

Online Discriminative Learning Theory and Applications

Nicolò Cesa-Bianchi

Università degli Studi di Milano



Summary

- 1 The online protocol
- 2 Structured classification
- 3 Active learning
- 4 Multiview learning



Summary

- 1 The online protocol
- 2 Structured classification
- 3 Active learning
- 4 Multiview learning



Theory of repeated games



James Hannan



David Blackwell

Learning to play a game (1956)

Play a game **repeatedly** against a **possibly suboptimal** opponent

Zero-sum 2-person games played more than once

	1	2	...	M
1	$\ell(1,1)$	$\ell(1,2)$...	
2	$\ell(2,1)$	$\ell(2,2)$...	
\vdots	\vdots	\vdots	\ddots	
N				

$N \times M$ known loss matrix

- Row player (**player**) has N actions
- Column player (**opponent**) has M actions

For each game round $t = 1, 2, \dots$

- Player chooses action i_t and opponent chooses action y_t
- The player suffers loss $\ell(i_t, y_t)$ (= gain of opponent)



Zero-sum 2-person games played more than once

	1	2	...	M
1	$\ell(1,1)$	$\ell(1,2)$...	
2	$\ell(2,1)$	$\ell(2,2)$...	
\vdots	\vdots	\vdots	\ddots	
N				

$N \times M$ known loss matrix

- Row player (**player**) has N actions
- Column player (**opponent**) has M actions

For each game round $t = 1, 2, \dots$

- Player chooses action i_t and opponent chooses action y_t
- The player suffers loss $\ell(i_t, y_t)$ (= gain of opponent)

Player can learn from opponent's history of past choices y_1, \dots, y_{t-1}



Prediction with expert advice, 1989



Volodya Vovk



Manfred Warmuth

- Opponent's moves y_1, y_2, \dots define a **nonstochastic sequential prediction problem**
- Loss matrix can **change with time**
- Design a player's strategy that predicts **any sequence** y_1, y_2, \dots nearly as well as the **single best action** for that sequence

Exponentially weighted forecaster

At time t pick action i with probability proportional to

$$\exp(-\eta \text{Loss}_{i,t})$$

where $\text{Loss}_{i,t}$ is **total loss** of action i up to now

Theorem

[C-B, 1997]

The average per-round expected loss of the forecaster converges to that of the **best action for the observed sequence** at rate

$$\sqrt{\frac{2}{T} \ln N}$$

where N is number of experts and T is the number of time steps



Exponentially weighted forecaster

At time t pick action i with probability proportional to

$$\exp(-\eta \text{Loss}_{i,t})$$

where $\text{Loss}_{i,t}$ is **total loss** of action i up to now

Theorem

[C-B, 1997]

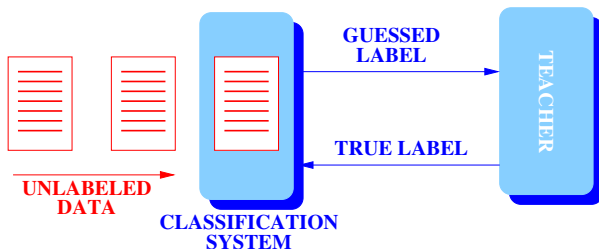
The average per-round expected loss of the forecaster converges to that of the **best action for the observed sequence** at rate

$$\sqrt{\frac{2}{T} \ln N}$$

where N is number of experts and T is the number of time steps

No dependence on number of opponent's actions!

From game theory to online learning



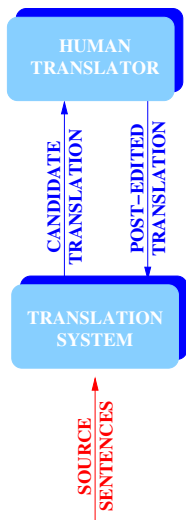
- Add **side info** to opponent's moves $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots$ $\mathbf{x}_t \in \mathbb{R}^d$
- Linear classifiers \mathbf{w} which predict using $\mathbf{w}^\top \mathbf{x}_t$
- A repeated game between the player choosing action $\mathbf{w}_t \in \mathbb{R}^d$ and the opponent choosing action (\mathbf{x}_t, y_t)
- Convergence to performance of **best linear classifier** under no statistical assumptions on data source

