Annals of Statistics, 2013

# Numerical Comparison

❏ Our adaptive best-response heuristic (KLZ):

- Explore using EJS

- If $\rho_i > 0.7$, select action from support (i) that maximizes $J_i(\boldsymbol{\rho}, u)$

$$J_i(\boldsymbol{\rho}, u) = \sum_{j \neq i} \frac{\rho_j}{1-\rho_i} D\left(p_i^u \| p_j^u\right)$$

❏ Compare to our final general strategy

❏ Compare these three strategies to DQN

- EJS work states conditions under which EJS is asymptotically optimal

- Example selected to violate those conditions

# Additional Queries

$$\epsilon = 10^{-7}$$

|       | $y = 0$ | $y = 1$ |
|-------|---------|---------|
| $h_0$ | 0.8     | 0.2     |
| $h_1$ | 0.2     | 0.8     |
| $h_2$ | 0.8     | 0.2     |

$u_1$

|       | $y = 0$ | $y = 1$ |
|-------|---------|---------|
| $h_0$ | 0.8     | 0.2     |
| $h_1$ | 0.8     | 0.2     |
| $h_2$ | 0.2     | 0.8     |

$u_2$

|       | $y = 0$      | $y = 1$    |
|-------|--------------|------------|
| $h_0$ | 0.8          | 0.2        |
| $h_1$ | $1 - \epsilon$ | $\epsilon$ |
| $h_2$ | 0.8          | 0.2        |

$u_3$

|       | $y = 0$      | $y = 1$    |
|-------|--------------|------------|
| $h_0$ | 0.8          | 0.2        |
| $h_1$ | 0.8          | 0.2        |
| $h_2$ | $1 - \epsilon$ | $\epsilon$ |

$u_4$

KL-divergence is asymmetric

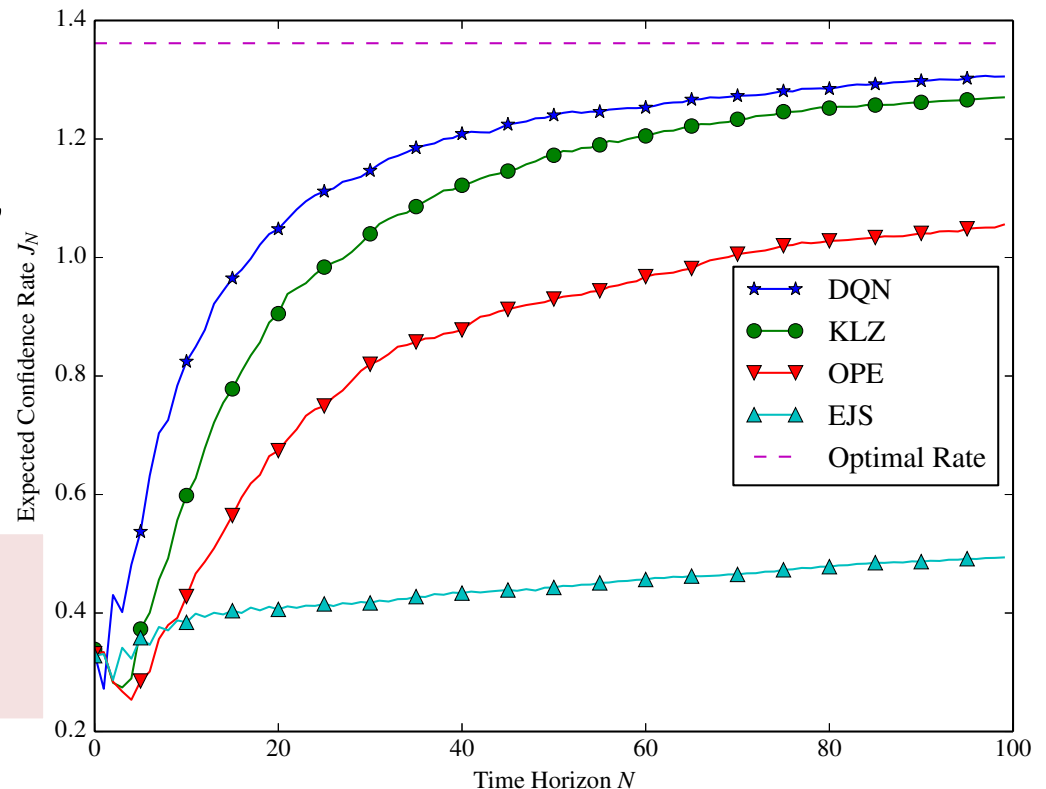# Deep Q Network for Active Classification

❑ Optimal strategy computationally expensive

▪ Infinite state space

❑ New measure from theoretical analysis:  structural properties

❑ **KLZ** close to optimal rate

❑ **OPE** asymptotically optimal, but very slow convergence
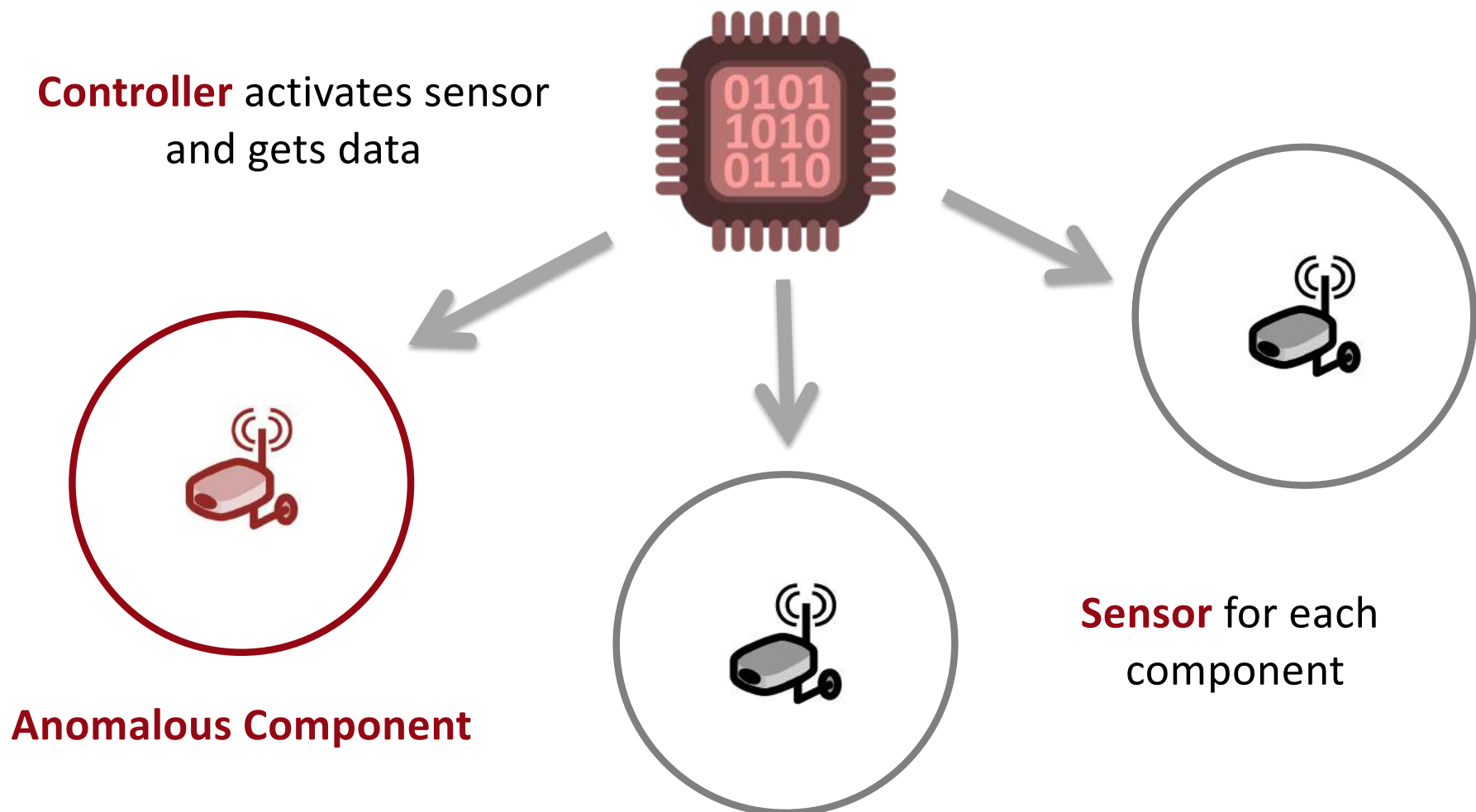
❑ **EJS** not optimal

**DQN** learns  the best policy
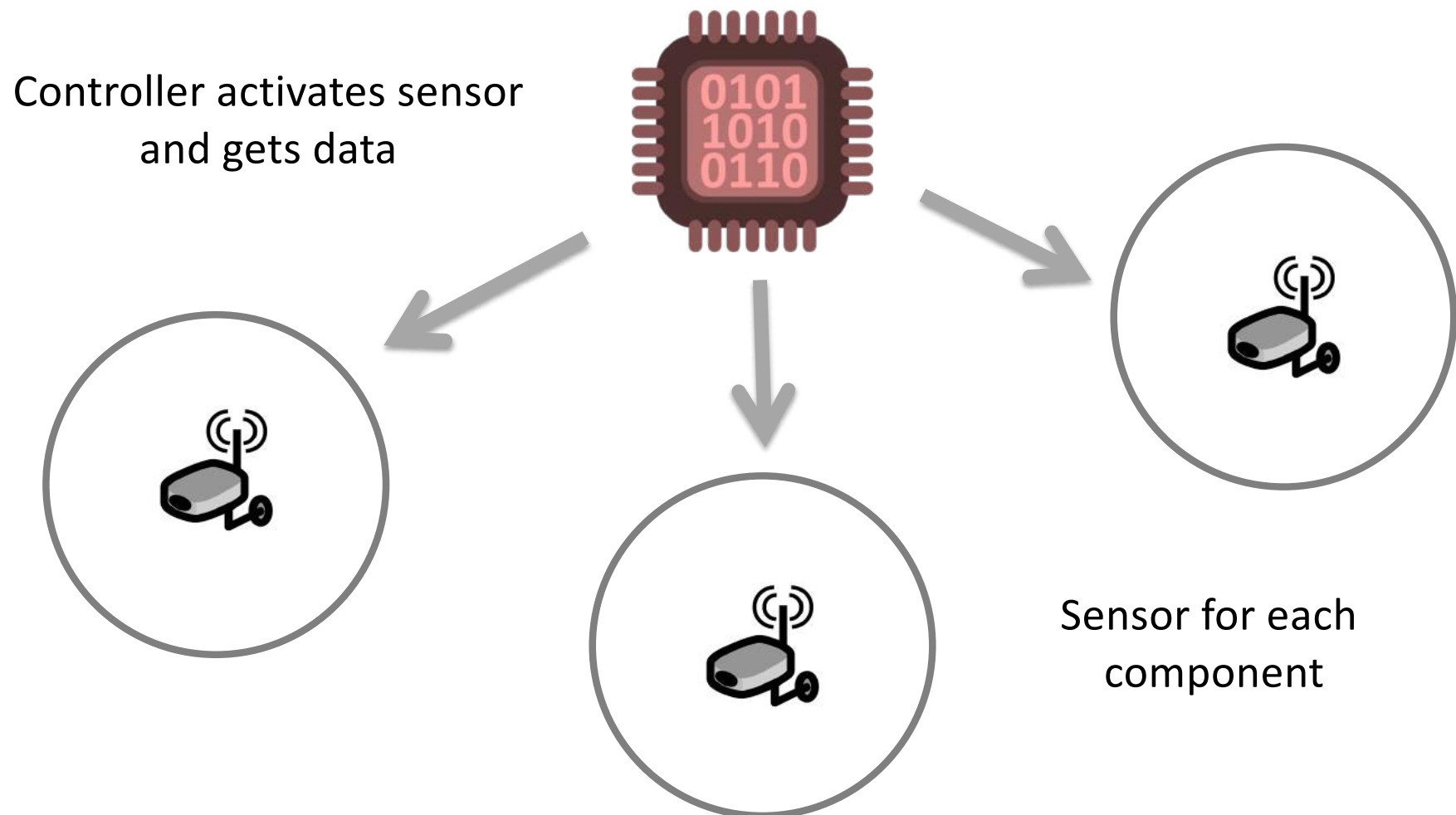
# (TIGHT) FINITE HORIZON?

# Testing for Anomalies

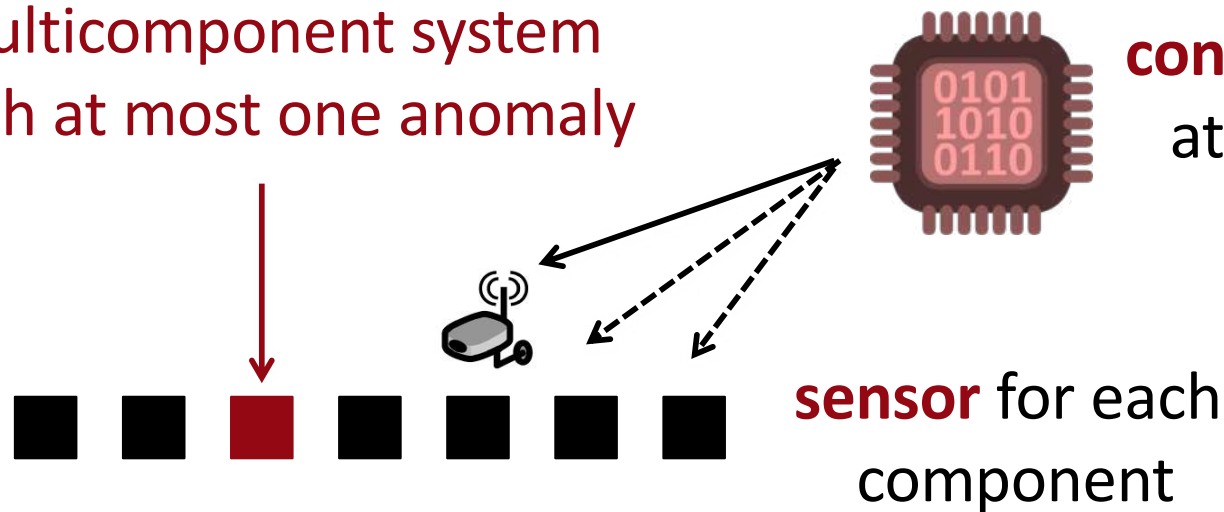## Multicomponent system with potential anomalies

