

Online Passive- Aggressive Algorithms

Jean-Baptiste Behuet
Tutor: Eneldo Loza Mencía

28/11/2007

Seminar aus Maschinellem Lernen

WS07/08

Overview

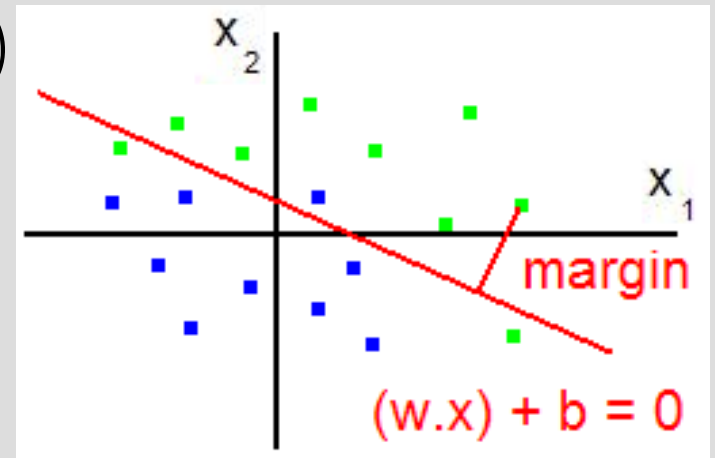
- Online algorithms
- Online Binary Classification Problem
 - Perceptron Algorithm
 - 3 versions of the Passive-Aggressive Algorithm
 - Loss bounds, Comparison with the Perceptron
- Other learning problems
- Experiments
- Conclusion

Online Algorithms

- Sequence of rounds t :
 - Instance x_t as input
 - Predicts \hat{y}_t as output
 - Receives correct output y_t
 - Updates prediction mechanism

Online Binary Classification: Perceptron Algorithm

- Round t :
 - instance $x_t \in \mathbb{R}^n$ with label $y_t \in \{-1, 1\}$
 - classification function based on weight vector $w_t \in \mathbb{R}^n \rightarrow$ defines hyperplane separating the 2 classes
 - prediction: $\hat{y}_t = \text{sign}((w_t \cdot x_t) + b)$
 - signed margin: $y_t((w_t \cdot x_t) + b)$
 - correct if $\text{margin} > 0$
- Goal: incrementally learn w_t
(incrementally modify hyperplane)



Online Binary Classification: Perceptron Algorithm (2)

- Start with random hyperplane (random w_0)
- At each round t of the algorithm:
 - receives x_t and predicts $\hat{y}_t = \text{sign}(w_t \cdot x_t + b_t)$
 - receives correct y_t and updates the hyperplane
- Update minimizes the distance of misclassified examples to the boundary
$$w_{t+1} = w_t + \rho (y_t - \hat{y}_t) x_t \quad \text{with } \rho > 0 \text{ learning rate}$$
(the hyperplane is updated when an error occurs)

Passive-Aggressive Algorithm for binary classification

- Want: $margin \geq 1$ as often as possible
(not only correctly classified examples)
- Hinge-loss function:

$$\ell(\mathbf{w}; (\mathbf{x}, y)) = \begin{cases} 0 & y(\mathbf{w} \cdot \mathbf{x}) \geq 1 \\ 1 - y(\mathbf{w} \cdot \mathbf{x}) & \text{otherwise} \end{cases}$$

- Loss suffered at round t : $\ell_t = \ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))$
- Number of prediction mistakes $\leq \sum l_t^2$

Passive-Aggressive Algorithm for binary classification (2)

- Initialization: $w_1 = (0, \dots, 0)$
- Update:
 - w_{t+1} solution of constrained optimization problem:

$$w_{t+1} = \operatorname{argmin}_{w \in \mathbb{R}^n} \frac{1}{2} \|w - w_t\|^2 \quad \text{s.t.} \quad \ell(w; (\mathbf{x}_t, y_t)) = 0.$$

