



Summer school
Online Learning
28.06.2015-02.07.2015
Copenhagen University

Participants



Participants & Lecturer

- ▶ Around 60 participants
- ▶ Lecturer
 - ▶ Shai Shalev-Shwartz, Hebrew University
 - ▶ Peter Auer, Leoben University
 - ▶ Nicola Cesa-Bianchi, Miland University
 - ▶ Csaba Szepesvari, Alberta University
 - ▶ Yevgeny Seldin, Copenhagen University



Topics

- ▶ Basics of online learning
 - ▶ Online convex optimization
 - ▶ Bandits (stochastic, adversarial....)
 - ▶ Online reinforcement learning
 - ▶ Space of online learning problems
-
- ▶ Theory only, proof of bounds

General online learning problem

- ▶ for $t = 1, 2, \dots, T$
 - ▶ receive question $x_t \in X$
 - ▶ predict $p_t \in D$
 - ▶ receive true answer $y_t \in Y$
 - ▶ suffer loss $l(p_t, y_t)$
 - ▶ (update model)
- ▶ Given a fixed hypothesis class H (e.g. halfspaces $h_t = \text{sign}(\langle w_t | x_t \rangle)$), $\forall t$ find $h_t \in H$, $h_t: X \rightarrow Y$ s.t. $\sum_{t=1}^T l(p_t, y_t)$ is minimized ($p_t = h_t(x_t)$).
- ▶ Loss functions
 - ▶ $l(p_t, y_t) = |p_t - y_t|$ "0 - 1 loss" or "absolut loss"
 - ▶ $l(p_t, y_t) = (p_t - y_t)^2$ "quadratic loss"

2 alternative restrictions/goals

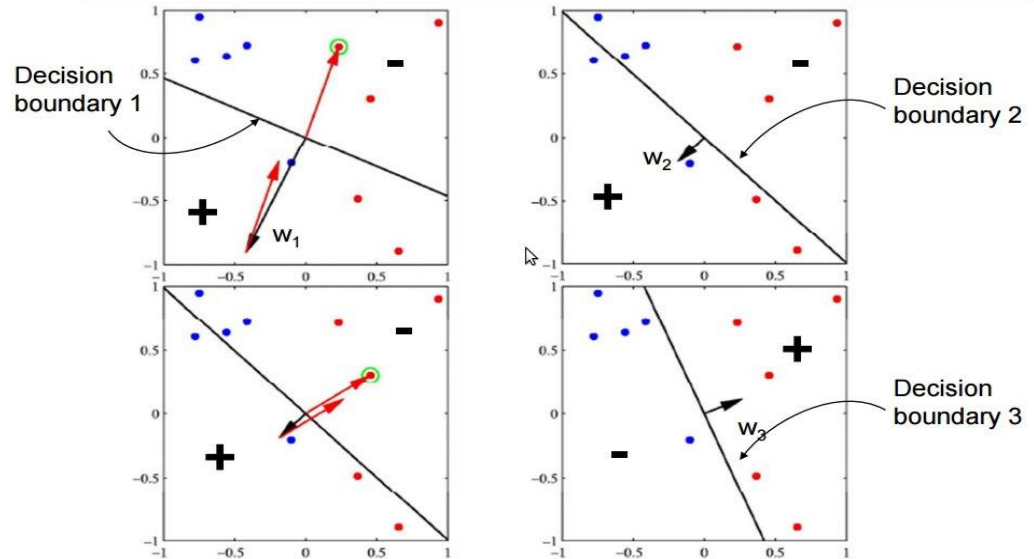
- ▶ No assumption about sequence (deterministic, stochastic, adversarial)
- ▶ Realizability
 - ▶ $\exists \hat{h}: X \rightarrow Y$ s.t. $y_t = \hat{h}(x_t) \forall t, \hat{h} \in H$
 - ▶ Goal : Find algorithm with minimal mistake bound (sublinear with T)
- ▶ No realizability
 - ▶ Goal: Find algorithm with minimal regret $R(\hat{h})$ compared to the best fixed predictor $\hat{h} \in H$
 - ▶ $R(\hat{h})_T = \sum_{t=1}^T l(p_t, y_t) - \sum_{t=1}^T l(\hat{h}(p_t), y_t)$
- ▶ Given an algorithm-> find & proof the corresponding mistake-bound/regret-bound