**Controller** activates sensor and gets data

**Anomalous Component**

**Sensor** for each component

# Testing for Anomalies

**Goal: Test whether there is anomaly or not**



Controller activates sensor and gets data

Sensor for each component

# Anomaly Detection – a problem with symmetries

multicomponent system with at most one anomaly

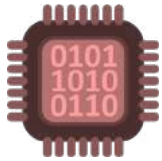**controller** activates sensors at different components at each time slot

**sensor** for each component

Number of components: $M(=7)$

True system state: $X(=3)$

$$X \in \{0, 1, \ldots, M\}$$

$$X = \begin{cases} 0 & \text{if no anomaly} \\ j & \text{if component } j \text{ anomalous} \end{cases}$$
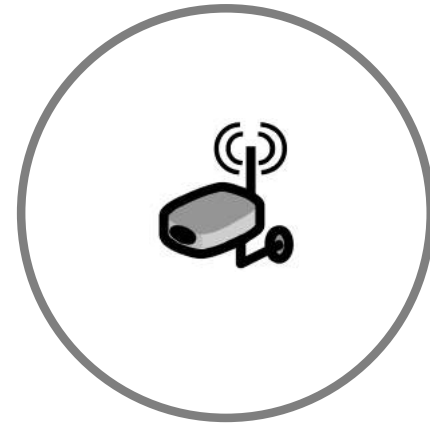
# System Model

**Component** $u$

**Observation** $y$

**Conditional Density**

$$p_1^u(y) \text{ if } X = u$$

Anomalous

$$p_0^u(y) \text{ if } X \neq u$$

Not Anomalous

**Symmetric** if density does not depend on $u$ $\qquad p_i^u(y) \; = \; p_i(y) \; \forall \, u$

# System Model

**Component** $u$
**Observation** $y$

$U_1 = 1, Y_1 \qquad U_2 = 5, Y_2 \qquad U_3 = 2, Y_3 \qquad U_4 = 3, Y_4$

$p_0(y) \qquad\qquad p_1(y)$ conditional density
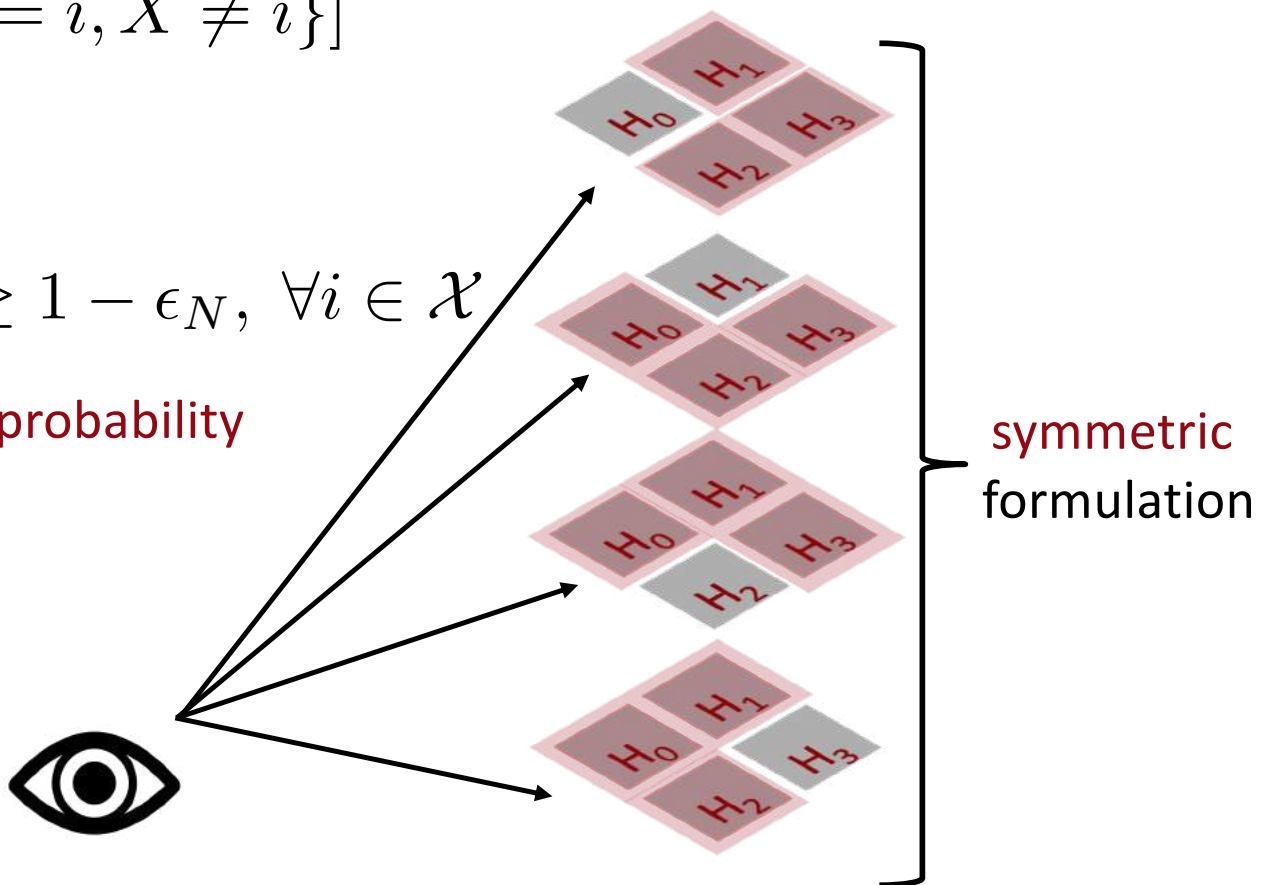
**Symmetric** if density does not depend on $u$

misclassification probability

$$\gamma_N = \mathbb{P}^{f,g}\left[\cup_{i \in \mathcal{X}}\{\hat{X}_N = i, X \neq i\}\right]$$

$$\min_{f \in \mathcal{F}, g \in \mathcal{G}} \gamma_N$$

subject to $\quad \psi_N(i) \geq 1 - \epsilon_N, \ \forall i \in \mathcal{X}$

correct inference probability



symmetric formulation
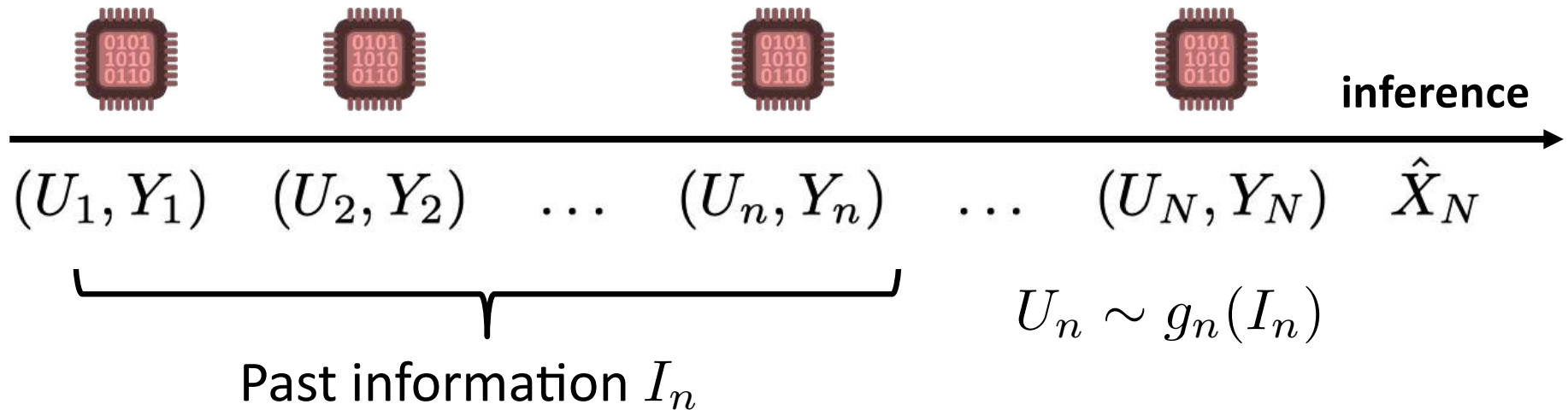
these all look the same!

# Same framework as before

❑ Experiment Selection Strategy:



$$(U_1, Y_1) \quad (U_2, Y_2) \quad \ldots \quad (U_n, Y_n) \quad \ldots \quad (U_N, Y_N) \quad \hat{X}_N$$

inference

Past information $I_n$

$$U_n \sim g_n(I_n)$$

Observation $Y_n$ independent of past given $U_n$ and $X$

❑ Inference Strategy: decide safe or not safe

binary valued inference

$$\hat{X}_N \sim f(I_{N+1})$$

also randomized

safe: $X = 0$

not safe: $X \neq 0$

# Contributions

❑ Pose fixed-horizon active Neyman-Pearson anomaly detector

- asymptotically optimal error rates

- For a symmetric system, even stronger non-asymptotic converse bounds

❑ Design deterministic experiment selection strategies

- Achieve asymptotic bounds

- Up to an additive logarithmic term (strong sense) in non-asymptotic regime → 2nd order optimal

❑ *Open loop strategies (asymptotically optimal) not strong in finite case*

# Neyman-Pearson Formulation

$$\psi_N \doteq \mathbb{P}^{f,g}[\hat{X}_N = 0 \mid X = 0]$$

correct detection probability

$$\phi_N \doteq \mathbb{P}^{f,g}[\hat{X}_N = 0 \mid X \neq 0]$$

incorrect detection probability

**Problem (P)**

$$\inf_{f \in \mathcal{F}, g \in \mathcal{G}} \phi_N$$
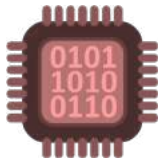
subject to $\quad \psi_N \geq 1 - \epsilon_N$

Infimum value: $\phi_N^*$

**minimize error subject to correct detection constraint**

❑ Incorrect safe declaration very expensive – can tolerate a few false alarms

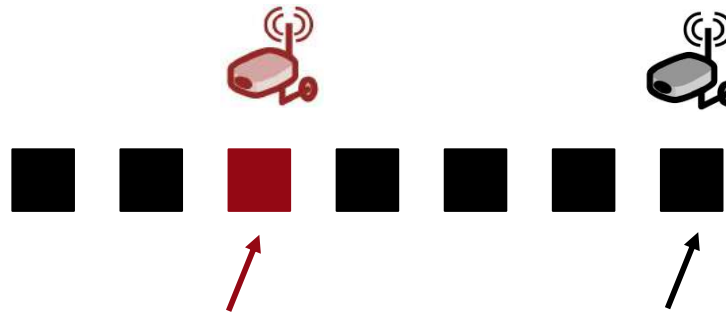❑ **GOAL**: Find detection/inference & experiment selection strategies to solve (P)

# Log-likelihood Ratios

**Component** $u$
**Observation** $y$

**Conditional Density**

$p_1^u(y)$ if $X = u$
Anomalous

$p_0^u(y)$ if $X \neq u$
Not Anomalous

$$L_j(u, y) \doteq \begin{cases} \log \frac{p_0^u(y)}{p_1^u(y)} & \text{if } u = j \\ 0 & \text{otherwise.} \end{cases}$$

