# Online Machine Learning

Simone Filice
filice.simone@gmail.com

University of Roma Tor Vergata
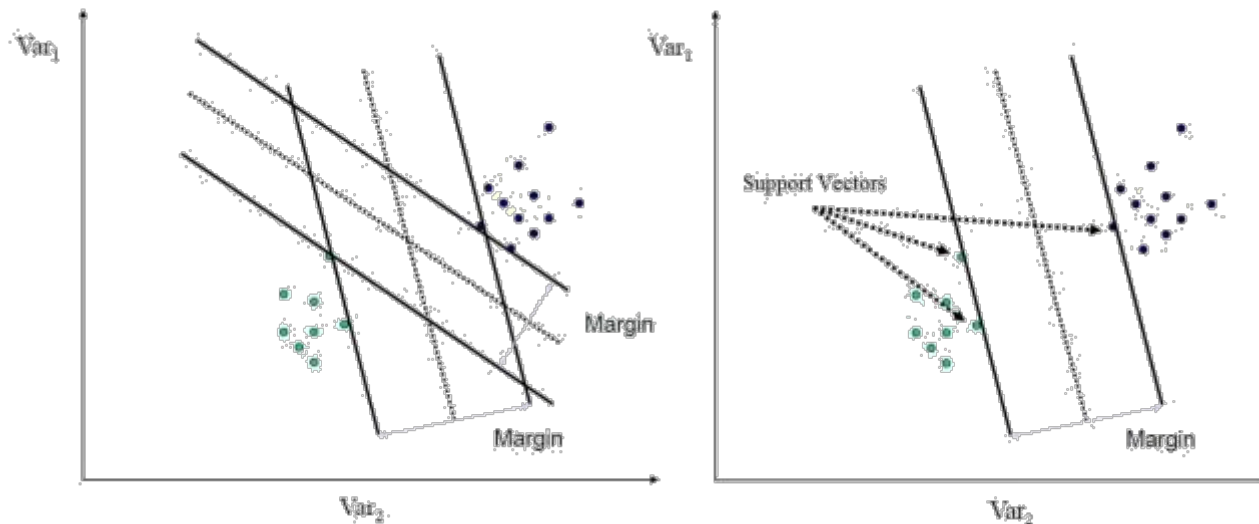
Web Mining e Retrieval 2013/2014

# Motivations

- Common ML algorithms simultaneously exploit a whole dataset. This process, referred as *batch learning*, is not practical when:
  - New data naturally arise over the time: exploiting new data means building from scratch a new model → usually not feasible!
  - The dataset is too large to be efficiently exploited: memory and computational problems!
  - The concept we need to learn changes over the time: batch learning provide a static solution that will surely degrade as time goes by

# Online Machine Learning

- Incremental Learning Paradigm:
  - Every time a new example is available, the learned hypothesis is updated

- Inherent Appealing Characteristics:
  - The model does not need to be re-generated from scratch when new data is available
  - Capability of tracking a Shifting Concept
  - Faster training process if compared to batch learners (e.g. SVM)

# Perceptron

- Perceptron is a simple discriminative classifier
    - Instances are feature vectors $x' \in \mathbb{R}^d$ with label $y \in [-1, +1]$
    - Classification function is an hyperplane in $\mathbb{R}^d$ : $f(x') = w' \cdot x' + b$



- Compact notation: $w = \{b, w'_1, w'_2, ..., w'_d\}$, $x = \{1, x'_1, x'_2, ..., x'_d\}$

# Batch Perceptron

□ IDEA : adjust the hyperplane until no training errors are done (input data must be linearly separable)

□ Batch perceptron learning procedure:

```
Start with w₁ = 0
do
        errors=false
        For all t=1…T
            Receive a new sample xₜ
            Compute y = wₜ · xₜ
            if y · yₜ < βₜ then wₜ₊₁ = γₜwₜ + αₜyₜxₜ  with αₜ > 0
                errors=true
             else
                    wₜ₊₁ = wₜ
while(errors)
return wₜ₊₁
```

Start with $\boldsymbol{w}_1 = 0$

do

      errors=false

      For all t=1…T

         Receive a new sample $\boldsymbol{x}_t$

         Compute $y = \boldsymbol{w}_t \cdot \boldsymbol{x}_t$

         if $y \cdot y_t < \beta_t$ then $\boldsymbol{w}_{t+1} = \gamma_t \boldsymbol{w}_t + \alpha_t y_t \boldsymbol{x}_t$ with $\alpha_t > 0$

           errors=true

          else

              $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t$

while(errors)

return $\boldsymbol{w}_{T+1}$

# Online Learning Perceptron

- IDEA : adjust the hyperplane after each classification ($\boldsymbol{w}_t =$ weight vector at time *t*) and never stop learning