Advantages of online algorithms

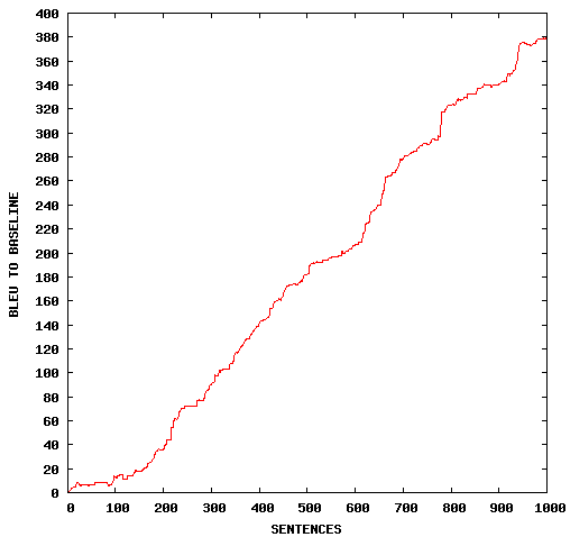
- Simple algorithms for learning **linear models**
- **Scalable**
- **Robust:** game-theoretic performance guarantees
- **Versatile:**
 - structured classification
 - active learning
 - matrix learning
 - tracking
 - bounded memory learning



Online protocol naturally exploits interaction



Computer Assisted Translation



Summary

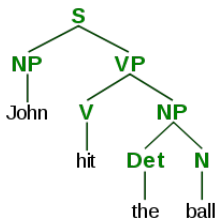
- 1 The online protocol
- 2 Structured classification
- 3 Active learning
- 4 Multiview learning



Structured Classification

A combinatorial label space (sequences, trees)

- **POS tagging:** sentence \rightarrow sequence of POS tags
- **Parsing:** sentence \rightarrow parse tree
- **Bilingual alignment:** sentence pair \rightarrow alignment (matching)
- **Letter to phoneme:** word \rightarrow phoneme sequence
- **Phrase-based translation:** source sentence \rightarrow target sentence



room	.	.	■	.	.	.
the	.	■
in	■
cold
too	■
is	■	.
it
en
la
habitacion
hace	.	■
demasiado
frio

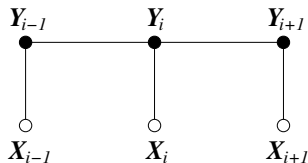


Conditional Random Fields [Lafferty, McCallum and Pereira, 2001]

- A graphical model $\mathbb{P}(\mathbf{y} \mid \mathbf{x})$ for predicting a set of **hidden variables** $\mathbf{y} = (y_1, \dots, y_m)$ given observation \mathbf{x}
- The conditional dependency of \mathbf{y} on \mathbf{x} is defined through **feature functions** $f_j(\mathbf{x}, \mathbf{y}_i) \in \mathbb{R}$
- $\mathbf{y}_1, \dots, \mathbf{y}_m$ are subparts of \mathbf{y} (e.g., single nodes, linked nodes)

Some simple examples

- (n-gram, POS tag)
- NP \rightarrow DT NN
- (le chat, the cat)



Conditional Random Fields – Inference

Loglinear model for the joint distribution of labels

$$\mathbb{P}(\mathbf{y} \mid \mathbf{x}) = \exp \left(\sum_{i,j} w_j f_j(\mathbf{x}, \mathbf{y}_i) - \ln Z_{\mathbf{x}} \right) = \exp \left(\mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y}) - \ln Z_{\mathbf{x}} \right)$$

Decoding

$$\begin{aligned} \hat{\mathbf{y}} &= \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mathbb{P}(\mathbf{y} \mid \mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \exp \left(\mathbf{w}^{\top} \mathbf{f}(\mathbf{x}, \mathbf{y}) \right) \\ &= \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \sum_{i,j} w_j f_j(\mathbf{x}, \mathbf{y}_i) \end{aligned}$$

$\mathcal{Y}(\mathbf{x})$ is the set of feasible labels for observation \mathbf{x}



Recall

Prediction via decoding $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$

Stochastic gradient ascent

Estimate \mathbf{w} by maximizing the log-likelihood of a training set $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{\mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbb{E}_{\mathbf{Y}}[\mathbf{f}(\mathbf{x}_t, \mathbf{Y})]}_{\nabla \ln \mathbb{P}(\mathbf{y}_t | \mathbf{x}_t)} \quad \text{for } t = 1, \dots, T$$