$$\boxed{D_j^u = \mathbb{E}[L_j(u, Y)]}$$
$$Y \sim p_0^u$$

log-likelihood ratios    $X = 0$ vs $X = j$

Kullback-Leibler Divergences

# Accumulated LLR and Confidence Level

❑ Accumulate log-likelihood ratios for each component

$$Z_n(j) \doteq \sum_{k=1}^{n} L_j(U_k, Y_k)$$

❑ Confidence level: is a log-likelihood ratio

$$\mathcal{C}(I_{n+1}, \rho_1) = -\log\left[\sum_{j \in \mathcal{U}} \exp\left(\log \tilde{\rho}_1(j) - Z_n(j)\right)\right]$$
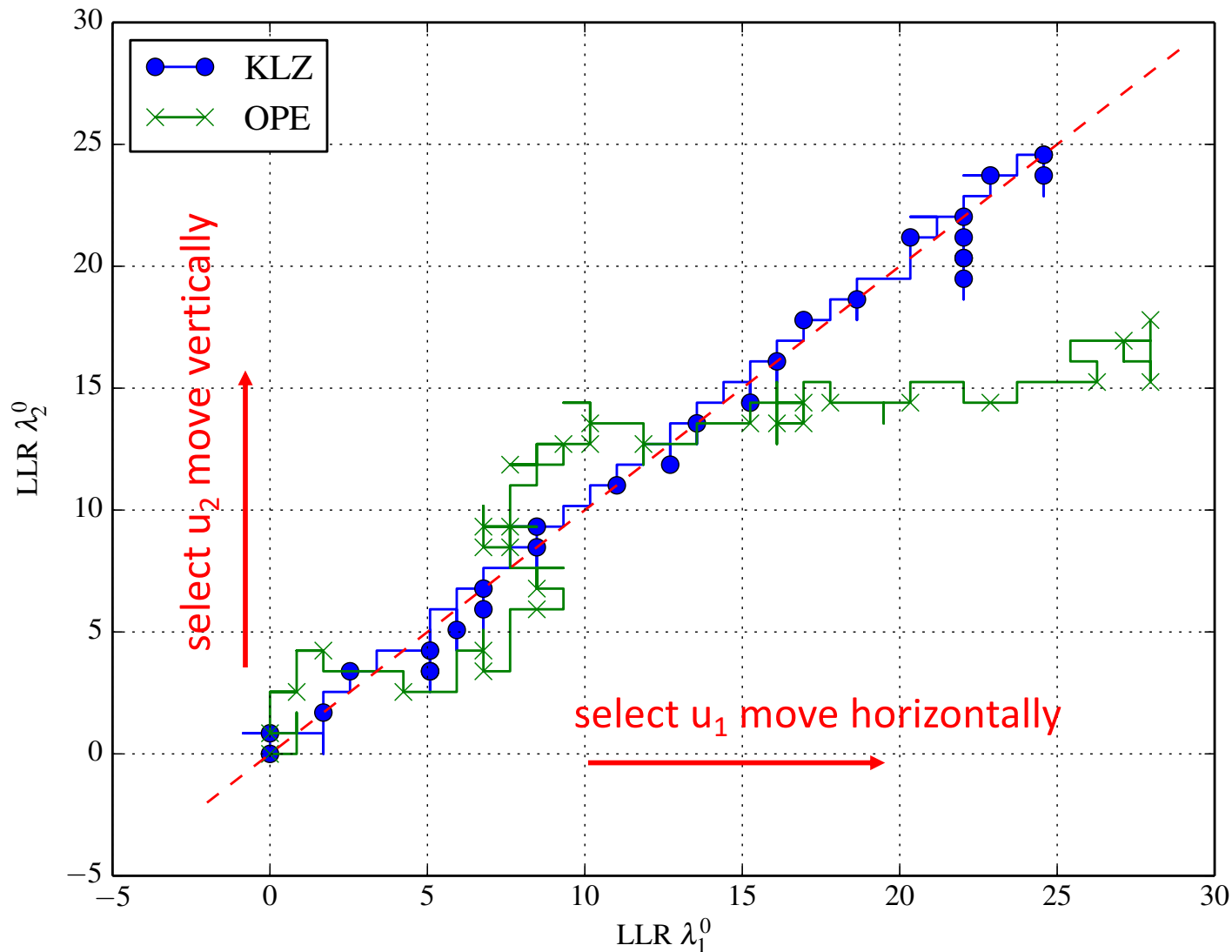
prior belief

$$\approx \min_{j \in \mathcal{U}}\{Z_n(j)\}$$

$$\tilde{\rho}_1(j) = \rho_1(j)/(1 - \rho_1(0))$$

# Accumulated LLR Evolution

evolution of LLRs under different experiment selection strategies



study the evolution of accumulated LLR vector

analysis easier for random walks difficult otherwise

# Interpreting the plot

- ❏ Kullback-Leibler Divergence:

$$D(p\|q) = \sum_y p(y) \log \frac{p(y)}{q(y)}$$

$$\mathbb{E}_0[L_n] = nD(p_0\|p_1)$$
$$\mathbb{E}_1[L_n] = -nD(p_1\|p_0)$$

Expectation of LLR is related to KL-Divergence

- ❏ Random walk

$$L_n \rightarrow nD(p_0\|p_1) \text{ under } H_0$$
$$L_n \rightarrow -nD(p_1\|p_0) \text{ under } H_1$$

# Recall Max-min KL-Divergence

- Define $\alpha, \beta$ distributions

$$D^* \doteq \max_{\alpha \in \Delta \mathcal{U}} \min_{j \in \mathcal{U}} \sum_{u \in \mathcal{U}} \alpha(u) D_j^u \qquad \text{argmax: } \alpha^*$$

$$= \min_{\beta \in \Delta \mathcal{U}} \max_{u \in \mathcal{U}} \sum_{j \in \mathcal{U}} \beta(j) D_j^u \qquad \text{argmin: } \beta^*$$

- Lemma: for anomaly detection/symmetric case

$$D^* = \left( \sum_{u \in \mathcal{U}} \frac{1}{D_u^u} \right)^{-1}$$

$$\alpha^*(u) = \beta^*(u) = D_u^u / D^*$$

recall $D_u^u \neq 0$ when anomaly

*uniform when symmetric*

# Asymptotic Results

❑ Weak converse: Based on Data Processing Inequality

$$-\frac{1}{N}\log \phi_N^* \leq \frac{D^*}{1-\epsilon_N} + \frac{O(1)}{N(1-\epsilon_N)}$$

<span style="color:#8B0000">error probability</span>

$$\psi_N \geq 1 - \epsilon_N$$

❑ Previous converse for the general case:

$$-\frac{1}{N}\log \gamma_N = -\frac{1}{N}\log \left( \sum_i \mathbb{P}[X \neq i]\phi_N(i) \right)$$

$$\leq \min_i D^*(i) + \Theta(1/N)$$

# Asymptotic Results

❑ Asymptotic achievability:

▪ Experiment selection strategy: randomly select component from distribution $\alpha^*$ (Open loop sufficient!)

▪ Inference strategy: decide safe only if confidence sufficiently large

$$\mathcal{C}(I_{n+1}, \rho_1) = -\log \left[ \sum_{j \in \mathcal{U}} \exp \left( \log \tilde{\rho}_1(j) - Z_n(j) \right) \right]$$

❑ Strategy essentially the same, but can decompose confidence function better due to symmetry of distributions

# Asymptotic Results

- Optimal error rate: under some minor assumptions

$$\lim_{N \to \infty} -\frac{1}{N} \log \phi_N^* = D^*$$

Generalization of Chernoff-Stein Lemma

$$D^* = \left( \sum_{u \in \mathcal{U}} \frac{1}{D_u^u} \right)^{-1}$$

$$L_j(u, y) \doteq \begin{cases} \log \frac{p_0^u(y)}{p_1^u(y)} & \text{if } u = j \\ 0 & \text{otherwise.} \end{cases}$$

$$D_j^u = \mathbb{E}[L_j(u, Y)]$$
$$Y \sim p_0^u$$

# NON-ASYMPTOTIC RESULTS

# Martingales

❑ Definition

$\{M_n\}_{n=0}^{\infty}$ is a Martingale wrt $\{X_n\}_{n=0}^{\infty}$ if $\forall\ n \geq 0$

1. $M_n = f(X_0, \cdots X_n)$

2. $\mathbb{E}\left[|M_n|\right] < \infty$

3. $\mathbb{E}\left[M_{n+1} \mid M_n, \cdots M_0\right] = M_n$  almost surely

   ▪ $\{X_n\}_{n=0}^{\infty}$ need not be specified, only items 2. and 3.

# Why Martingales?

❑ Prove bounds/convergence

  ▪ Estimation and control

❑ Can generalize LLN and CLT

  ▪ Sums of random variables

❑ Martingale difference sequences

  ▪ Exploited in prediction/control

❑ Foster-Lyapunov drift

  ▪ Explore the stability of Markov processes

❑ *Martingale theory allows for a lack of Markovity and linearity*

# Example

191

$\{X_n\}$ iid with $M_n = \sum_{k=0}^{n} X_k$ such that

$$\mathbb{E}[X_0] = 0$$

$$\mathbb{E}[|X_k|] < \infty$$

1. $M_n = f(X_0, \cdots X_n)$

2.

$$\mathbb{E}[|M_n|] = \left[\left|\sum_k X_k\right|\right]$$

$$< \left[\sum_k |X_k|\right] < \infty$$

# Martingale property

$$
\begin{aligned}
\mathbb{E}\left[M_{n+1}\middle|M_n,\cdots M_0\right] &= \mathbb{E}\left[\sum_{k=0}^{n+1} X_k \middle| X_n,\cdots X_0\right] \\
&= \sum_{k=0}^{n+1} \mathbb{E}\left[X_k\middle| X_n,\cdots X_0\right] \\
&= \mathbb{E}\left[X_{n+1}\middle| X_n,\cdots X_0\right] + \sum_{k=0}^{n} \mathbb{E}\left[X_k\middle| X_n,\cdots X_0\right] \\
&= \mathbb{E}\left[X_{n+1}\right] + \sum_{k=0}^{n} X_k \\
&= 0 + M_n = M_n
\end{aligned}
$$

# Concentration Inequalities

❑ **Azuma-Hoeffding** inequality (1963/1967)

$\{M_n\}$ is Martingale, if $\exists \{\delta_i\} \in \mathbb{R}$ such that

$$\mathbb{P}\left[|M_n - M_{n-1}| \leq \delta_i\right] = 1 \quad \forall \ n$$

then

$$\mathbb{P}\left[|M_n - M_0| \geq C\right] \leq 2\exp\left(-\frac{C^2}{2\sum \delta_i^2}\right) \quad C > 0$$

▪ If increments bounded, probability of a large deviation is small

▪ Samples *concentrate* about a point as *n* gets large

# Proof Ingredients

❑ Proof of AH

- Chernoff bound/Markov inequality

- Convexity/Jensen's inequality

- Martingale property

- Minimize over Chernoff variable

❑ AH versus us…

- General Martingales, bounded increments

- We will exploit conditional independence, but possibly unbounded increments

- BIG PICTURE, very similar

# Key Decomposition Lemma

$$\tilde{\rho}_1(j) = \rho_1(j)/(1 - \rho_1(0)) \qquad \tilde{\rho}_{n+1}(j) = \frac{\tilde{\rho}_1(j)e^{-Z_n(j)}}{\sum_{k \in \mathcal{U}} \tilde{\rho}_1(k)e^{-Z_n(k)}}$$

$$\mathcal{C}(I_{n+1}, \rho_1) = \left[ \bar{Z}_n + D(\beta^* \| \tilde{\rho}_1) \right] + \left[ -D(\beta^* \| \tilde{\rho}_{n+1}) \right]$$

arg min max

$$\bar{Z}_n \doteq \sum_{j \in \mathcal{U}} \beta^*(j) Z_n(j)$$

**sub-martingale** in general

symmetric case: **i.i.d. sum** and **strategy independent**

# Key Decomposition Lemma

$$\tilde{\rho}_1(j) = \rho_1(j)/(1 - \rho_1(0)) \qquad \tilde{\rho}_{n+1}(j) = \frac{\tilde{\rho}_1(j)e^{-Z_n(j)}}{\sum_{k \in \mathcal{U}} \tilde{\rho}_1(k)e^{-Z_n(k)}}$$

$$\mathcal{C}(I_{n+1}, \rho_1) = \left[\bar{Z}_n + D(\beta^* || \tilde{\rho}_1)\right] + \left[-D(\beta^* || \tilde{\rho}_{n+1})\right]$$

non-positive

$$\bar{Z}_n \doteq \sum_{j \in \mathcal{U}} \beta^*(j) Z_n(j)$$

**sub-martingale** in general

symmetric case: **i.i.d. sum** and **strategy independent**

# Non-asymptotic Bounds - Symmetric

❑ **Theorem**

Strong converse: follows from decomposition and
strong converse in Polyanskiy, Poor and Verdu, IT Transactions 2010

$$-\log \phi_N^* \leq \mathsf{INV}_N\left(\epsilon_N + \frac{\epsilon_N}{\eta}\right) + \log \frac{\eta}{\epsilon_N}$$

$$-\log \phi_N^* \geq \mathsf{INV}_N\left(\epsilon_N - \frac{\epsilon_N}{\eta}\right) - O\left(\log \frac{\eta}{\epsilon_N}\right)$$

Strong achievability: based on decomposition, an adaptive experiment
selection strategy and a Chernoff bound

$\eta > 0$: may depend on $N$

$\mathsf{INV}_N$: quantile function of $\bar{Z}_N + D(\beta^*||\tilde{\rho}_1)$

# Berry-Esseen Theorem

❑ Consider the empirical mean of i.i.d. variables

$$Y_n = \frac{X_1 + \cdots + X_n}{n}$$

$$\mathbb{E}[X_1] = 0$$
$$\mathbb{E}[X_1^2] = \sigma^2$$
$$\mathbb{E}[|X_1|^3] = \rho$$

❑ Then

$$|F_n(x) - \Phi(x)| \leq \frac{C\rho}{\sigma^3 \sqrt{n}}$$

CDF of $\dfrac{Y_n \sqrt{n}}{\sigma}$

CDF of standard normal

# Berry-Esseen Approximation

❑ Corollary: straightforward application of the Berry-Esseen theorem (approximate everything as Gaussian from CLT)

$$-\log \phi_N^* \leq ND^* - \sqrt{NV}Q^{-1}\left(\epsilon_N + \frac{\epsilon_N}{\eta} + \frac{6T}{\sqrt{NV^3}}\right) + O\left(\log \frac{\eta}{\epsilon_N}\right)$$

$$-\log \phi_N^* \geq ND^* - \sqrt{NV}Q^{-1}\left(\epsilon_N - \frac{\epsilon_N}{\eta} - \frac{6T}{\sqrt{NV^3}}\right) - O\left(\log \frac{\eta}{\epsilon_N}\right)$$

$V$ : variance of LLR

$T$ : centered absolute third moment of LLR

$Q$ : tail distribution of standard normal

# Two Experiment Selection Strategies

- ❏ Open-loop randomized: asymptotically optimal

  randomly select component from distribution $\alpha^*$

- ❏ Adaptive deterministic: also asymptotically optimal
  at each time $n$, select the component $j$
  that minimizes $Z_{n-1}(j) - \log \tilde{\rho}_1(j)$

  confidence $\quad C(I_{n+1}, \rho_1) = -\log \left[ \sum_{j \in \mathcal{U}} \exp \left( \log \tilde{\rho}_1(j) - Z_n(j) \right) \right]$