- Online perceptron learning procedure:

```
Start with w₁ = 0
For all t=1…
        Receive a new sample xₜ
        Compute y = wₜ · xₜ
        Receive a feedback yₜ
        if y · yₜ < βₜ then wₜ₊₁ = γₜwₜ + αₜyₜxₜ   with αₜ > 0
        else   wₜ₊₁ = wₜ
endfor
```

# Shifting Perceptron

- IDEA: weak dependance from the past in order to obtain a tracking ability

- Shifting Perceptron learning procedure (Cavallanti et al 2006):

```
Start with w₁ = 0 , k=0
For all t=1…
        Receive a new sample xₜ
        Compute y = sign(wₜ · xₜ)
        Receive a feedback yₜ
        if y ≠ yₜ then
```
$$\lambda_k = \frac{\lambda}{\lambda+k} \quad \text{with} \quad \lambda > 0$$
$$\boldsymbol{w}_{t+1} = (1 - \lambda_k)\boldsymbol{w}_t + \lambda_k y_t \boldsymbol{x}_t$$
```
                k=k+1
        else    w_{t+1} = w_t
endfor
```

# Online Linear Passive Aggressive

☐ IDEA: Every time a new example $\langle x_t, y_t \rangle$ is available the current classification function is modified as less as possible to correctly classify the new example

☐ Passive Aggressive learning procedure (Crammer et al 2006):

Start with $\boldsymbol{w}_1 = 0$ , k=0

For all t=1…

    Receive a new sample $\boldsymbol{x_t}$

    Compute $y = \mathrm{sign}(\boldsymbol{w}_t \cdot \boldsymbol{x_t})$

    Receive a feedback $y_t$

    Measure a classification loss (divergence between $y_t$ and $y$)
    Modify the model to get zero loss, preserving what was
    learned from previous examples

# Online Linear Passive Aggressive

☐ Loss measure:

Hinge loss: $l(w; (x_t, y_t)) = \max(0; 1 - y_t(w \cdot x_t))$

☐ Model variation:

$$\|w_{t+1} - w_t\|^2$$

☐ Passive Aggressive Optimization Problem:

$$w_{t+1} = argmin_w \frac{1}{2}\|w - w_t\|^2 \text{ such that } l(w; (x_t, y_t)) = 0$$

☐ Closed form solution:

$$w_{t+1} = w_t + \tau_t y_t x_t \text{ where } \tau_t = \frac{l(w_t; (x_t, y_t))}{\|x_t\|^2}$$

# Online Linear Passive Aggressive

- The previous formulation is a hard margin version that has a problem:
  - a single outlier could produce a high hyperplane shifting, making the model forget the previous learning
- Soft version solution:
  - control the algorithm aggressiveness through a parameter C

- PA-I formulation:

$$\boldsymbol{w}_{t+1} = argmin_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{w} - \boldsymbol{w}_t\|^2 + C\xi \quad \text{s.t. } l\big(\boldsymbol{w}; (\boldsymbol{x}_t, y_t)\big) \leq \xi \text{ with } \xi \geq 0$$

$$\Longrightarrow \boldsymbol{w}_{t+1} = \boldsymbol{w}_t + \tau_t y_t \boldsymbol{x}_t \quad \text{where } \tau_t = \min\left\{C; \frac{l(\boldsymbol{w}_t;(\boldsymbol{x}_t,y_t))}{\|\boldsymbol{x}_t\|^2}\right\}$$

- PA-II model:

$$\boldsymbol{w}_{t+1} = argmin_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{w} - \boldsymbol{w}_t\|^2 + C\xi^2 \quad \text{s.t. } l\big(\boldsymbol{w}; (\boldsymbol{x}_t, y_t)\big) \leq \xi \text{ with } \xi \geq 0$$

$$\Longrightarrow \boldsymbol{w}_{t+1} = \boldsymbol{w}_t + \tau_t y_t \boldsymbol{x}_t \quad \text{where } \tau_t = \frac{l(\boldsymbol{w}_t;(\boldsymbol{x}_t,y_t))}{\|\boldsymbol{x}_t\|^2 + \frac{1}{2}C}$$

