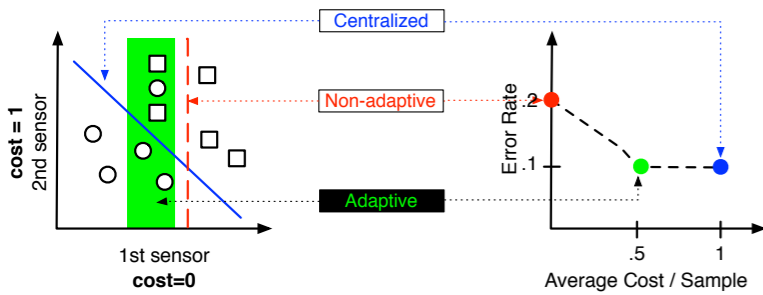


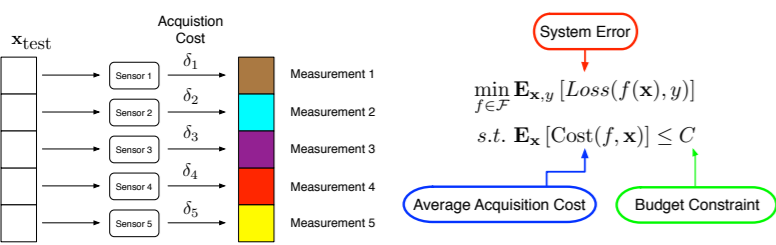
Objective

Discriminative learning framework for sequential decision making under budget constraints

Error Budget Trade-off



Sensor Acquisition Cost



- Sensors:
 - physical measurement in some modalities
 - or computing features of various complexity
 - feature=measurement (possibly high dimensional)
- Typically: higher cost = more informative
- Cost: resources, time, computation ...

Previous Work

- Generative & Parametric Modeling
 - MDP [Ji and Carin, 2007, Kapoor and Horvitz, 2009]
 - Trees [Sheng and Ling, 2006, Bilgic and Getoor, 2007, Zubek and Dietterich, 2002],
 - Utility [Kanani and Melville, 2008]
 - estimate/model $P(x_k | x_j)$
 - not possible in our setting due to high dim.
- Discriminative Methods
 - Detection Cascades: partial binary decisions [Zhang and Zhang, 2010, Chen et al., 2012, Viola and Jones, 2001]
 - TEFE: myopic [Liu et al., 2008]
- Reject Classification [Chow, 1970, Yuan and Casasent, 2003, Bartlett and Wegkamp, 2008, Rodriguez-Díaz and Castañón, 2009]

K. Trapeznikov, V. Saligrama, D. Castañón, *Multi-Stage Classifier Design*, ACML 2012

System Risk

$$R(f^1, f^2, \dots, f^K, \mathbf{x}, y) = \sum_{k=1}^K \underbrace{S^k(\mathbf{x}^k)}_{\text{State Var.}} \underbrace{R_k(\mathbf{x}^k, y, f^k)}_{\text{Stage Risk}}$$

$$\text{Stage: } R_k(\mathbf{x}^k, y, f^k) = \begin{cases} \delta_{k+1}, & f^k(\mathbf{x}^k) = \text{reject} \\ 1, & f^k(\mathbf{x}^k) \neq y \wedge f^k(\mathbf{x}^k) \neq \text{reject} \end{cases}$$

$$\text{State: } S^k(\mathbf{x}^k) = \begin{cases} S^{k-1}(\mathbf{x}^{k-1}), & f^k(\mathbf{x}^k) = \text{rejects} \\ 0, & \text{else} \end{cases}, S^0 = 1$$

Markov Decision Process (MDP)

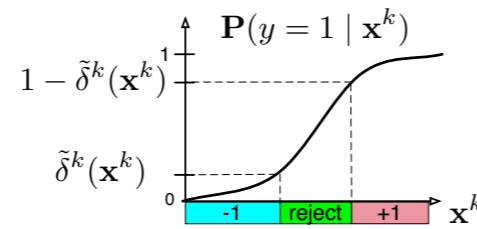
Analyze MDP solution \rightarrow structure to aid in ERM

- If $\mathbf{P}(\mathbf{x}, y)$ known then objective:

$$\min_{f^1, f^2, \dots, f^K} \mathbf{E}[R(\cdot) | \mathbf{x}]$$
- Solution through a Dynamic Program (DP) for $f^k(\mathbf{x}^k)$:

$$\tilde{\delta}^k(\mathbf{x}^k) = \min_{f^{k+1}, \dots, f^K} \mathbf{E} \left[\sum_{t=k+1}^K S^t(\mathbf{x}^t) R_t(\mathbf{x}^t, y, f^t) \mid \mathbf{x}^k, S^k = 1 \right] + \underbrace{\delta_{k+1}}_{\text{meas. cost}}$$

expected risk of stages $k+1 \dots K$



- Modified Stage Risk

$$\tilde{R}_k(\mathbf{x}^k, y, f^k, \tilde{\delta}^k) = \begin{cases} \tilde{\delta}^k(\mathbf{x}^k), & f^k(\mathbf{x}^k) = r \\ 1, & f^k(\mathbf{x}^k) \neq y \wedge f^k(\mathbf{x}^k) \neq r \end{cases}$$