- w_{t+1} has the form: $w_{t+1} = w_t + \tau_t y_t \mathbf{x}_t$ where $\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2}$

Passive-Aggressive Algorithm for binary classification (3)

- Trade-off:
 - w_{t+1} required to have no loss on current example
 - w_{t+1} as close as possible to w_t
- " Passive-Aggressive " :
 - " passive " when $l_t = 0$
 - " aggressive " otherwise: w_{t+1} forced to satisfy the constraint $l_t = 0$ on the current example

Two variations of the PA algorithm

- Problem of aggressiveness in case of noise

- PA-I:
$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi \quad \text{s.t.} \quad \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) \leq \xi \quad \text{and} \quad \xi \geq 0$$

- PA-II:
$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 + C\xi^2 \quad \text{s.t.} \quad \ell(\mathbf{w}; (\mathbf{x}_t, y_t)) \leq \xi$$

with C aggressiveness parameter

- same update form as for PA:
$$\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$$

$$\tau_t = \min \left\{ C, \frac{\ell_t}{\|\mathbf{x}_t\|^2} \right\} \quad \text{for PA-I} \quad \tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} \quad \text{for PA-II}$$

Relative loss bounds

- Number of prediction mistakes $\leq \sum l_t^2$
- Comparison of the loss attained by PA with the loss attained by a fixed classifier $\text{sign}(u \cdot x)$

$$l_t = \ell(\mathbf{w}_t; (\mathbf{x}_t, y_t))$$

$$l_t^* = \ell(\mathbf{u}; (\mathbf{x}_t, y_t))$$

- For the original PA algorithm:

$$\sum_{t=1}^T l_t^2 \leq \|\mathbf{u}\|^2 R^2$$

$$\forall t \quad \|x_t\| \leq R \quad \text{and} \quad l_t^* = 0$$

Relative loss bounds (2)

- For the original PA algorithm:

$$\sum_{t=1}^T \ell_t^2 \leq \left(\|\mathbf{u}\| + 2\sqrt{\sum_{t=1}^T (\ell_t^*)^2} \right)^2 \quad \forall t \quad \|x_t\| = 1, \quad \forall u \in \mathbb{R}^n$$

- For PA-I:

$$M \leq \max \{R^2, 1/C\} \left(\|\mathbf{u}\|^2 + 2C \sum_{t=1}^T \ell_t^* \right) \quad \forall t \quad \|x_t\| \leq R, \quad \forall u \in \mathbb{R}^n$$

- For PA-II:

$$\sum_{t=1}^T \ell_t^2 \leq \left(R^2 + \frac{1}{2C} \right) \left(\|\mathbf{u}\|^2 + 2C \sum_{t=1}^T (\ell_t^*)^2 \right) \quad \forall t \quad \|x_t^2\| \leq R^2, \quad \forall u \in \mathbb{R}^n$$

Comparison with the Perceptron Algorithm

Passive-Aggressive Algorithms

INPUT: aggressiveness parameter $C > 0$

INITIALIZE: $\mathbf{w}_1 = (0, \dots, 0)$

For $t = 1, 2, \dots$

- receive instance: $\mathbf{x}_t \in \mathbb{R}^n$
- predict: $\hat{y}_t = \text{sign}(\mathbf{w}_t \cdot \mathbf{x}_t)$
- receive correct label: $y_t \in \{-1, +1\}$
- suffer loss: $\ell_t = \max\{0, 1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)\}$
- update:

1. set:

$$\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2} \quad (\text{PA})$$

$$\tau_t = \min\left\{C, \frac{\ell_t}{\|\mathbf{x}_t\|^2}\right\} \quad (\text{PA-I})$$

$$\tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} \quad (\text{PA-II})$$

2. update: $\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t y_t \mathbf{x}_t$