Recall

Prediction via decoding $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} \mathbf{w}^\top \mathbf{f}(\mathbf{x}, \mathbf{y})$

Stochastic gradient ascent

Estimate \mathbf{w} by maximizing the log-likelihood of a training set $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$

$$\mathbf{w} \leftarrow \mathbf{w} + \underbrace{\mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbb{E}_{\mathbf{Y}}[\mathbf{f}(\mathbf{x}_t, \mathbf{Y})]}_{\nabla \ln \mathbb{P}(\mathbf{y}_t | \mathbf{x}_t)} \quad \text{for } t = 1, \dots, T$$

Structured Perceptron

A Viterbi approximation $\mathbf{f}(\mathbf{x}_t, \hat{\mathbf{y}}_t)$ of the expectation $\mathbb{E}_{\mathbf{Y}}[\mathbf{f}(\mathbf{x}_t, \mathbf{Y})]$

$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \hat{\mathbf{y}}_t)$$

Structured Perceptron

$$\hat{y}_t = \operatorname{argmax}_{y \in \mathcal{Y}(x)} \mathbf{w}_t^\top \mathbf{f}(x_t, y_t)$$

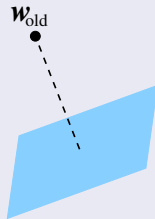
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \mathbf{f}(x_t, y_t) - \mathbf{f}(x_t, \hat{y}_t)$$

- Simplest example of online linear algorithm
- Generates an ensemble $\mathbf{0} = \mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_T$ of classifiers by performing multiple epochs over the training set
- Average classifier: $\frac{1}{T} \sum_t \mathbf{w}_t$ has typically low risk
- If training set is i.i.d., then the **statistical risk** can be provably bounded \rightarrow online to batch conversion



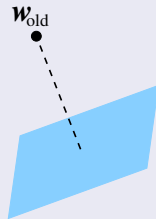
Margin optimization associated with the 1-best label for x_t

$$\begin{aligned} \min_{\mathbf{w}} \quad & (\|\mathbf{w} - \mathbf{w}_{\text{old}}\|^2 + C \xi_t) \\ \text{s.t.} \quad & \underbrace{\mathbf{w}^\top \mathbf{f}(x_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(x_t, \hat{\mathbf{y}}_t)}_{\text{linear constraint}} \geq 1 - \xi_t \end{aligned}$$



Margin optimization associated with the 1-best label for x_t

$$\begin{aligned} \min_{\mathbf{w}} \quad & (\|\mathbf{w} - \mathbf{w}_{\text{old}}\|^2 + C \xi_t) \\ \text{s.t.} \quad & \underbrace{\mathbf{w}^\top \mathbf{f}(x_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(x_t, \hat{\mathbf{y}}_t)}_{\text{linear constraint}} \geq 1 - \xi_t \end{aligned}$$



The solution is the Passive-Aggressive algorithm

$$\mathbf{w} = \mathbf{w}_{\text{old}} + \eta_t (\mathbf{f}(x_t, \mathbf{y}_t) - \mathbf{f}(x_t, \hat{\mathbf{y}}_t))$$

the learning rate η_t has a closed form expression



More global updates

- Gradient ascent update tends to enforce all **margin constraints**

$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}) \geq 1 \quad \text{for all } \mathbf{y} \neq \mathbf{y}_t$$

- Perceptron addresses only $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \hat{\mathbf{y}}_t) \geq 1$
- Consider the N -best labels \mathbf{y} in the ranking induced by $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y})$



More global updates

- Gradient ascent update tends to enforce all **margin constraints**
$$\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}) \geq 1 \quad \text{for all } \mathbf{y} \neq \mathbf{y}_t$$
- Perceptron addresses only $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \hat{\mathbf{y}}_t) \geq 1$
- Consider the N -best labels \mathbf{y} in the ranking induced by $\mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y})$

MIRA algorithm

[Crammer and Singer, 2003]

- Address all the constraints in the N -best list

$$\min_{\mathbf{w}} \quad (\|\mathbf{w} - \mathbf{w}_{\text{old}}\|^2 + C \xi_t)$$

$$\text{s.t.} \quad \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}^{(i)}) \geq 1 - \xi_t \quad \text{for } i = 1, \dots, N$$

- Use a **SVM solver** on a pseudo training set with instances

$$\mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{f}(\mathbf{x}_t, \mathbf{y}^{(i)}) \quad i = 1, \dots, N$$