# Data Separability

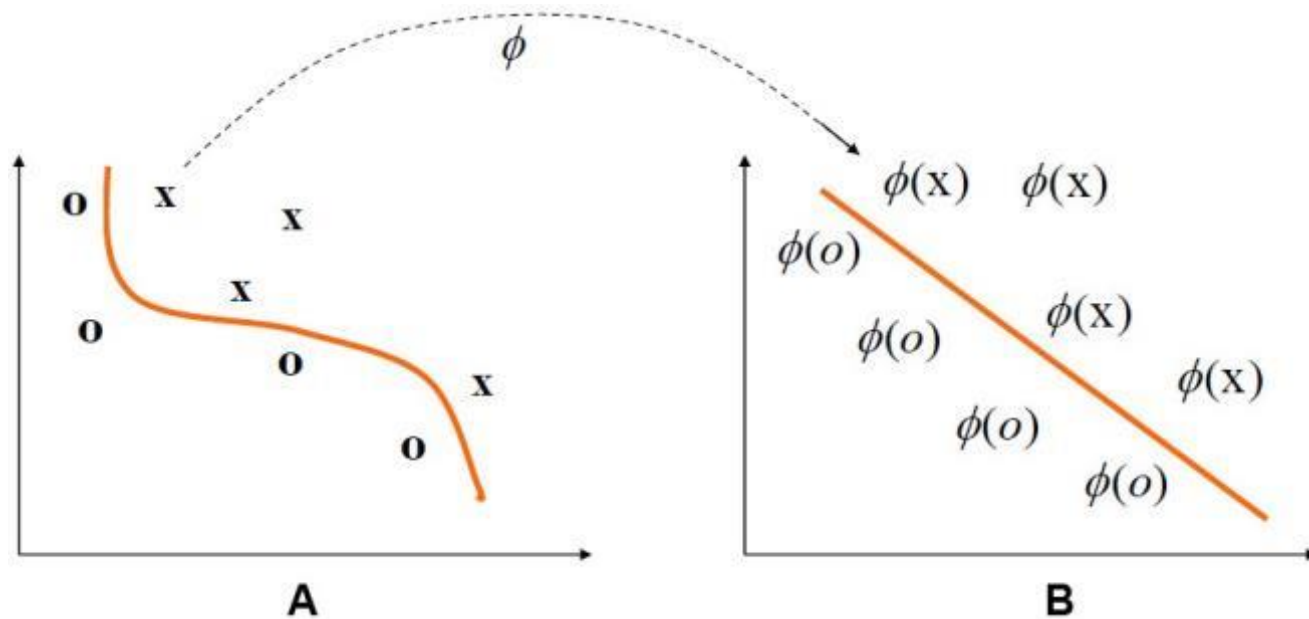□ Training data could not be separable

□ Possible solutions:

    □ Use a more complex classification function → Risk of overfitting!

    □ Define a new set of feature that makes the problem linearly separable

    □ Project the current examples in a space in which they are separable...

# Kernel Methods

- Training data can be projected in a space in which they are more easily separable



- Kernel Trick: any kernel function K performs the dot product in the kernel space without explicitly project the input vectors in that space

- Structured data (tree, graph, high order tensor…) can be exploited

# Kernelized Passive Aggressive

☐ In kernelized Online Learning algorithms a new support vector is added every time a misclassification occurs

| LINEAR VERSION | KERNELIZED VERSION |
|---|---|
| **Classification function** | |
| $f_t(\boldsymbol{x}) = \boldsymbol{w}_t{}^T \boldsymbol{x}$ | $f_t(x) = \sum_{i \in S} \alpha_i k(x, x_i)$ |
| **Optimization Problem (PA-I)** | |
| $\boldsymbol{w}_{t+1} = argmin_{\boldsymbol{w}} \frac{1}{2} \|\boldsymbol{w} - \boldsymbol{w}_t\|^2 + C\xi$ <br><br> Such that $1 - y_t f_t(\boldsymbol{x}_t) \leq \xi, \xi \geq 0$ | $f_{t+1}(x) = \text{argmin}_{\text{f}} \frac{1}{2} \|f(x) - f_t(x)\|^2{}_{\mathcal{H}} + C\xi$ <br><br> Such that $1 - y_t f_t(x_t) \leq \xi, \xi \geq 0$ |
| **Closed form solution** | |
| $\boldsymbol{w}_{t+1} = \boldsymbol{w}_t + \tau_t y_t \boldsymbol{x}_t$ <br><br> where $\tau_t = \min\left\{C; \frac{\max(0, 1 - y_t f_t(\boldsymbol{x}_t))}{\|\boldsymbol{x}_t\|^2}\right\}$ | $f_{t+1}(x) = \text{f}_t(x) + \alpha_t k(x, x_t)$ <br><br> where $\alpha_t = y_t \cdot \min\left\{C; \frac{\max(0, 1 - y_t f_t(x_t))}{\|x_t\|^2{}_{\mathcal{H}}}\right\}$ |