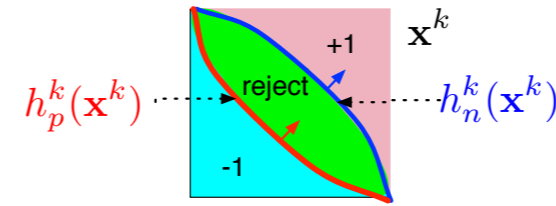
$$f^k = \arg \min_f \mathbf{E} [\tilde{R}_k(\mathbf{x}^k, y, f, \tilde{\delta}^k) \mid \mathbf{x}^k] = \arg \min_f \mathbf{E} [R(\cdot) \mid \mathbf{x}^k]$$

Given $\tilde{\delta}^k(\mathbf{x}^k)$: Multi-Stage Risk Minimization \rightarrow Single Stage

Parameterization of Reject Option

$f^k(\mathbf{x}^k)$ is a classifier with reject option \rightarrow reduce reject decision to supervised binary classification

Biased Classifiers

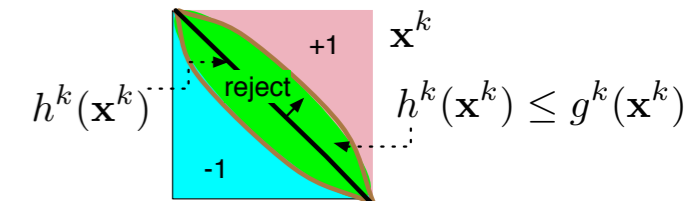


$$f^k(\mathbf{x}^k) = \begin{cases} \text{sgn} [h_p^k(\mathbf{x}^k)], & \text{sgn} [h_p^k(\mathbf{x}^k)] = \text{sgn} [h_n^k(\mathbf{x}^k)] \\ \text{reject}, & \text{sgn} [h_p^k(\mathbf{x}^k)] \neq \text{sgn} [h_n^k(\mathbf{x}^k)] \end{cases}$$

$$\tilde{R}_k(\mathbf{x}^k, y_i, h_p, h_n, \tilde{\delta}_i^k) = \underbrace{\mathbb{1}_{[h_p(\mathbf{x}^k) y_i \leq 0]} \mathbb{1}_{[h_n(\mathbf{x}^k) y_i \leq 0]}}_{\text{error penalty if not reject}} + \underbrace{\tilde{\delta}_i^k}_{\text{cost-to-go}} \underbrace{\left\{ \mathbb{1}_{[h_p(\mathbf{x}^k) y_i \leq 0]} + \mathbb{1}_{[h_n(\mathbf{x}^k) y_i \leq 0]} - 2 \mathbb{1}_{[h_p(\mathbf{x}^k) y_i \leq 0]} \mathbb{1}_{[h_n(\mathbf{x}^k) y_i \leq 0]} \right\}}_{\text{rejected}}$$

- restricted to binary setting

Margin Based



$$f^k(\mathbf{x}^k) = \begin{cases} \text{sgn} [h^k(\mathbf{x}^k)], & |h^k(\mathbf{x}^k)| > g^k(\mathbf{x}^k) \\ \text{reject}, & |h^k(\mathbf{x}^k)| \leq g^k(\mathbf{x}^k) \end{cases}$$

$$\tilde{R}_k(\mathbf{x}^k, y_i, h^k, g^k) = \underbrace{\mathbb{1}_{[h(\mathbf{x}^k) y_i \leq 0]} \mathbb{1}_{[|h(\mathbf{x}^k)| > g(\mathbf{x}^k)]}}_{\text{error penalty not rejected}} + \underbrace{\tilde{\delta}_i^k}_{\text{cost to go}} \underbrace{\mathbb{1}_{[|h(\mathbf{x}^k)| \leq g(\mathbf{x}^k)]}}_{\text{rejected}}$$

- extends to multi-class setting (MC to Binary)