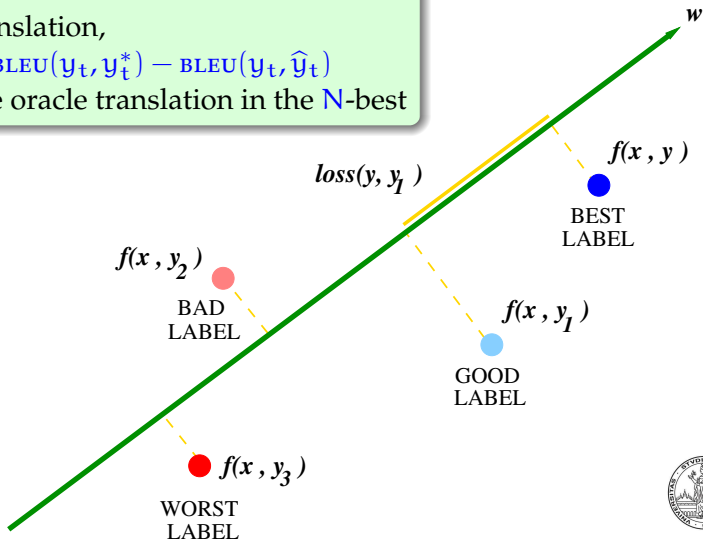
Cost-sensitive margin-based learning

$$\mathbf{w}^\top \mathbf{f}(x_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(x_t, \hat{\mathbf{y}}_t) \geq \text{loss}(\mathbf{y}_t, \hat{\mathbf{y}}_t) - \xi_t$$

In machine translation,

$$\text{loss}(\mathbf{y}_t, \hat{\mathbf{y}}_t) = \text{BLEU}(\mathbf{y}_t, \mathbf{y}_t^*) - \text{BLEU}(\mathbf{y}_t, \hat{\mathbf{y}}_t)$$

where \mathbf{y}_t^* is the oracle translation in the N -best



Using duality

Cost-sensitive passive-aggressive optimization at \mathbf{x}_t

$$\begin{aligned} \min_{\mathbf{w}} \quad & (\|\mathbf{w} - \mathbf{w}_{\text{old}}\|^2 + C \xi_t) \\ \text{s.t.} \quad & \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \mathbf{y}_t) - \mathbf{w}^\top \mathbf{f}(\mathbf{x}_t, \hat{\mathbf{y}}_t) \geq \text{loss}(\mathbf{y}_t, \hat{\mathbf{y}}_t) - \xi_t \end{aligned}$$

Learning rate is the dual-maximizing Lagrange multiplier

The dual Lagrange function $\mathcal{D}_t(\alpha)$ **always** satisfies

$$\begin{aligned} \operatorname{argmax}_{0 \leq \alpha \leq C} \mathcal{D}_t(\alpha) &= \eta_t \end{aligned}$$



Theorem

For *any* sequence $(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$

$$\sum_{t=1}^T \text{loss}(\mathbf{y}_t, \hat{\mathbf{y}}_t) \leq \sum_{t=1}^T \mathcal{D}_t(\eta_t) \leq \min_{\mathbf{w}} \mathcal{P}(\mathbf{w}, (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T))$$

- \mathcal{P} is a convex SVM-style objective that trades-off $\|\mathbf{w}\|^2$ with the overall extent by which the linear margin constraints in the N -best lists are violated by \mathbf{w}
- Bound extends to any algorithm that optimizes over more constraints (e.g., MIRA)

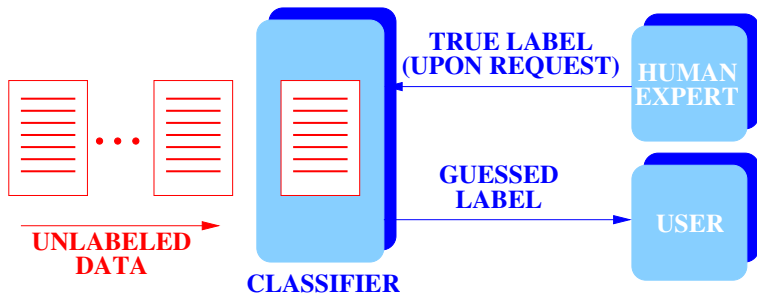


Summary

- 1 The online protocol
- 2 Structured classification
- 3 Active learning**
- 4 Multiview learning



Online active learning



Learning a binary classifier

Setup

- **Data process** $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$
- **Label process** $y_1, y_2, \dots \in \{-1, +1\}$

Assumptions

- Observing the data process is “cheap”
- Observing the label process is “expensive”
→ need to query the human expert



Learning a binary classifier

Setup

- **Data process** $\mathbf{x}_1, \mathbf{x}_2, \dots \in \mathbb{R}^d$
- **Label process** $y_1, y_2, \dots \in \{-1, +1\}$

Assumptions

- Observing the data process is “cheap”
- Observing the label process is “expensive”
→ need to query the human expert

Question

How much better can we do by subsampling adaptively the label process?