Perceptron Algorithm

INPUT: learning rate $\rho > 0$

INITIALIZE: \mathbf{w}_1 random

For $t = 1, 2, \dots$

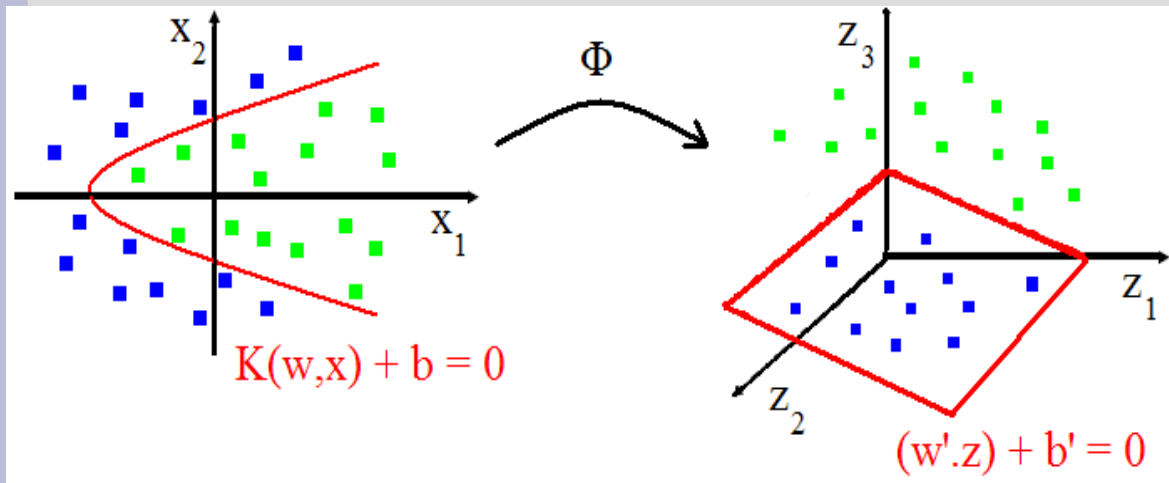
- receive instance: $\mathbf{x}_t \in \mathbb{R}^n$
- predict: $\hat{y}_t = \text{sign}(\mathbf{w}_t \cdot \mathbf{x}_t)$
- receive correct label: $y_t \in \{-1, +1\}$
- update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \rho(y_t - \hat{y}_t) \mathbf{x}_t$$

- Bounds are comparable both in separable (PA) and non-separable (PA-I, PA-II) cases

Generalization to the non-linear case: Principle

- map the data space into a feature space where the data is now linearly separable



- feature map $\Phi: \chi \rightarrow H$
- replace $(w.x)$ by Mercer Kernel $K(w,x)$ (non-linear function)

- $K(w,x)$ is the inner product of the vectors $\Phi(w)$ and $\Phi(x)$
- algorithm learns w'_t (weight vector in feature space H) and predicts $\hat{y}_t = \text{sign}(w'^T_t \cdot \Phi(x_t))$

Other problems

- Regression
- Uniclass prediction
- Multiclass problems

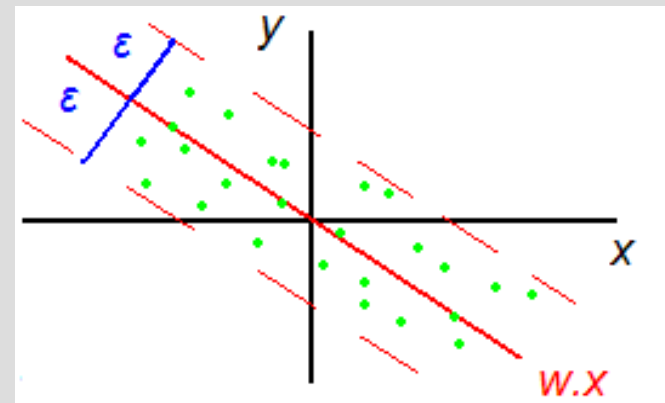
Regression

- main difference with the binary problem:

$$y \notin \{-1, 1\}, y \in \mathbb{R}$$

- instance $x_t \in \mathbb{R}^n$

→ prediction $\hat{y}_t = (w_t \cdot x_t)$



- ε -sensitive hinge loss function:

$$l_{\varepsilon}(\mathbf{w}; (\mathbf{x}, y)) = \begin{cases} 0 & |\mathbf{w} \cdot \mathbf{x} - y| \leq \varepsilon \\ |\mathbf{w} \cdot \mathbf{x} - y| - \varepsilon & \text{otherwise} \end{cases}$$