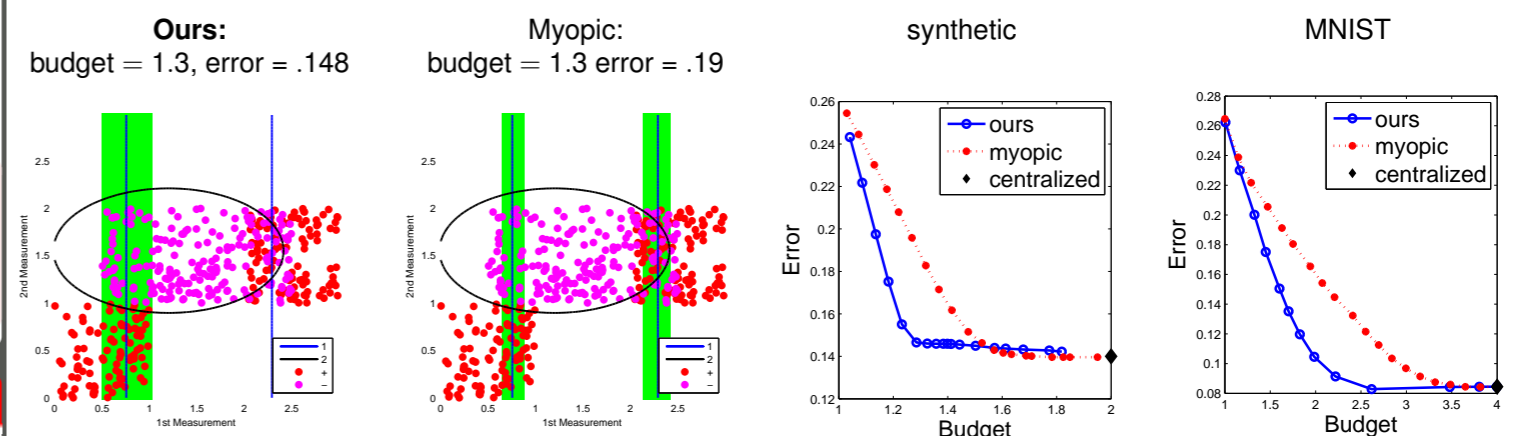
Alternating Minimization

- $\min \tilde{R}_k(\cdot) \rightarrow$ series of supervised learning problems
- cyclical optimization of one stage k at a time
- surrogate loss: $\mathbb{1}_{[z]} \rightarrow \mathcal{L}[z]$
- smooth global objective \implies coordinate descent converges to a local minimum

Generalization

- Polynomial Kernel Classifiers: complexity is bounded $K \log K \times$ most complex stage
- Boosted Classifiers: margin based bound for a two stage system

Numerical Experiments



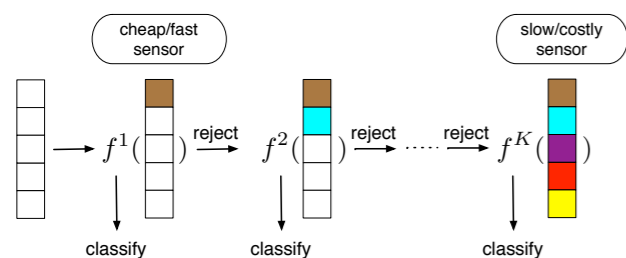
- Myopic: threshold margin of stage classifier to reject constant fraction
- Utility: threshold expected utility of stage classifier

Dataset	Size	Stage 1	Stage 2	Stage 3	Stage 4	# Classes	Target Error	Myopic	Ours	Utility
synthetic	4,000	Sensor 1	Sensor 2	2	.147	52%	28%	
pima	768	weight, age, ..	glucose test	insulin test	..	2	.245	41%	15%	
threat	1230	PMMW image	IR image	AMMW image	..	2	.16	89%	71%	
coverttype	581012	soils	wild. areas	elev. aspect,	7	.285	79%	40%	
letter	20000	pixel counts	moments	edge feat's	..	26	.25	81%	51%	
mnist	70000	4 x 4 image	7 x 7	14 x 14	28 x 28	10	.085	90%	52%	
landsat	6435	Band 1	Band 2	Band 3	Band 4	7	.17	56%	31%	
mammogram	830	CAD feat's	expert rating	2	.173		25%	65%

- Performance: % of the maximum budget required to achieve the target error rate
- Target rate is chosen to be close to the error of the centralized strategy

Multi-Stage Sequential Reject Classifier

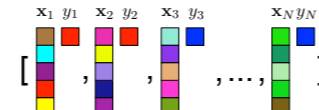
- Sample: $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_K]$, True label: $y \in \{1, 2, \dots, C\}$



- Order of stages/sensors is fixed
- k th stage:
 - $f^k(\mathbf{x}^k)$: full decision with a reject option
 - acquires k th meas. for a cost δ_k
 - \mathbf{x}^k : first k meas. of \mathbf{x}

Stage-Wise Empirical Risk Minimization

- Training Data with full measurements:



- Point-wise Cost-to-go Empirical Estimate:

$$\tilde{\delta}_i^{k-1} = S_i^k \tilde{R}_k(\mathbf{x}_i^k, y_i, f^k, \tilde{\delta}_i^k) + \delta_i^k, \quad i = 1, 2, \dots, N$$

Instead of learning $\tilde{\delta}^k(\mathbf{x}^k)$, use $\tilde{\delta}_i^k$ to learn **decision boundaries directly**

- Empirical Risk Minimization for stage k :

$$f^k(\mathbf{x}^k) = \arg \min_{f \in \mathcal{F}^k} \frac{1}{N} \sum_{i=1}^N S_i^k \tilde{R}_k(y_i, \mathbf{x}_i^k, f, \tilde{\delta}_i^k)$$