- Adaptive design in Statistics [Zacks, 2009]

- Uncertainty sampling [Cohn, Atlas and Ladner, 1990]

- Query by committee [Freund, Seung, Shamir and Tishby, 1997]
Modified Perceptron [Dasgupta, Kalai and Monteleoni, 2005]
→ exponential advantage in the noise-free case

- General strategies for the noisy case
[Balcan, Beygelzimer and Langford, 2006]
[Balcan, Broder and Zhang, 2007]
[Dasgupta, Hsu and Monteleoni, 2008]



A parametric classifier

Regularized least squares (slightly modified)

- $\hat{f}_t(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$ where

$$\mathbf{w} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left(\|S^\top \mathbf{w} - \mathbf{y}\|^2 + \|\mathbf{w}\|^2 + (\mathbf{w}^\top \mathbf{x})^2 \right)$$

- $S = [\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_N}]$ is matrix of queried data points
- $\mathbf{y} = (y_{t_1}, \dots, y_{t_N})$ is vector of queried labels

Remarks

- Involves only inner products
→ replaceable by kernel functions
- Time quadratic (in number of queries) for incremental update

- View $\hat{f}_t(\mathbf{x})$ as a biased estimator of $f^*(\mathbf{x})$
- Use large deviation analysis to obtain **confidence interval**

$$\hat{f}_t(\mathbf{x}_t) \pm c \sqrt{\frac{\ln t}{N_t}} \quad \text{for confidence level } 1 - t^{-1}$$

- N_t is number of queries up to time t



- View $\hat{f}_t(\mathbf{x})$ as a biased estimator of $f^*(\mathbf{x})$
- Use large deviation analysis to obtain **confidence interval**

$$\hat{f}_t(\mathbf{x}_t) \pm c \sqrt{\frac{\ln t}{N_t}} \quad \text{for confidence level } 1 - t^{-1}$$