Regression: PA algorithms

- Initialization: $w_1 = (0, \dots, 0)$
- Update:
 - w_{t+1} solution of constrained optimization problem:

$$w_{t+1} = \operatorname{argmin}_{w \in \mathbb{R}^n} \frac{1}{2} \|w - w_t\|^2 \quad \text{s.t.} \quad \ell_{\mathcal{E}}(w; (\mathbf{x}_t, y_t)) = 0$$

- w_{t+1} has the form: $w_{t+1} = w_t + \operatorname{sign}(y_t - \hat{y}_t) \tau_t \mathbf{x}_t$

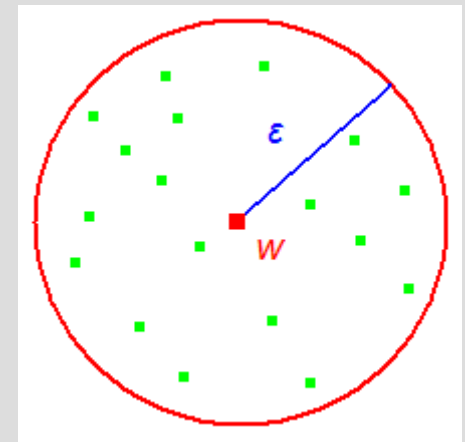
$$\tau_t = \ell_t / \|\mathbf{x}_t\|^2 \quad (PA) \quad \tau_t = \min \left\{ C, \frac{\ell_t}{\|\mathbf{x}_t\|^2} \right\} \quad (PA-I) \quad \tau_t = \frac{\ell_t}{\|\mathbf{x}_t\|^2 + \frac{1}{2C}} \quad (PA-II)$$

- Same loss bounds as for binary classification

Uniclass prediction

- Principle of a round:
 - no input x_t
 - predicts the next element of the sequence to be w_t
 - receives y_t and suffers loss:

$$l_\varepsilon(\mathbf{w}; \mathbf{y}) = \begin{cases} 0 & \|\mathbf{w} - \mathbf{y}\| \leq \varepsilon \\ \|\mathbf{w} - \mathbf{y}\| - \varepsilon & \text{otherwise} \end{cases}$$

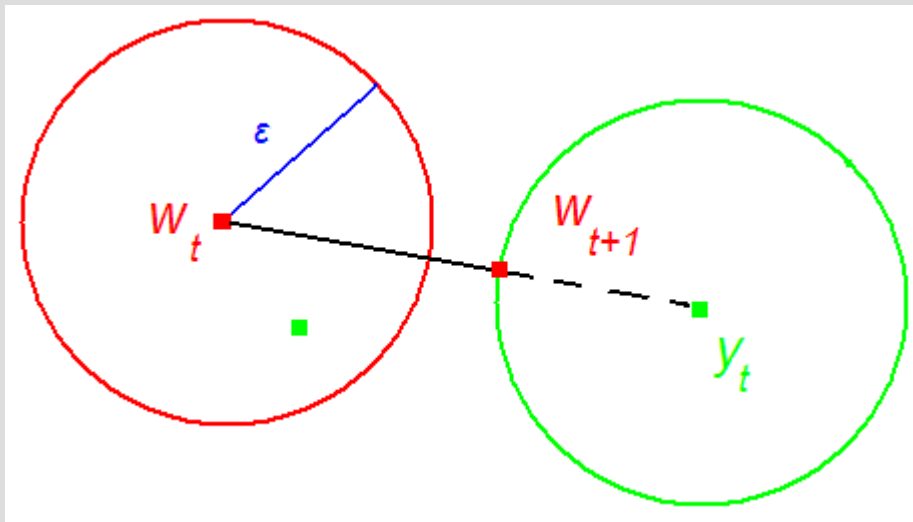


- Equivalent:
 - find the center \rightarrow elements are within a radius of ε

Uniclass prediction: PA algorithms

- Update: w_{t+1} solution of optimization problem:

$$w_{t+1} = \operatorname{argmin}_{w \in \mathbb{R}^n} \frac{1}{2} \|w - w_t\|^2 \quad \text{s.t.} \quad l_\varepsilon(w; y_t) = 0$$



w_{t+1} has the form:

$$w_{t+1} = w_t + \tau_t \frac{y_t - w_t}{\|y_t - w_t\|}$$

$$\tau_t = l_t \quad (\text{PA})$$

$$\tau_t = \min \{ C, l_t \} \quad (\text{PA-I})$$

$$\tau_t = \frac{l_t}{1 + \frac{1}{2C}} \quad (\text{PA-II})$$

Multiclass multilabel classification

- Principle:
 - set of all possible labels $Y = \{1, \dots, k\}$
 - receives instance x_t (associated with relevant labels)
 - outputs a score for each of the k labels
 - *prediction vector* $\in \mathbb{R}^k$
 - receives the set of "relevant" labels Y_t for x_t
 - "relevant" must be ranked higher than "irrelevant"
 - updates the prediction mechanism

Multiclass multilabel: Problem settings

- feature vector: $\Phi(x, y) = (\Phi_1(x, y), \dots, \Phi_d(x, y))$
(set of features: Φ_1, \dots, Φ_d)

- Prediction vector:

$$\left((w_t \cdot \Phi(\mathbf{x}_t, 1)), \dots, (w_t \cdot \Phi(\mathbf{x}_t, k)) \right) \quad w_t \in \mathbb{R}^d$$

- Margin of the exemple (x_t, Y_t) :

$$\gamma(w_t; (\mathbf{x}_t, Y_t)) = \min_{r \in Y_t} w_t \cdot \Phi(\mathbf{x}_t, r) - \max_{s \notin Y_t} w_t \cdot \Phi(\mathbf{x}_t, s)$$

Multiclass multilabel: Problem settings (2)

- Margin: difference between
 - score of the lowest ranked relevant label
 - score of the highest ranked irrelevant label

$$r_t = \operatorname{argmin}_{r \in Y_t} \mathbf{w}_t \cdot \Phi(\mathbf{x}_t, r) \quad \text{and} \quad s_t = \operatorname{argmax}_{s \notin Y_t} \mathbf{w}_t \cdot \Phi(\mathbf{x}_t, s)$$

- Hinge-loss function:

$$\ell_{\text{MC}}(\mathbf{w}; (\mathbf{x}, Y)) = \begin{cases} 0 & \gamma(\mathbf{w}; (\mathbf{x}, Y)) \geq 1 \\ 1 - \gamma(\mathbf{w}; (\mathbf{x}, Y)) & \text{otherwise} \end{cases}$$

Multiclass multilabel: PA algorithms

- Equivalence: $\ell_{MC}(\mathbf{w}_t; (\mathbf{x}_t, Y_t)) = \ell(\mathbf{w}_t; (\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t), +1))$

- \mathbf{w}_{t+1} solution of optimization problem:

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \|\mathbf{w} - \mathbf{w}_t\|^2 \quad \text{s.t.} \quad \mathbf{w} \cdot (\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t)) \geq 1$$

- \mathbf{w}_{t+1} has the form: $\mathbf{w}_{t+1} = \mathbf{w}_t + \tau_t (\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t))$

$$\tau_t = \frac{\ell_t}{\|\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t)\|^2} \quad (\text{PA}) \quad \tau_t = \min \left\{ C, \frac{\ell_t}{\|\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t)\|^2} \right\} \quad (\text{PA-I})$$