- ❏ Example setting: two-component and binary observations

# Individual Sampling Results

Open-loop strategies:
1. Uniform random selection
2. Round robin
3. Open-loop sampling cannot get better than this

Adaptive Selection

Note that computed strong lower bounds are fairly tight

128 component system with Gaussian likelihoods and individual sampling
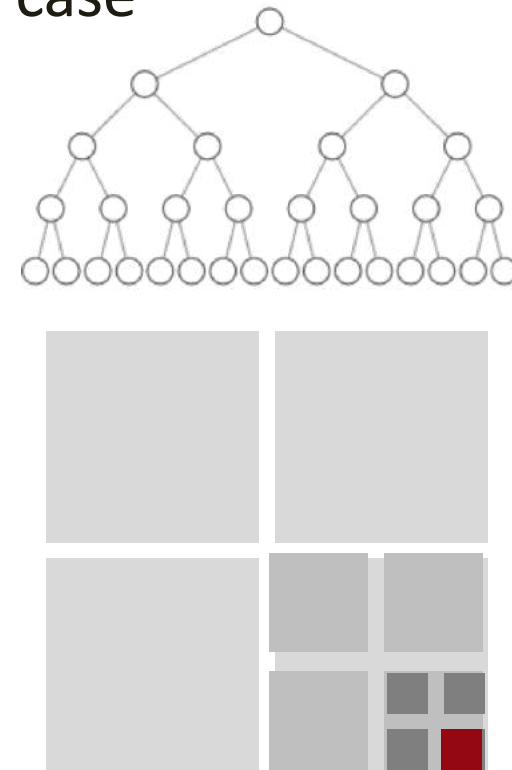
# Exploration Phase

❑ Exploration important for symmetric case

  ▪ Search for anomaly based using grouped observations

Classical approaches suggest lawnmower-type exhaustive search
Chernoff, 1959; Nitinawarat, Atia, Veeravalli 2013

Binary-search type approaches more efficient
Naghshvar, Javidi 2012, 2013; Chiu, Javidi 2020

# Exploration Time

Exploration time:  $T \doteq \min\{n' : \mathcal{C}_X(\rho_n) \geq 0 \ \forall n \geq n'\}$
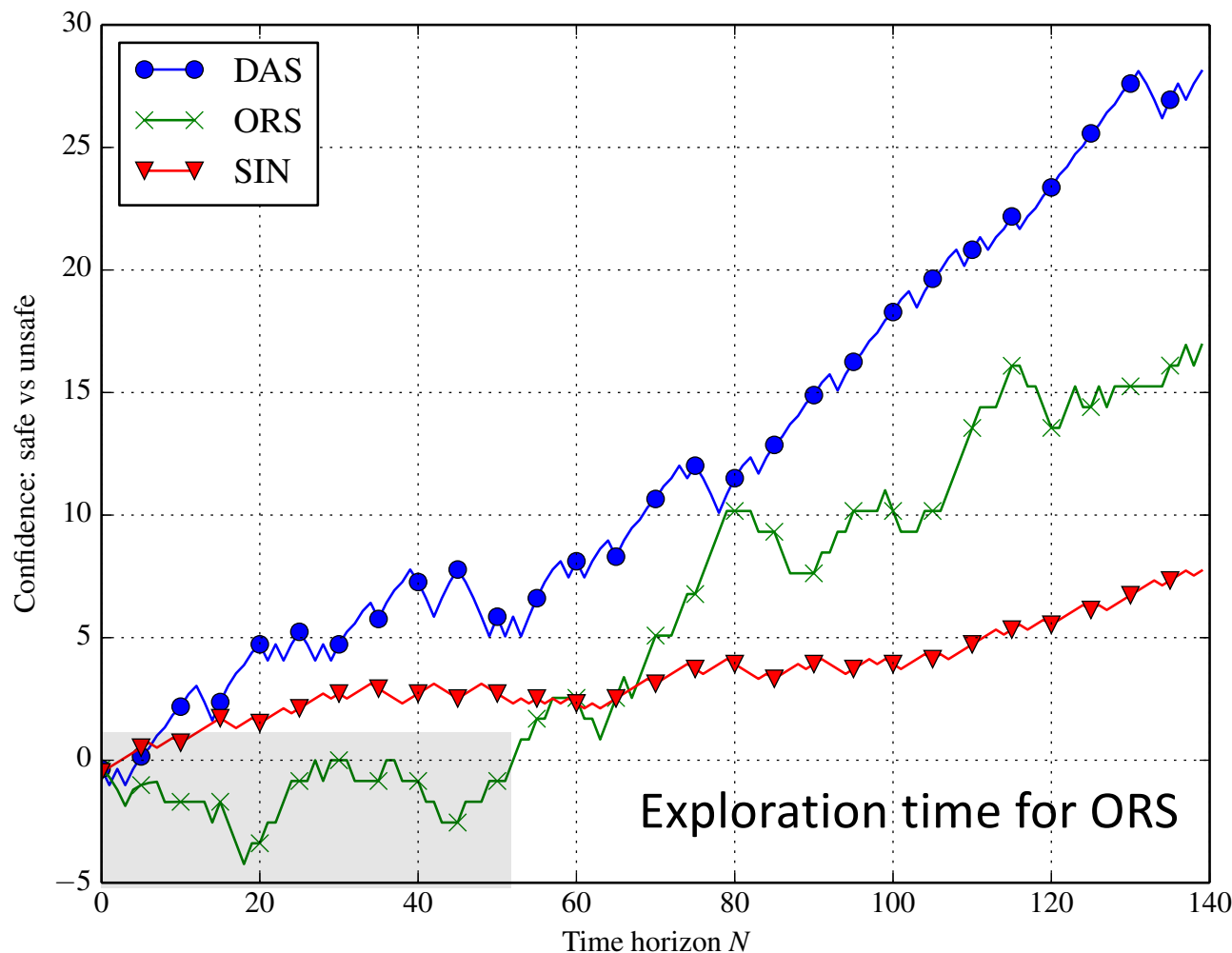


after exploration time, our most likely hypothesis is **always** the true hypothesis

exploration strategy should ensure exploration time is small – we derive high probability upper bounds on this

# Exploration Time

Exploration time: $T \doteq \min\{n' : \mathcal{C}_X(\rho_n) \geq 0 \ \forall n \geq n'\}$



Exploration time for ORS

After exploration time our most likely hypothesis is **always** the true hypothesis
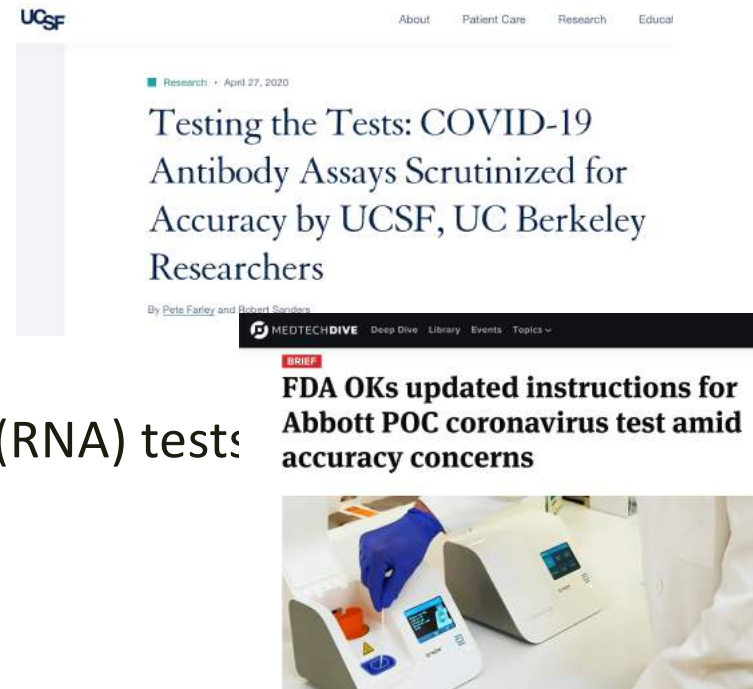
compute exploration time only in hindsight

Exploration strategy should ensure exploration time is small – we derive high probability upper bounds on this

# SARS-CoV-2 Testing

- ❑ A few realities have emerged
  - ▪ Insufficient number of tests
  - ▪ Tests have different efficacies
  - ▪ Timing of test administration matters
    - – Both for serological (antibody) and PCR (RNA) tests

- ❑ The future should enable
  - ▪ Heterogeneous tests
  - ▪ Regular testing

- ❑ How can active methods help?

UCSF
About    Patient Care    Research    Educat

Research · April 27, 2020

**Testing the Tests: COVID-19 Antibody Assays Scrutinized for Accuracy by UCSF, UC Berkeley Researchers**

By Pete Farley and Robert Sanders

MEDTECHDIVE    Deep Dive    Library    Events    Topics

BRIEF

**FDA OKs updated instructions for Abbott POC coronavirus test amid accuracy concerns**

NEWS | CORONAVIRUS (COVID-19) | JUNE 10, 2020

**COVID-19 Genetic PCR Tests Give False Negative Results if Used Too Early**

A new study confirms what many suspected, that PCR testing even 8 days after infection shows 20 percent false positives
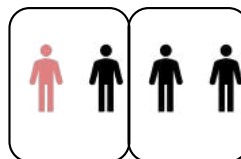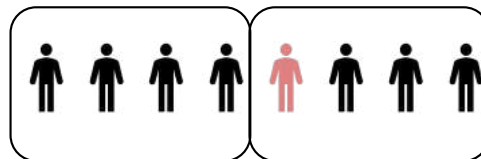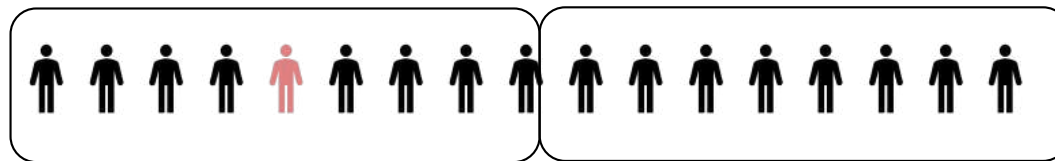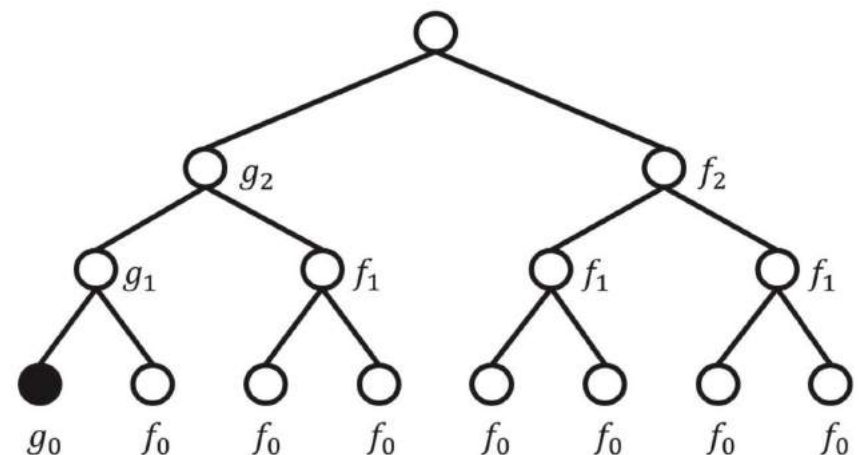
# Recall Group Testing