# Linear Vs Kernel Based Learning

| LINEAR VERSION | KERNELIZED VERSION |
|---|---|
| **Classification function** | |
| explicit hyperlplane in the original space<br>☹ Only linear functions can be learnt | implicit hyperplane in the RKHS<br>☺ Non linear functions can be learnt |
| **Example form** | |
| ☹ Only feature vectors can be exploited | ☺ Structured representations can be exploited |
| **Computational complexity** | |
| ☺ A classification is a single dot product | ☹ A classification involves $|S|$ kernel computations |
| **Memory usage** | |
| ☺ Only a the explicit hyperplane must be stored | ☹ All the support vectors and their weights must be stored |

# Learning on a Budget

- In kernelized online learning algorithm the set of support vectors can grow without limits

- Possible solution: Limit the number of support vector, defining a budget $B$

- This solution has the following advantages:

  - The memory occupation is upperbounded by $B$ support vectors

  - Each classification needs at most $B$ kernel computations

  - In shifting concept tasks, budget algorithms can outperform non-budget counterparts because they are faster in adapting

# Limit the number of Support Vectors

- In order to respect the budget B, different policies can be formulated:

  - Stop learning when budget is exceeded: *Stoptron*

  - Delete a random support vector: *Randomized Perceptron*

  - Delete the more redundant support vector: *Fixed Budget Conscious Perceptron*

  - Delete the oldest support vector: *Least recent Budget Perceptron* and *Forgetron*

  - Modify the Support Vectors weights in order to adapt the classification hypothesis to the new sample: *Projectron*

  - *Online Passive-Aggressive on a Budget*

# Stoptron

- Baseline of the online learning on a budget algorithms: Fix a budget $B$ and stop learning when the number of support vectors is equal to $B$

- Stoptron algorithm (Orabona et al 2008):

```
Start with S = Ø
For all t=1…
        Receive a new sample x_t
        Compute   y = Σ_{i∈S} α_i y_i K(x_i, x_t )
        Receive a feedback y_t
        if   yy_t < β   and  |S| < B  then
                S = S ∪ {t}
                α_t = 1
        endif
endfor
```

# Randomized Perceptron

☐ <u>Simplest deleting policy:</u> when the budget $B$ is exceeded remove a random support vector

☐ Randomized Perceptron algorithm (Cavallanti et al 2007):

```
Start with S = ∅
For all t=1…
        Receive a new sample x_t
        Compute   y = ∑_{i∈S} α_i y_i K(x_i, x_t )
        Receive a feedback y_t
        if   yy_t < β
            if |S| = B
                select randomly s ∈ S, S = S \ {s}
            endif
            S = S ∪ {t}  α_t = 1
        endif
endfor
```

# Forgetron

- <u>Deleting policy:</u> Every time a new support vector is added, the weights of the others are reduced. Thus SVs lose weight with aging and removing the older SV should assure a minimum impact to the classification function.

- Forgetron algorithm (Dekel et al 2008):

```
Start with S = ∅
For all t=1…
        Receive a new sample 𝑥ₜ
        Compute   𝑦 = ∑ᵢ∈ₛ αᵢyᵢK(𝒙ᵢ, 𝒙ₜ)
        Receive a feedback yₜ
        if   yyₜ < β
            if |S| = B
                  S=S\min{S} //the oldest Support vector is removed
            endif
            S = S ∪ {t}  αₜ = 1,  αᵢ = φₜαᵢ  ∀i ∈ S \ {t} //adding a new Sv and shrinking
        endif
endfor
```

# Summary

- Online learning methods can:
  - Incrementally learn from new samples
  - Dinamically adapt to problem variations
  - Reduce the computational cost of building a new model

- Online learning methods can be used with kernels but they suffer from the "*curse of kernelization*":
  - The number of support vectors can grow without bounds

- Several number of budgeted solutions have been proposed