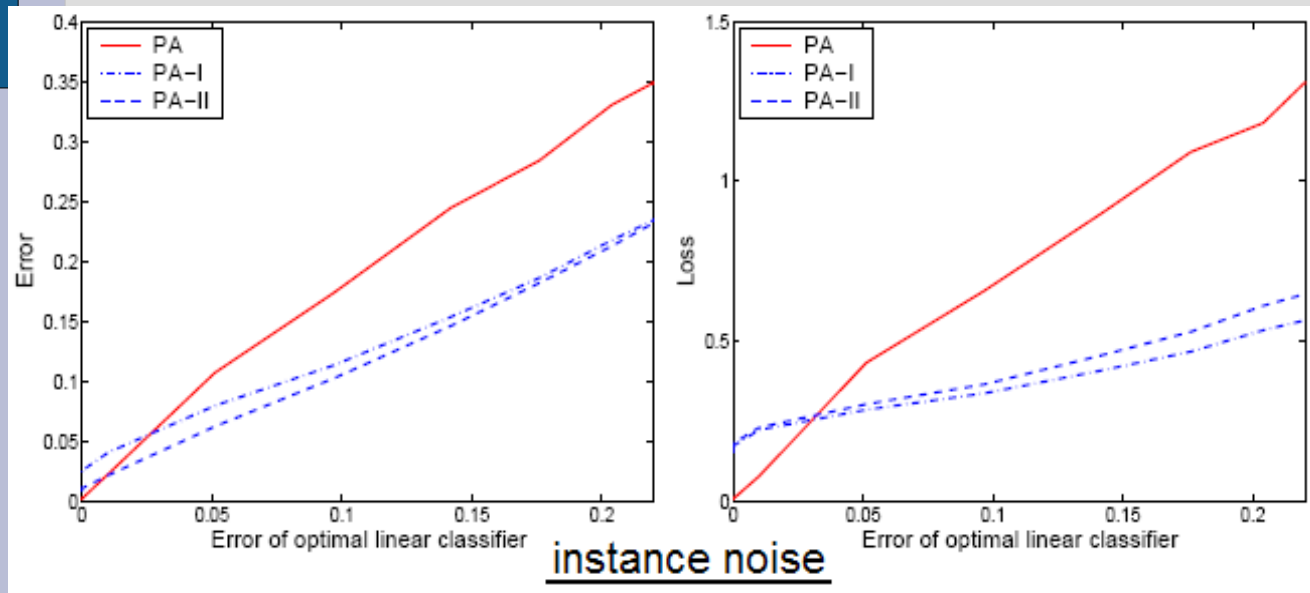
$$\tau_t = \frac{\ell_t}{\|\Phi(\mathbf{x}_t, r_t) - \Phi(\mathbf{x}_t, s_t)\|^2 + \frac{1}{2C}} \quad (\text{PA-II})$$

Experiments

1. Robustness to noise
2. Effect of the aggressiveness parameter C
3. Multiclass problems:
comparison with other online algorithms