- Used in WW2 to test soldiers for syphilis
  - R. Dorfman, "The Detection of Defective Members of Large Populations," The Annals of Mathematical Statistics, 1943
  - Binary search
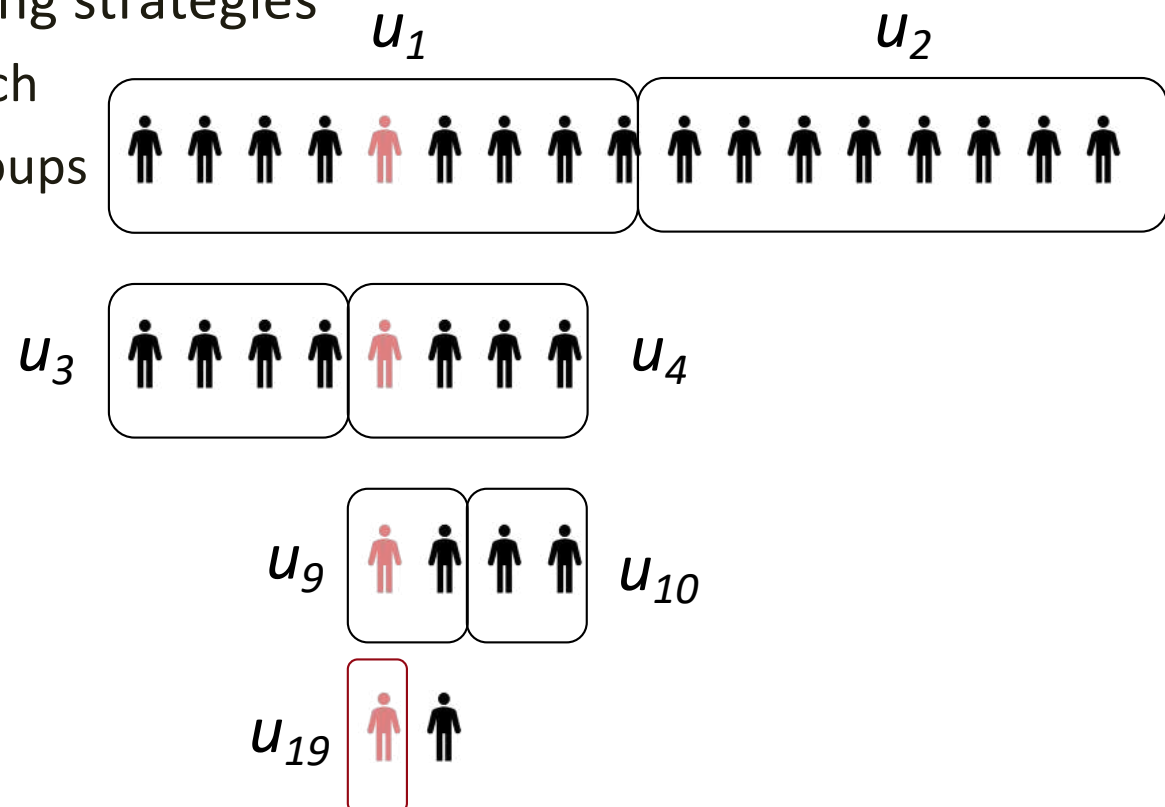
- N tests → log(N) tests
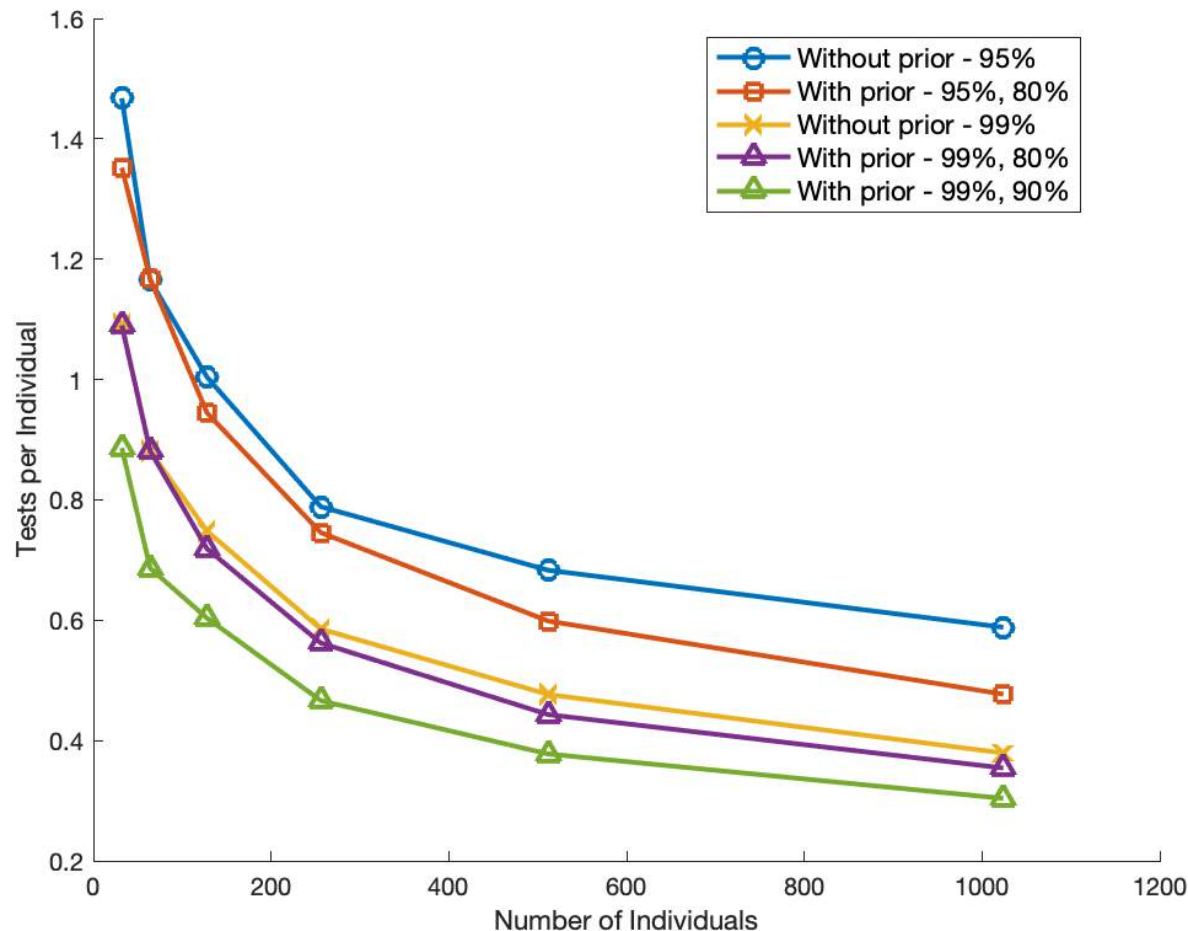
# Mapping to Active Testing

- ❑ A variety of formulations
    - ▪ Form all possible groups, each distinct group is an experiment
        - – Computationally expensive
    - ▪ Pre-select grouping strategies
        - – E.g. Binary search
        - – Time-varying groups
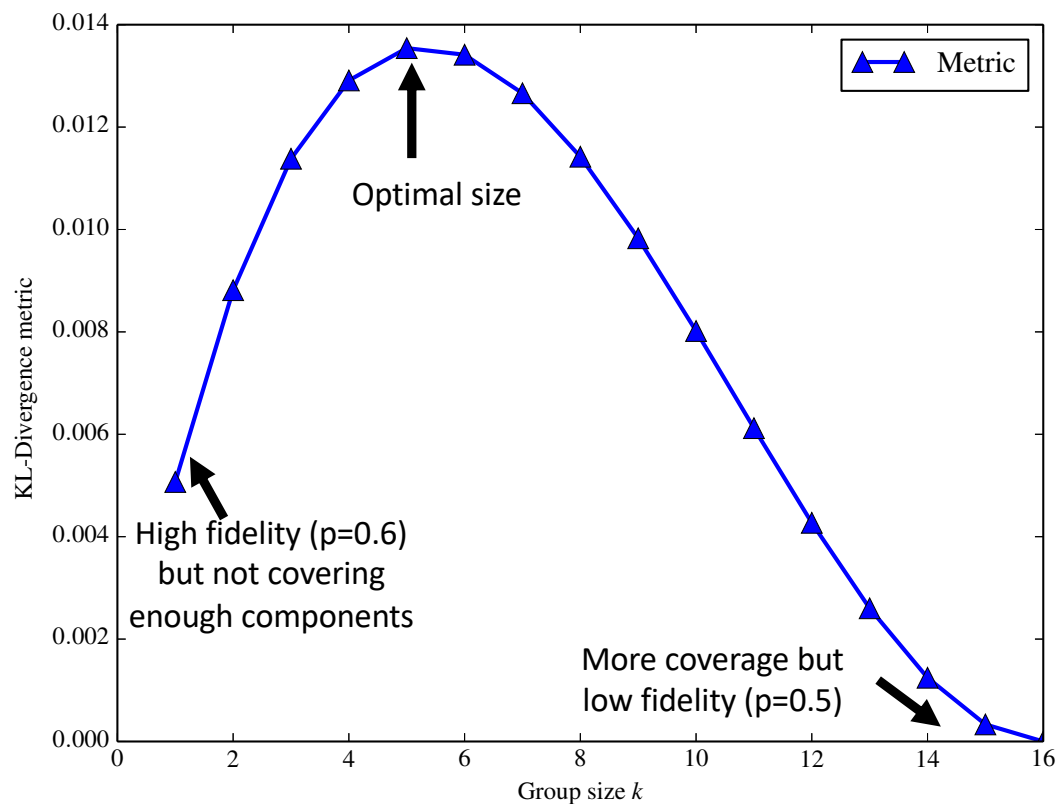
# Fully-adaptive Tests

- Perform a cheap test first on each individual – we consider tests with 80% and 90% accuracy
- Use the prior for group testing subsequently
- Can reduce number of group tests by 20%

- Performing cheap tests first better when the cost of cheap test is about 10-15 times smaller

fully adaptive tests can take a lot of time – need to parallelize

# Group Sampling Results

Selecting optimal group size



Optimal rate

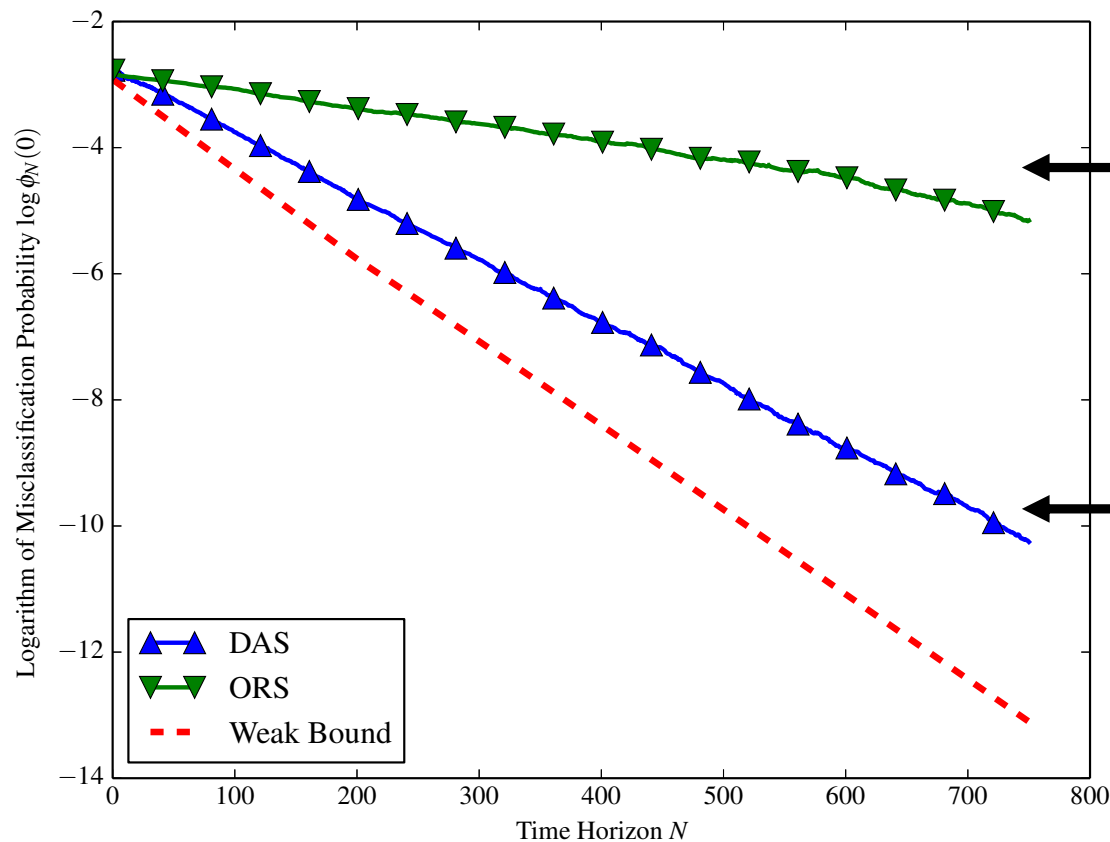$$D^* \doteq \max_{1 \leq k \leq M} \frac{k D_{\{1\}}^{\{1,\ldots,k\}}}{M}$$

$$k^* \doteq \arg \max_{1 \leq k \leq M} \frac{k D_{\{1\}}^{\{1,\ldots,k\}}}{M}$$

Optimal size

A 16-component system with linear dilution: binary symmetric noise goes from 0.6 to 0.5 (indistinguishable)

# Group Sampling Results

Open-loop strategy: randomly select a subset with size $k^*$

Adaptive Selection: select $k^*$ most likely elements

Dramatic performance gap between open-loop and adaptive selection
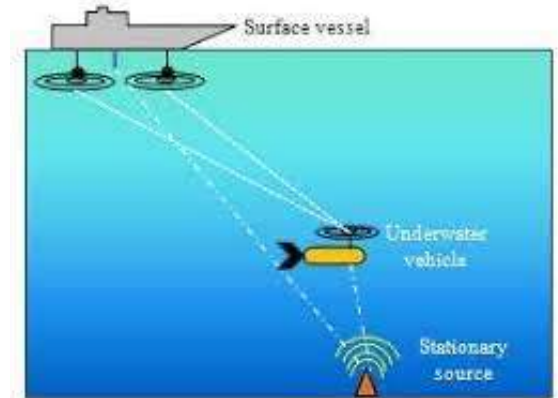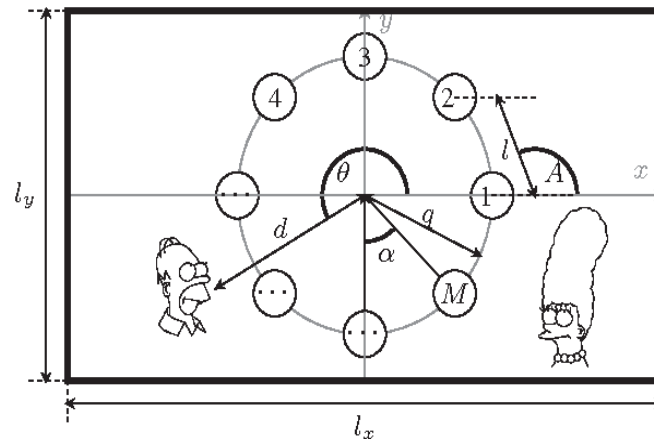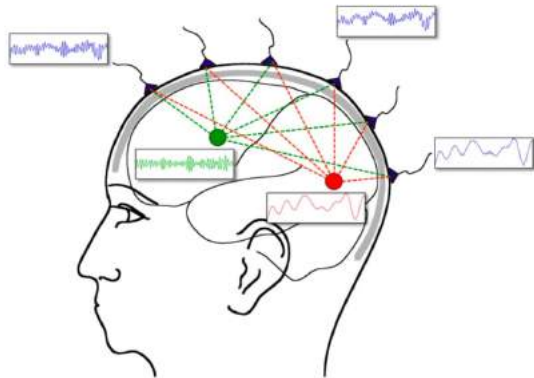
A 16-component system with linear dilution

# ONE LAST APPLICATION

# Source Localization

❑ classical signal processing problem

❑ *Applications:*



❑ Drawbacks of existing works:

- ▪ Parametric methods – model mismatch issues

- ▪ Model parameters  hard to estimate

- ▪ Model-free approaches coarse localization

- ▪ ML-based approaches require lots of training data

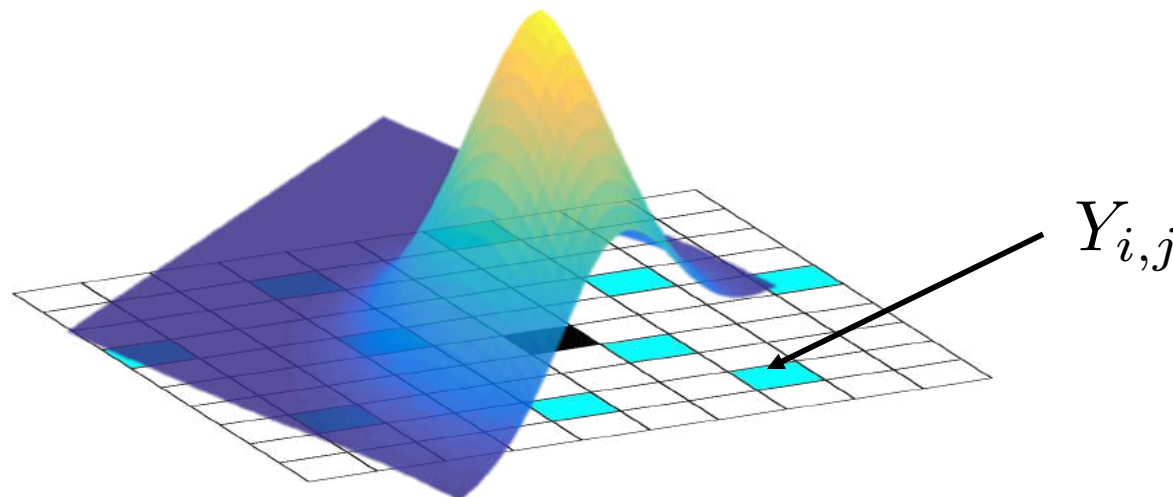# Localization Challenge

- Source location $s^* \in \mathbb{R}^2$ (unknown) ■

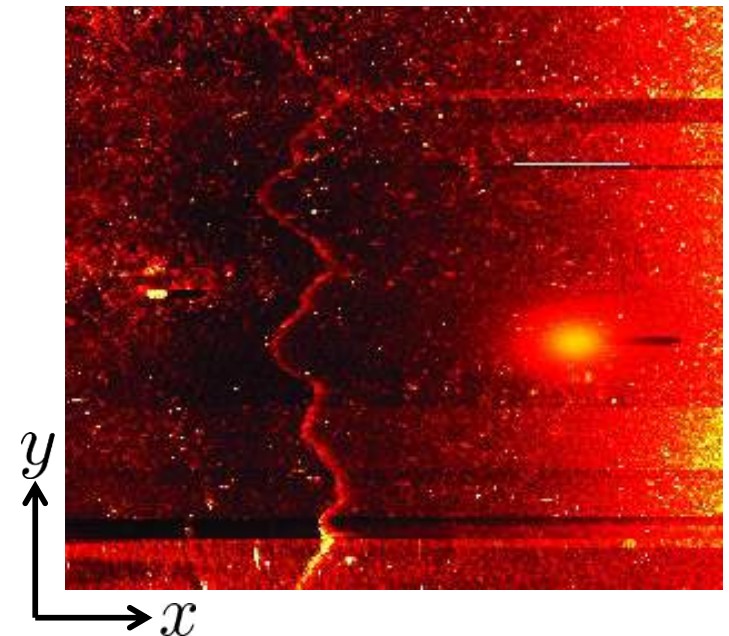$$Y \doteq H(s^*) + Z$$

- If $\mathbf{Y} \in \mathbb{R}^{N \times N}$, $N^2$ hypothesis testing problem
  - Trade-off known distributions for signal structure

- Random samples at locations ■

- Only knowledge about target signal is that it is unimodal



$Y_{i,j}$

# What is a good model?

- ❑ Real sidescan sonar data

- ❑ Any other structural properties to exploit?



Intensity $= \boldsymbol{H}(x, y)$

# Review Singular Value Decomposition

$$\mathbf{X} \in \mathbb{C}^{m \times n}$$

$$\mathbf{X} = \mathbf{U\Sigma V}^H$$

$$\left.\begin{array}{rcl} \mathbf{UU}^H &=& \mathbf{I} \\ \mathbf{VV}^H &=& \mathbf{I} \end{array}\right\} \text{unitary} \quad \Sigma = \text{singular value matrix}$$

$$\text{rank}\left(\mathbf{X}\right) = r$$

$$\Sigma_{i,i} = \sigma_i > 0 \ \ i \leq r$$

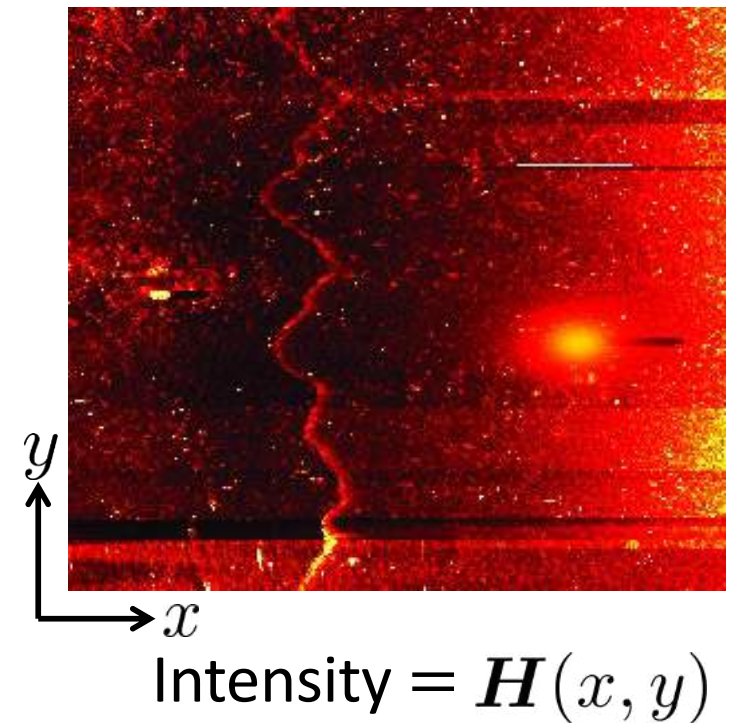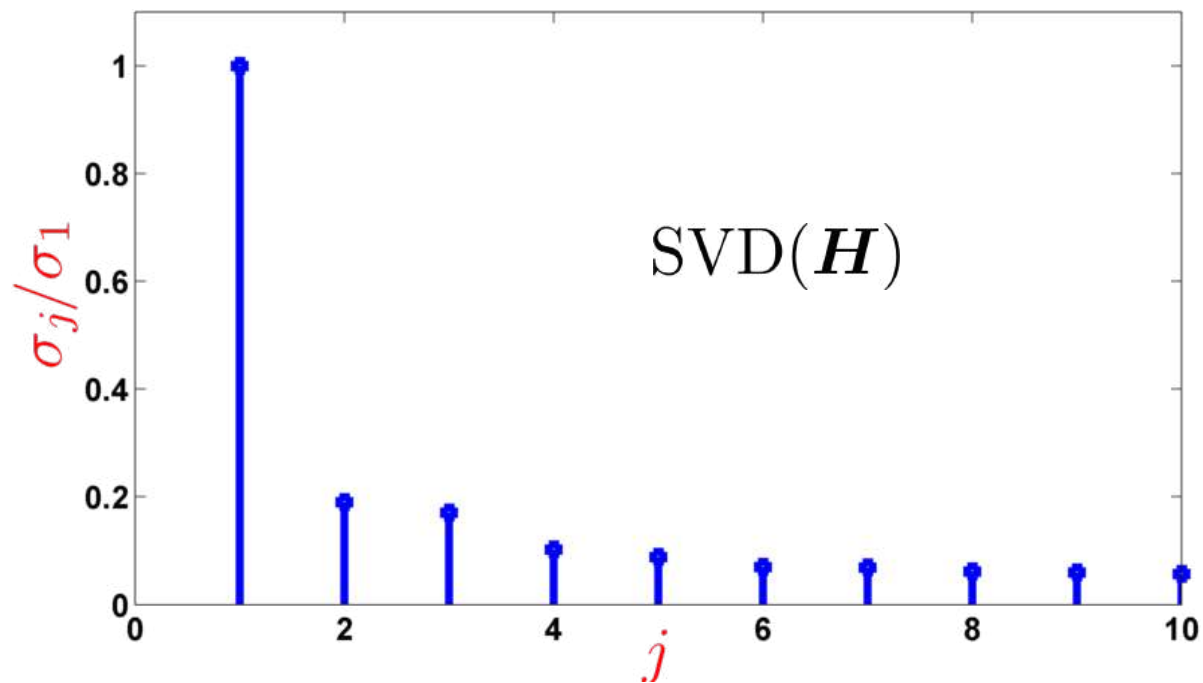$$\mathbf{X}^{m \times n} \ = \ \mathbf{U}^{m \times m} \quad \begin{bmatrix} \sigma_1 & 0 & 0 & 0 & 0 \\ 0 & \sigma_2 & 0 & 0 & 0 \\ 0 & 0 & \sigma_3 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad \mathbf{V}^{n \times n}$$

# Low Rank!

❑ Real sidescan sonar data

❑ Approximate target as **rank one matrix** in image space



$$\text{SVD}(\boldsymbol{H})$$

$$\text{Intensity} = \boldsymbol{H}(x, y)$$

# Low rank approximation

- Largest singular value

$$(\sigma_1, \mathbf{u}_1, \mathbf{v}_1)$$

- Best rank one approximation
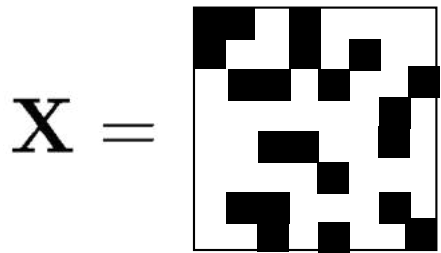
$$\hat{\mathbf{X}} = \arg\min_{\hat{\mathbf{X}}} \|\mathbf{X} - \hat{\mathbf{X}}\|_F$$

$$\text{subject to} \quad \text{rank}\left(\hat{\mathbf{X}}\right) = 1$$

$$= \sigma_1 \mathbf{u}_1 \mathbf{v}_1^H$$

$$\|\cdot\|_F : \text{Frobenius norm}$$

- $\mathbf{u}_1, \mathbf{v}_1$ are also unimodal, if $\mathbf{X}$ unimodal [Chen & M TSP'19]

# Review of Matrix Completion

$\mathbf{X} =$ [matrix with black and white cells]

$\mathbf{X}(i,j)$ known for black cells
unknown for white cells (missing data)

If $\mathbf{X}$ low-rank, we can recover missing data

$$\min_{\mathbf{Z}} \operatorname{rank}(\mathbf{Z})$$
$$\text{for} \quad \mathcal{P}(\mathbf{Z}) = \mathcal{P}(\mathbf{X})$$

NP-hard

aka *nuclear norm*

$$\min_{\mathbf{Z}} \sum_i \sigma_i$$
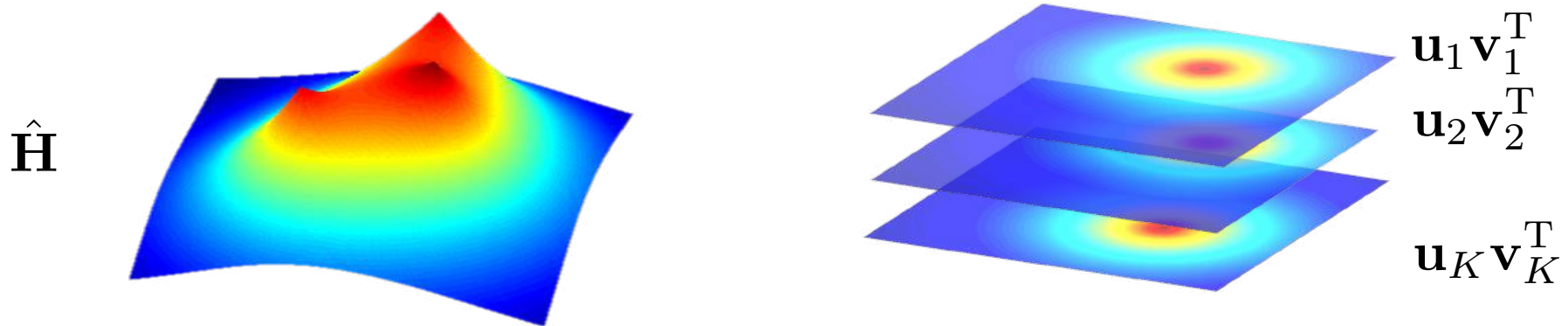$$\text{for} \quad \mathcal{P}(\mathbf{Z}) = \mathcal{P}(\mathbf{X})$$

convex relaxation

$$\sigma_i = \text{singular values of } \mathbf{X}$$

E Candes & B Recht, Found of Computational Math 3/2009
D Gross IEEE Trans on Information Theory 3/2011

# Our Prior Work

- ❑ Multisource localization from random samples

  - ▪ Exploit unimodality of each source signal



$$\underset{\{\alpha_k, \mathbf{u}_k, \mathbf{v}_k\}}{\text{minimize}} \left\| \mathcal{P}_\Omega \left( \hat{\mathbf{H}} - \sum_{k=1}^{K} \alpha_k \mathbf{u}_k \mathbf{v}_k^{\mathrm{T}} \right) \right\|_F^2$$