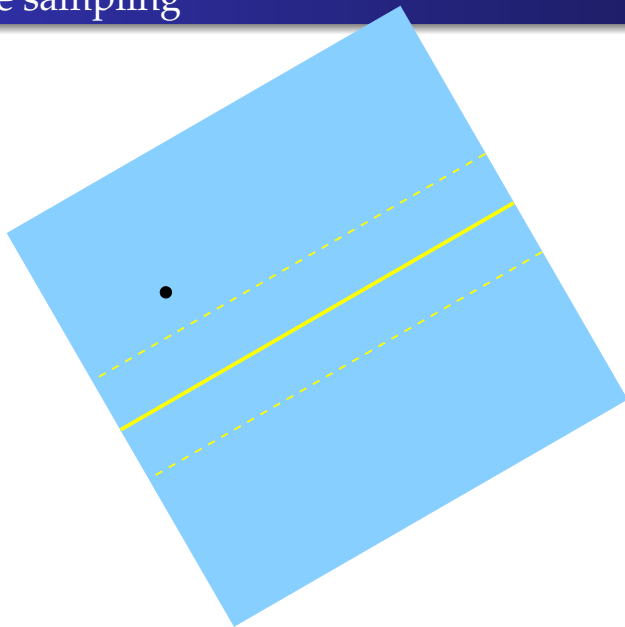
- N_t is number of queries up to time t

Query if *margin* of current point is small

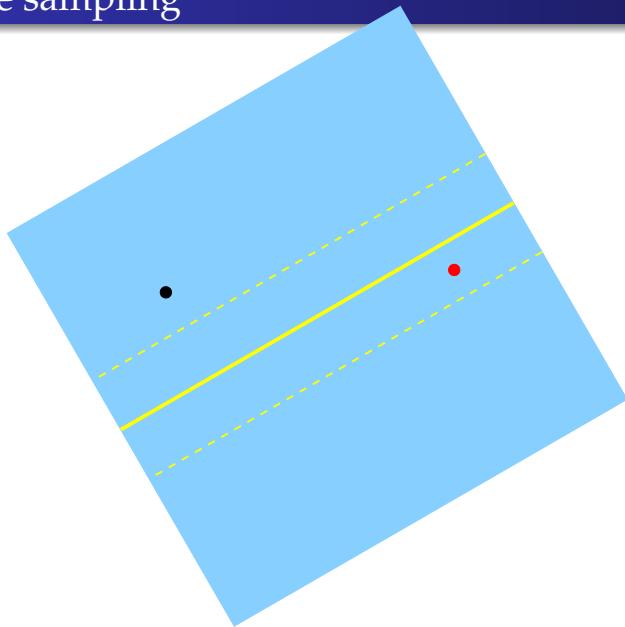
$$\text{If } |\hat{f}_t(\mathbf{x}_t)| \leq c \sqrt{\frac{\ln t}{N}} \quad \text{then query } \mathbf{x}_t$$