Experiment 1: Robustness to noise

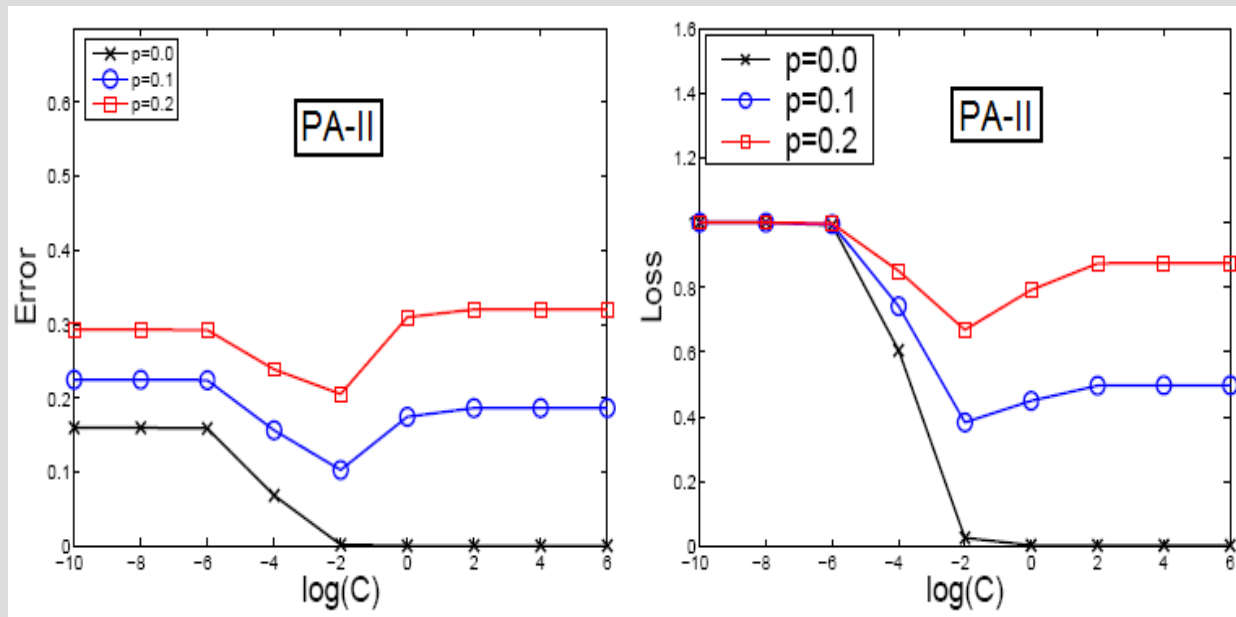
- Binary classification, 4000 generated examples (results averaged on 10 repetitions)



- Instance noise
- label noise
- Find optimal fixed linear classifier (brute force)
- $C = 0.001$

- Low noise level: 3 make similar number of errors
- High noise level: PA-I and PA-II outperform PA

Experiment 2: Effect of C



C " aggressiveness
parameter "

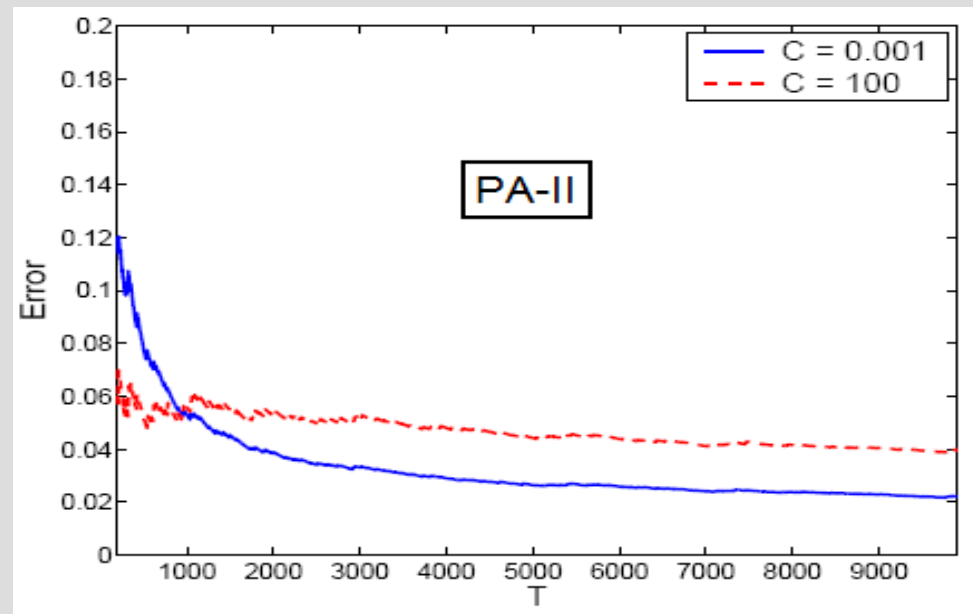
- Results meet the theoretic loss bounds

$$\sum_{t=1}^T \ell_t^2 \leq \left(R^2 + \frac{1}{2C} \right) \left(\|\mathbf{u}\|^2 + 2C \sum_{t=1}^T (\ell_t^*)^2 \right)$$