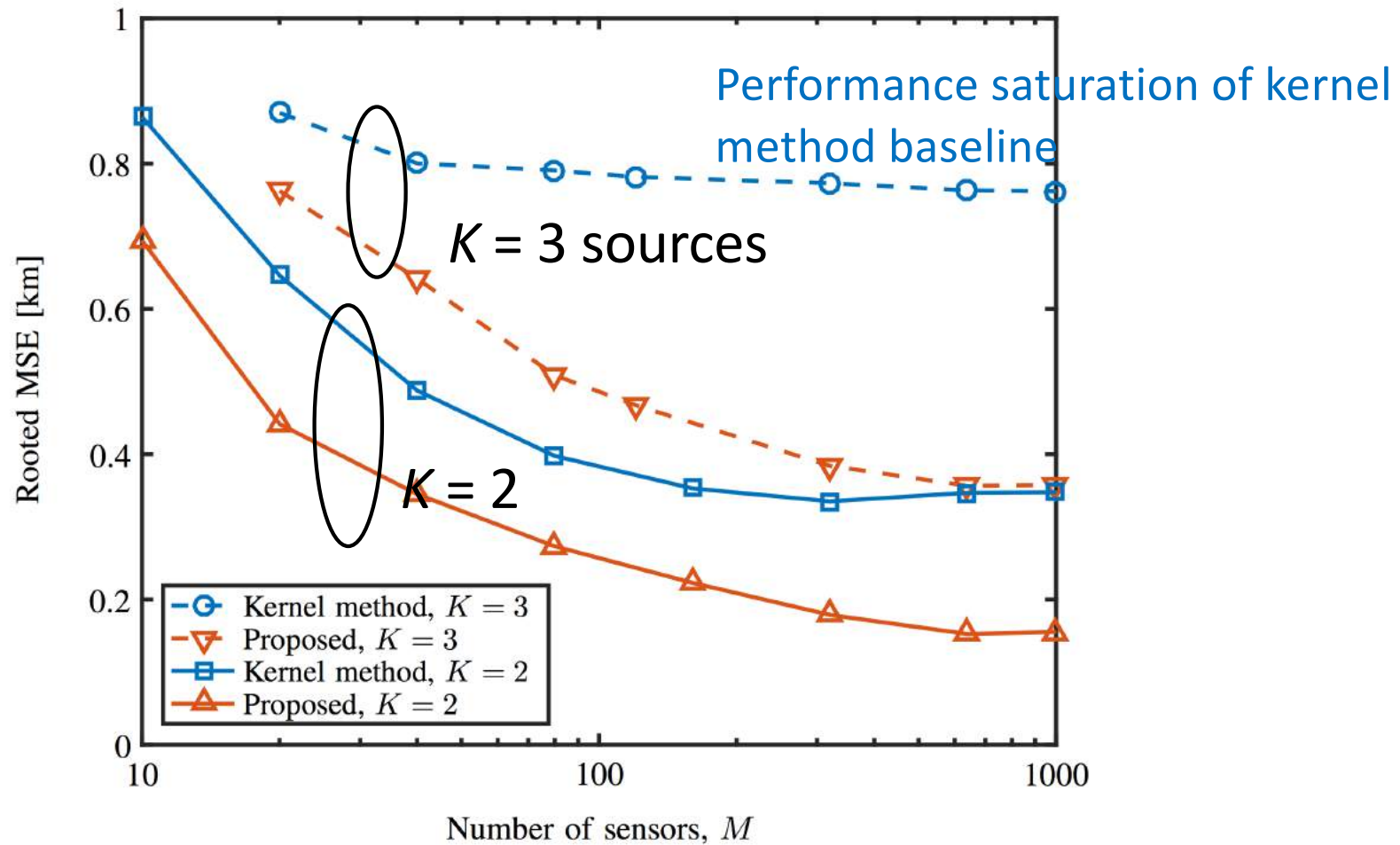$$\text{subject to} \quad \mathbf{u}_k, \mathbf{v}_k \text{ are unimodal}$$

- ❑ Can be solved via projected gradient methods

Chen & **M**, DSP'17, Asilomar'17, ICASP'18, ICASSP'19, TSP'19
Zhang, Chen, Xie, Shapiro &**M**, SPL'21

# Multiple Sources



Performance saturation of kernel method baseline

*K* = 3 sources

*K* = 2

do not need to complete matrix first

# CAN WE MAKE THIS ACTIVE?

# Signal Model

- ❑ Source at location $s^* \in \mathbb{R}^2$ (unknown)

$$Y \doteq H(s^*) + Z$$

- ❑ $H(s^*)$ unimodal
  - ▪ For a single source $H(s^*)$ is rank 1 [Chen & M TSP'19]

- ❑ Definition: Matrix is unimodal with mode at $(i^*, j^*)$ if
  - ▪ $M_{1,j} \leq M_{2,j} \cdots \leq M_{i^*,j} \geq M_{i^*+1,j} \geq \cdots \geq M_{n,j} \quad \forall j$

    $M_{i,1} \leq M_{i,2} \cdots \leq M_{i,j^*} \geq M_{i,j^*+1} \geq \cdots \geq M_{i,n} \quad \forall i$

- ❑ No assumptions except unimodality (non-parametric) $\Longrightarrow$

  convergence + optimal error bounds HARD!

Narayanamurthy & **M**, ISIT'22, Asilomar'22

# Algorithm - Exploration

❏ Initial Exploration: Latin Squares



❏ choose each row, column exactly once, with equal probability

– widely used in experiment design, cryptography, board games

❏ Randomized initialization insufficient

❏ complete rank-1 matrix to get initial row, col estimate

▪ Recall from matrix completion, SVD, $\mathbf{u}_1, \mathbf{v}_1$ are also unimodal if $\mathbf{X}$ unimodal

# Adaptive Sampling - Exploitation

- ❑ Given initialization/exploration, how we do we exploit?
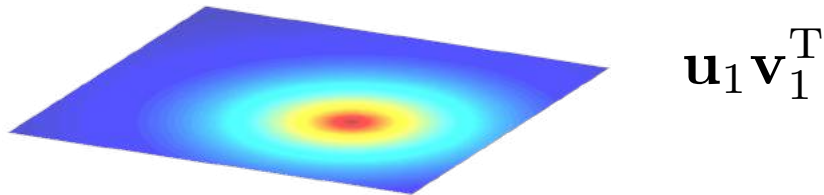    - ▪ Uncertainty-Based approach: query max entropy location

- ❑ **Theorem:** (Uncertainty Quantification for MC, Chen et al. '21)
    - ▪ Consider a rank r matrix $\boldsymbol{Y} \stackrel{\mathrm{SVD}}{=} \boldsymbol{U}\boldsymbol{\Sigma}_y \boldsymbol{V}^\top$
    - ▪ given $\mathcal{O}(nr^5\mathrm{polylog}(n))$ entries sampled uniformly at random
    - ▪ let $\hat{\boldsymbol{Y}}$ denote output of ANY matrix completion algorithm
    - ▪ With probability at least $1 - n^{-3}$

$$\hat{\boldsymbol{Y}}_{i,j} \sim \mathcal{N}(\boldsymbol{Y}_{i,j}, C\sqrt{r/n}(\|\boldsymbol{U}^{(i)}\|^2 + \|\boldsymbol{V}^{(j)}\|^2))$$

# Decomposing the problem

- ❑ Consider our single source
  - ▪ The two singular vectors are individually unimodal

     $\mathbf{u}_1 \mathbf{v}_1^{\mathrm{T}}$

  - ▪ We can look in "each" direction independently
  - ▪ Recall unimodality definition:

$$M_{1,j} \leq M_{2,j} \cdots \leq M_{i^*,j} \geq M_{i^*+1,j} \geq \cdots \geq M_{n,j}$$



which component is maximum?

# Decomposing the problem

❑ Consider our single source

- The two singular vectors are individually unimodal

$$\mathbf{u}_1 \mathbf{v}_1^{\mathrm{T}}$$

- We can look in "each" direction independently
- Recall unimodality definition:

$$M_{1,j} \leq M_{2,j} \cdots \leq M_{i^*,j} \geq M_{i^*+1,j} \geq \cdots \geq M_{n,j}$$
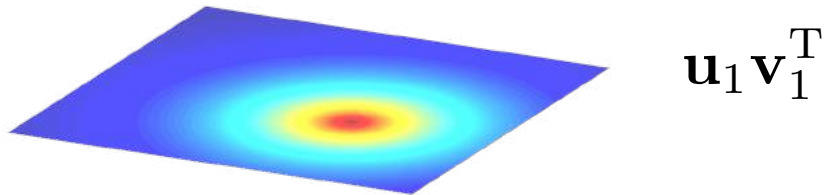
which component is maximum?

# Stochastic Multi Armed Bandits I
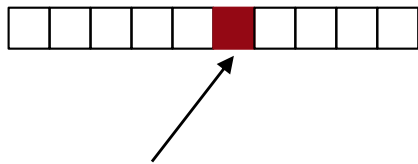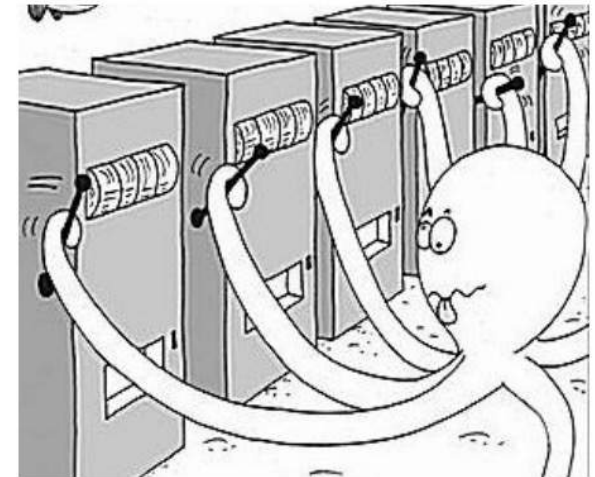
- For each $t$, agent chooses one of $K$ arms and plays it
- The $i\text{-th}$ arm produces reward $r_{i,t} \sim \mathcal{P}_i$ with mean $\mu_i$ (unknown)



source: Microsoft Research

- Agent's objective: maximize cumulative rewards
  - or, find $i^* \doteq \arg\max_i \mu_i$

- Several variants studied based on differing $\mathcal{P}_i$

# Stochastic Multi Armed Bandits II

❑ Example: Stochastic Bernoulli Bandit -- $\mathcal{P}_i$ are Bernoulli

- Let $r_{i,t} \in \{0,1\}$ and $\mathbb{E}[r_{i,t}] = \mu_i$

- If $\mu_i$ were known, optimal policy is to play fixed action
  $$i^* \doteq \arg\max_i \mu_i$$

- If unknown, need to do something better

❑ Regret: $R_n \doteq n \max_i \mu_i - \mathbb{E}[\sum_{t=1}^{n} r_{i,t}]$

- Q: how does $R_n$ scale with $n$ ?

- A: a "good learner" attains sub-linear regret, i.e., $\lim_{n \to \infty} \dfrac{R_n}{n} = 0$

❑ For Bernoulli bandits (our example), $R_n = \Theta(\sqrt{n})$

- [Lattimore and Szepesvari] Bandit Algorithms, '20

# Stochastic Multi Armed Bandits II

- ❑ Example: Stochastic Bernoulli Bandit -- $\mathcal{P}_i$ are Bernoulli
  - Let $r_{i,t} \in \{0, 1\}$ and $\mathbb{E}[r_{i,t}] = \mu_i$
  - If $\mu_i$ were known, optimal policy is to play fixed action
    $$i^* \doteq \arg\max_i \mu_i$$
  - If unknown, need to do something better
- ❑ Regret: $R_n \doteq n \max_i \mu_i - \mathbb{E}[\sum_{t=1}^{n} r_{i,t}]$
  - Q: how does $R_n$ scale with $n$ ?
  - A: a "good learner" attains sub-linear regret, i.e., $\lim_{n\to\infty} \frac{R_n}{n} = 0$
- ❑ For Bernoulli bandits (our example), $R_n = \Theta(\sqrt{n})$
  - [Lattimore and Szepesvari] Bandit Algorithms, '20

# Algorithms: ETC

❑ Explore-then-Commit (ETC):

  ▪ Play each arm a fixed number of times, $m$ (Exploration)

  ▪ After $Km$ rounds, always play "best" arm (Exploitation)

    – Recall that we have $K$ arms

# Algorithms: UCB

❑ **Upper Confidence Bound** (UCB): optimism in the face of uncertainty

- UCB of arm $i$ , in round $t$ is

$$\mathrm{UCB}_i(t-1, \delta) = \hat{\mu}_i(t-1) + \sqrt{\frac{2 \log(1/\delta)}{T_i(t-1)}}$$

- $\delta$ confidence parameter – controls exploration vs exploitation tradeoff
- $T_i(t-1)$ number of times arm $i$ has been played till round $t$
  - If arm has been tried many times, second term will be small (less uncertainty)
- $\hat{\mu}_i(t-1)$ empirical reward of arm $i$ at round $t$ (averaging)
- In each round, pick the arm with largest UCB

- $\delta$ large ➡ a lot of initial exploration (limited optimism)