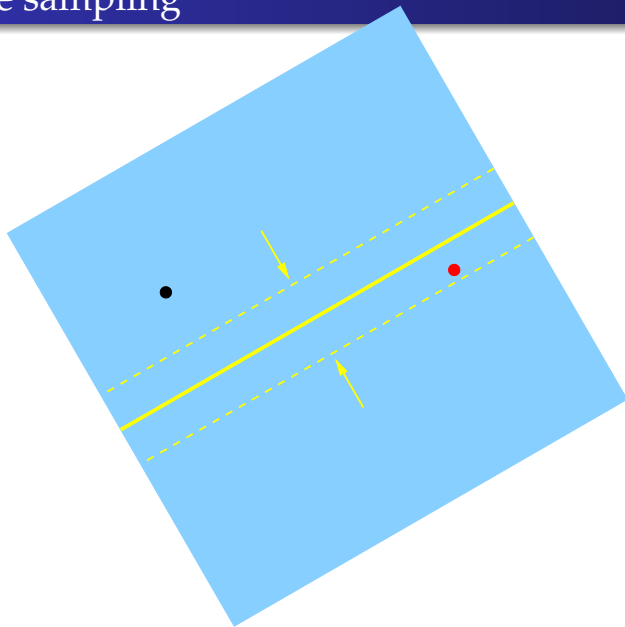
Selective sampling



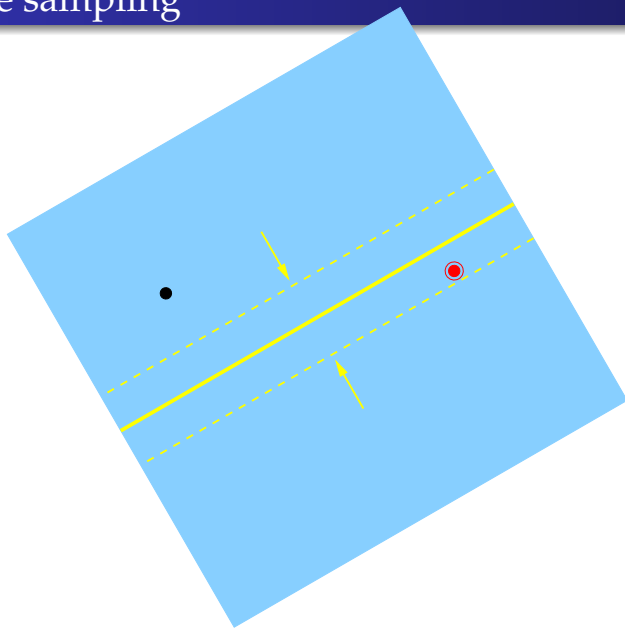
Selective sampling



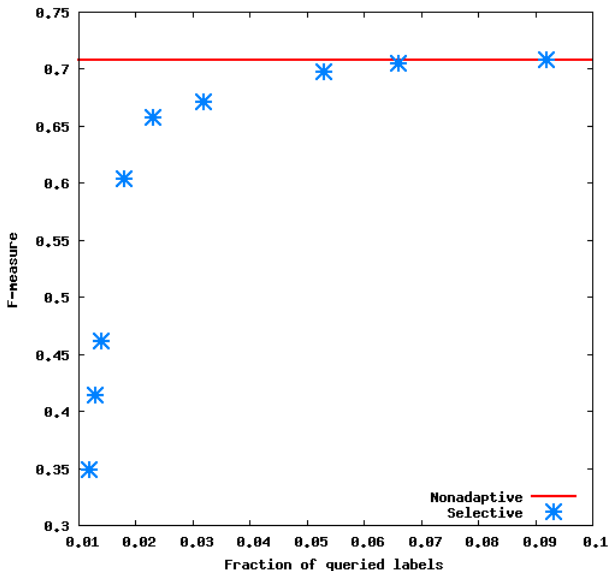
Selective sampling



Selective sampling



Results on text categorization



Summary

- 1 The online protocol
- 2 Structured classification
- 3 Active learning
- 4 Multiview learning



Some applications

- Learn a CD-HMM (emission densities are Gaussian mixtures)
[Cheng, Sha and Saul, 2009]
- Learn the best subspace projection (online PCA)
[Warmuth and Kuzmin, 2008]
- Multitask/multiview learning
[Cavallanti, C-B and Gentile, 2008]



Multiview binary classification

Definitions

- Each data element (instance) is described by a set of tuples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \in \mathbb{R}^d$ (views) each obtained from different measurements \rightarrow e.g., **multimodal learning**
- Build $d \times K$ matrix of views
- Run online linear classification algorithms in the linear space with inner product $\langle \mathbf{W}, \mathbf{X} \rangle = \text{tr}(\mathbf{W}^\top \mathbf{X})$
- Binary prediction $\text{sgn}(\langle \mathbf{W}, \mathbf{X}_t \rangle)$

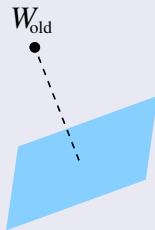
- How to do provably better than combining all views in a vector?
- Transfer information across views in order to reduce the learning complexity



Multiview binary classification

Margin optimization (binary case) for X_t

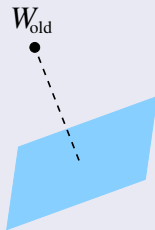
$$\begin{aligned} \min_W \quad & (\|W - W_{\text{old}}\|^2 + C \xi_t) \\ \text{s.t.} \quad & \underbrace{y_t \langle W, X_t \rangle}_{\text{linear constraint}} \geq 1 - \xi_t \end{aligned}$$



Multiview binary classification

Margin optimization (binary case) for X_t

$$\begin{aligned} \min_W \quad & (\|W - W_{\text{old}}\|^2 + C \xi_t) \\ \text{s.t.} \quad & \underbrace{y_t \langle W, X_t \rangle}_{\text{linear constraint}} \geq 1 - \xi_t \end{aligned}$$



Spectral co-regularization

- $\|\cdot\|$ = Frobenius norm: columns of W updated independently
→ Perceptron/Passive-Aggressive on combined views
- $\|\cdot\|$ = any **unitarily invariant** norm: columns may interact

Matrix p -norm Perceptron

- Schatten p -norm of W is the p -norm of the SVD vector
- Prediction is sign of $(V^T V)^{p-2} V^T X$
- V is updated using the Perceptron rule $V \leftarrow V + y X$



Matrix p -norm Perceptron

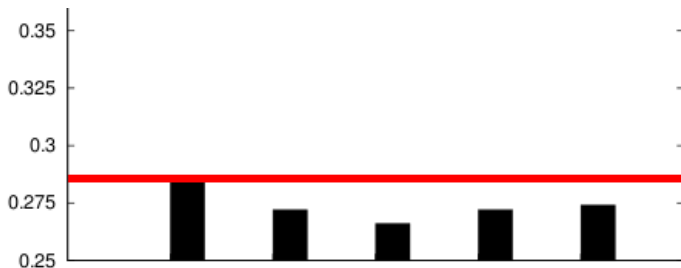
- Schatten p -norm of W is the p -norm of the SVD vector
- Prediction is sign of $(V^T V)^{p-2} V^T X$
- V is updated using the Perceptron rule $V \leftarrow V + y X$

Theory

- 1 If the best model W is simple $\text{rank}(W) \ll K$
 - 2 and views X_t are informative $\text{rank}(X_t) \approx K$
- improve over combined views by a factor of K



Experimental results



- Gene function classification using six views
- Values $p = 2, 4, 6, 8, 10$
- Red line ($p = 2$) is error of Perceptron on combined views



THANK YOU