- Rule: when there is noise in data, C should be small

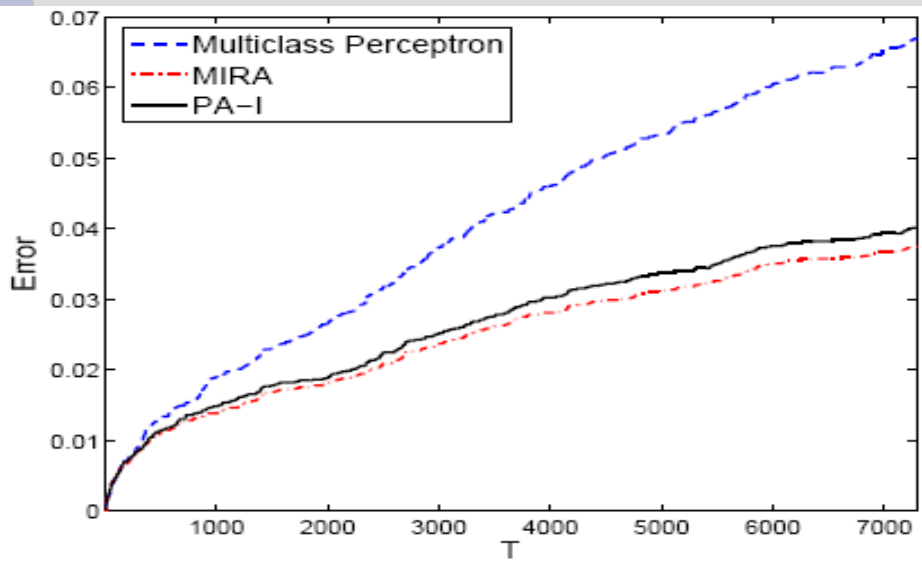
Experiment 2: Effect of C (2)

- Evolution of error rate with the number of examples



Experiment 3: Multiclass problems

- Use standard multiclass datasets: USPS, MNIST
- Comparison of the multiclass PA algorithms with:
 - multiclass versions of the Perceptron algorithm
 - MIRA (Margin Infused Relaxed Algorithm)



- PA-I and MIRA comparable
- but MIRA solves a complex optimisation problem for each update
≠ PA: simple expression

Conclusion

- Further research:
 - extension to other problems
 - conversion to batch algorithms
 - PA with bounded memory constraints
(memory requirements imposed when using Mercer Kernels)

References

- Crammer, Koby. Dekel, Ofer. Keshet, Joseph. Shalev-Shwartz, Shai. Singer, Yoram. " Online Passive-Aggressive Algorithms ". Jerusalem, 2006
<<http://jmlr.csail.mit.edu/papers/volume7/crammer06a/crammer06a.pdf>>
(20 Oct. 2007)
- Schiele, Bernt. "Maschinelles Lernen - Statistische Verfahren". Darmstadt: Technische Universität Darmstadt, 18. Mai 2007
<<http://www.mis.informatik.tu-darmstadt.de/Education/Courses/ml/slides/ml-2007-0518-svm2-v1.pdf>>
(16 Nov. 2007)
- Rojas, Raul. "Perceptron Learning". Neural Networks - A Systematic Introduction. Springer-Verlag, Berlin, 1996
<<http://page.mi.fu-berlin.de/rojas/neural/chapter/K4.pdf>>
(22 Nov. 2007)
- Rodriguez, Carlos C. "The Kernel Trick". October 25, 2004
<omega.albany.edu:8008/machine-learning-dir/notes-dir/ker1/ker1.pdf>
(26 Nov. 2007)

Thank you for your attention

Questions?