# UCB intuition I

- Consider 2-arm bandit problem with $\mu_1 = 0, \mu_2 = -0.5$

- Initially, variance ⬆
  "confidence" ⬇

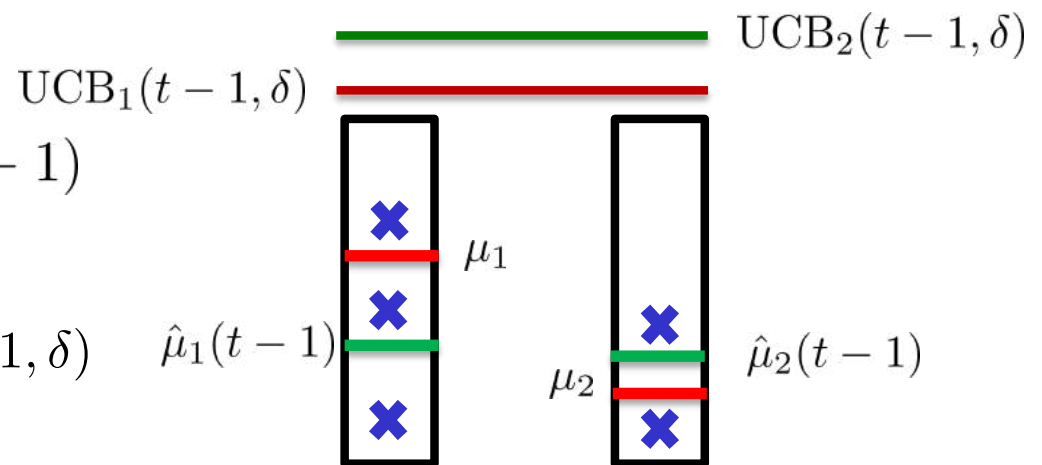- although $\hat{\mu}_1(t-1) \approx \hat{\mu}_2(t-1)$
  arm 2 picked next since

  $$\text{UCB}_2(t-1,\delta) \quad > \quad \text{UCB}_1(t-1,\delta)$$

- hope is that as time progresses,

  $$\text{UCB}_1(t-1,\delta) \gg \text{UCB}_2(t-1,\delta)$$

# UCB intuition II

❑ as time progresses, LLN/CLT says

$$\hat{\mu}_i(t) \to \mu_i$$

❑ CLT also provides "Gaussian like"

tails and thus (informally)

$$\mathbb{P}\left(|\hat{\mu}_i - \mu_i| \geq \sqrt{\frac{2\log(1/\delta)}{T_i(t-1)}}\right) \leq \delta$$



$\mathrm{UCB}_1(t-1,\delta)$

$\mathrm{UCB}_2(t-1,\delta)$

$\hat{\mu}_1(t-1)$  $\mu_1$

$\mu_2$  $\hat{\mu}_2(t-1)$

❑ UCB picks the "correct" arm and guarantees sub-linear regret

❑ Actual regret bounds depend on

▪ choice of $\delta$

▪ sub-optimality gaps, i.e., $\Delta_i \doteq (\max_i \mu_i) - \mu_i$

▪ …

# MAB: Algorithms II

❑ Gaussian rewards, 10 arm problem

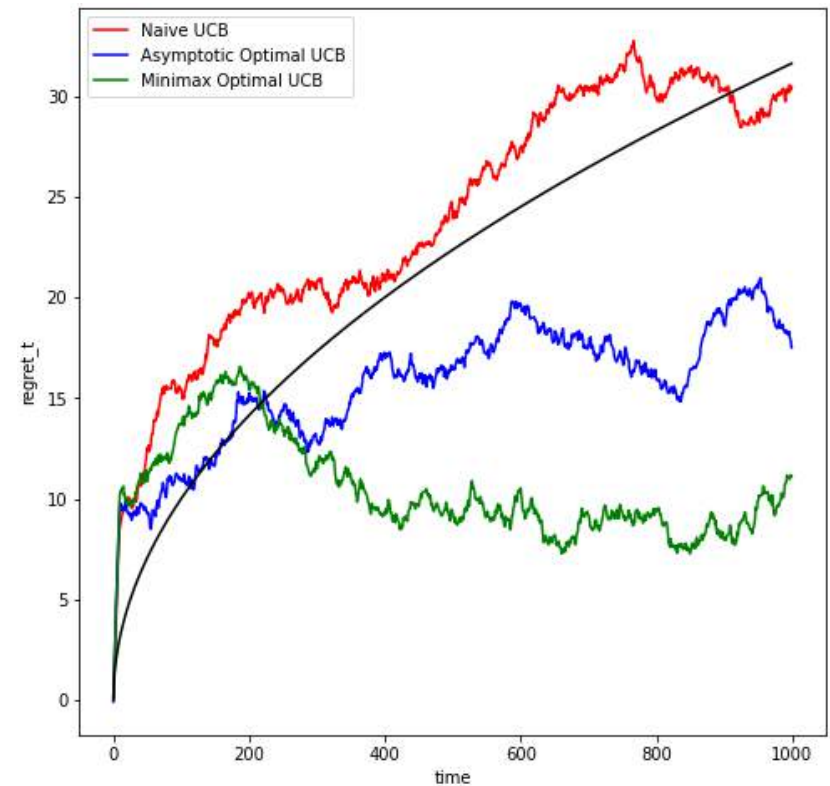- Naïve UCB, Asymptotic UCB, Minimax UCB vary only in choice of $\delta$
- Black line is $y = c\sqrt{t}$

# What is our main result?

- ❑ With our Latin Squares exploration, followed by UCB-based active sampling, we have

- ❑ **Theorem:** With probability at least $1 - o(1)$

$$\mathbb{E}\left[\text{regret}\right] \leq C \sum_{k,l} \frac{\text{correct}_{k,l}^{u}}{\text{sub opt gap}_{k,l}^{u}\,^2} \frac{\text{correct}_{k,l}^{v}}{\text{sub opt gap}_{k,l}^{v}\,^2} \|\text{coord err}\|^2 \frac{\log^2 m}{m}$$

- ▪ Terms for each direction independently – 2 MABs

- ▪ Can exploit prior results on MAB with sub-Gaussian random variables (bounds on regret)
  - – *sub-Gaussianity and concentration inequalities again*

# Main Result

❑ Define

- $\boldsymbol{Y} \doteq \lambda_y^2 \boldsymbol{u}\boldsymbol{v}^\top$ with $\|\boldsymbol{u}\| = \|\boldsymbol{v}\| = 1$ (SVD)

- $b \doteq \max_{i,j} \boldsymbol{Y}_{i,j}$ (max value)

- $\Delta_{k|l}^u \doteq \boldsymbol{Y}_{i*,l} - \boldsymbol{Y}_{k,l}$ and $\Delta_{l|k}^v \doteq \boldsymbol{Y}_{k,j*} - \boldsymbol{Y}_{k,l}$ (sub-optimality gaps)

- $\gamma_{k,l}^u \doteq \boldsymbol{u}_k + 2b\Delta_{k|l}^u$ and $\gamma_{k,l}^v \doteq \boldsymbol{v}_l + 2b\Delta_{l|k}^v$ ($\approx$ correction terms)

- $\boldsymbol{c}_{k,l} \doteq (k,l)^\top$ and $\boldsymbol{c}^* \doteq (i^*, j^*)^\top$ (coordinates)

- $\boldsymbol{R}_m \doteq \frac{1}{m} \sum_{\tau=1}^m \|\hat{\boldsymbol{s}}_\tau - \boldsymbol{s}^*\|^2$ (regret)

❑ **Theorem:** With probability at least $1 - o(1)$

$$\mathbb{E}[\boldsymbol{R}_m] \leq C \sum_{k,l=1}^n \frac{\gamma_{k,l}^u}{(\Delta_{k|l}^u)^2} \cdot \frac{\gamma_{k,l}^v}{(\Delta_{l|k}^v)^2} \cdot \|\boldsymbol{c}_{k,l} - \boldsymbol{c}^*\|^2 \frac{\log^2 m}{m}$$

# Discussion of Result

- ❏ $\Delta_{k|l}^u, \Delta_{l|k}^v$ are "sub-optimality" gaps
  - − as in multi-armed bandit literature, regret $\propto \dfrac{1}{(\Delta_{k|l}^u)^2}$
  - − can potentially be improved to $\dfrac{1}{(\Delta_{k|l}^u)}$ (better stopping time analysis)

- ❏ $\gamma_{k,l}^u, \gamma_{k,l}^v$ are "correction" factors
  - − typical results in MAB consider equal, known variance
  - − our problem – potentially distinct variance estimates

- ❏ $\dfrac{\log^2 m}{m}$ factor standard in MAB regret bounds
  - − best known results (for equal variance case) scale as $\dfrac{\log m}{m}$
  - − Q: can we adapt to our problem? (likely need "better" variance estimates)

$$\mathbb{E}[\boldsymbol{R}_m] \leq C \sum_{k,l=1}^{n} \frac{\gamma_{k,l}^u}{(\Delta_{k|l}^u)^2} \cdot \frac{\gamma_{k,l}^v}{(\Delta_{l|k}^v)^2} \cdot \|\boldsymbol{c}_{k,l} - \boldsymbol{c}^*\|^2 \frac{\log^2 m}{m}$$

# Special Cases

- For Gaussian Energy $h(\boldsymbol{x}, \boldsymbol{y}) := \frac{1}{\sqrt{2\pi\nu^2}} \exp\left(-\frac{\|\boldsymbol{x}-\boldsymbol{y}\|_2^2}{2\nu^2}\right)$

$$\mathbb{E}[\boldsymbol{R}_m] \leq C\nu^2 \sum_{k,l=1}^{n} \frac{\|\boldsymbol{c}_{k,l} - \boldsymbol{c}^*\|^2}{\exp\left(-\frac{\|\boldsymbol{c}_{k,l}-\boldsymbol{c}^*\|^2}{2n\nu^2}\right)} \frac{\log^2 m}{m}$$

- For Laplacian Energy $h(\boldsymbol{x}, \boldsymbol{y}) := \frac{1}{\gamma} \exp\left(-\frac{\|\boldsymbol{x}-\boldsymbol{y}\|_1}{\gamma}\right)$

$$\mathbb{E}[\boldsymbol{R}_m] \leq C\gamma \sum_{k,l=1}^{n} \frac{\|\boldsymbol{c}_{k,l} - \boldsymbol{c}^*\|^2}{\exp\left(-\frac{\|\boldsymbol{c}_{k,l}-\boldsymbol{c}^*\|}{2n\gamma}\right)} \frac{\log^2 m}{m}$$
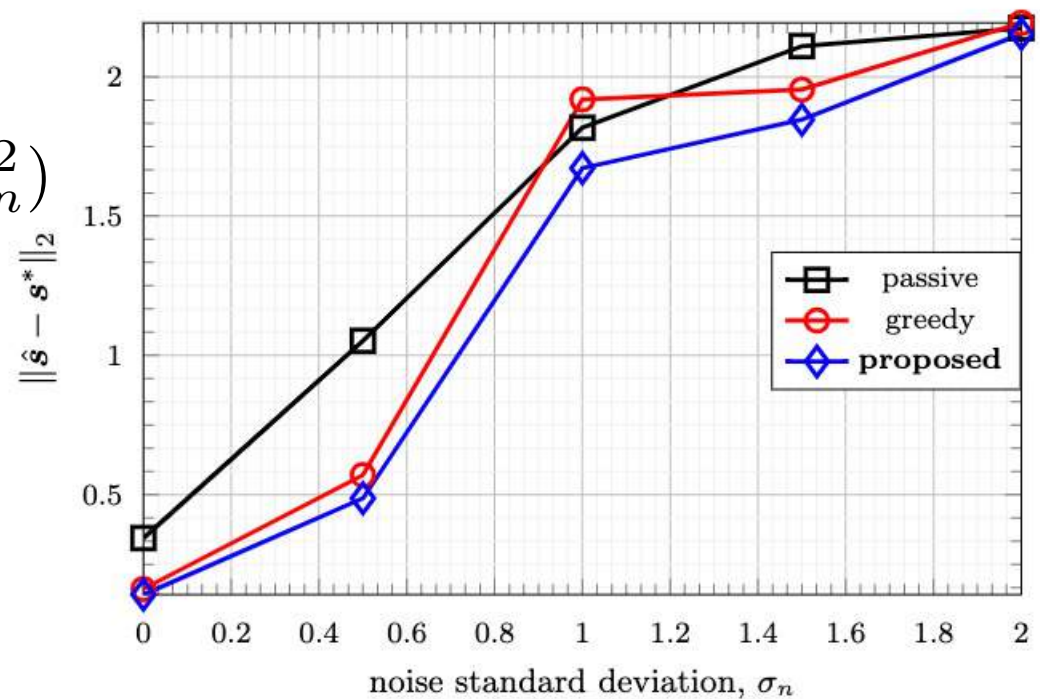
# Variance Parameter

- Laplacian energy function, vary $\gamma$

- As $\gamma$ increases, proposed method better

- Greedy, proposed methods uniformly better than passive

# Measurement Noise

- Gaussian energy function
- add noise, $z_{i,j} \sim \mathcal{N}(0, \sigma_n^2)$
- Proposed method more noise tolerant
- outperforms passive and greedy approaches as expected

# Summary + Future Work

- ☐ Proposed method for active non-parametric peak location
- ☐ Showed experimental improvement for several energy functions
- ☐ Provide preliminary theoretical guarantees

- ☐ Improve error bounds
- ☐ Consider multiple sources
- ☐ Apply to zeroth-order optimization problems

# BIG PICTURE

❑ Active hypothesis testing

▪ So many applications!

▪ Information theory in the wild

❑ Important questions

▪ How do you build your tree of actions/observations?

▪ What is the right measure of informativeness that allows you to prune the tree?

❑ Martingales, concentration inequalities

▪ Very useful tools for a wide-range of applications (need more than the CLT)

❑ The classics still matter

▪ Chernoff, Stein, Wald, Blackwell, Fisher, Bayes, Neyman, Pearson

# thanks

USC Viterbi

Sunav Choudhary

Nicolo Michelusi

Gautam Thatte

Maxime Ferreira da Costa

Praneeth Narayanamurthy

Jianxiu Li

Joni Shaska

Madhavi Rajiv

Mustafa Can Gursoy

Talha Bozkus

Chen Peng

Jeongmin Chae

Daria Riabukhina