Perceptron for realizability case

- ▶ Separate data points for binary classification
- ▶ $Y = \{-1, 1\}$, hypothesis class $H =$ all halfspaces in \mathbb{R}^n
- ▶ $l(w_t) = \max(1 - y_t \langle w_t | x_t \rangle, 0)$, similar to "hinge - loss"

When an error is made, moves the weight in a direction that corrects the error

- ▶ initialize: $w_1 = 0$
- ▶ for $t = 1, 2, \dots, T$
 - ▶ receive x_t
 - ▶ predict $p_t = \text{sign}(\langle w_t | x_t \rangle)$
 - ▶ if $y_t \langle w_t | x_t \rangle \leq 0$
 - ▶ $w_{t+1} = w_t + y_t x_t$
 - ▶ else $w_{t+1} = w_t + y_t x_t$



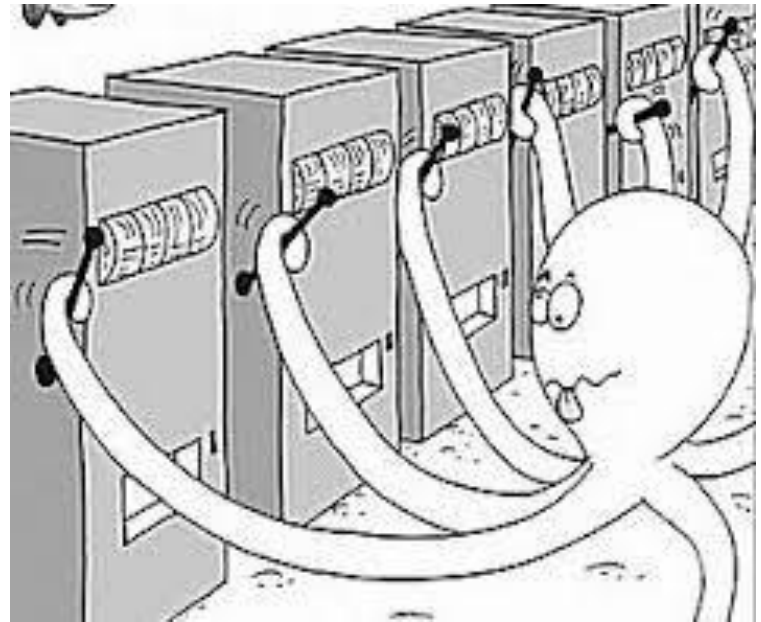
Red points belong to the positive class,
blue points belong to the negative class

- ▶ Suppose that $\|x_t\| \leq D$ and $\exists \hat{w} (\|\hat{w}\|_2 = 1) \text{ s.t. } y_t \langle \hat{w} | x_t \rangle \geq \gamma \forall t$

$$\rightarrow M \leq \left(\frac{D}{\gamma}\right)^2$$

Multi-armed Bandits

- ▶ for $t = 1, 2, \dots, T$
 - ▶ play action $y_t \in Y = \{1, \dots, K\}$
 - ▶ receive reward $r_t(y_t) \in [0, 1]$
- ▶ Only reward of played action is seen
- ▶ Goal: maximize reward!
 - ▶ minimize regret $R = T\hat{\mu} - \sum_{t=1}^T r_t$
- ▶ Exploitation vs. Exploration
- ▶ Old, but very popular nowadays....Why?!
 - Google uses it



Multi-armed Bandits

- ▶ Various versions
 - ▶ Stochastic stationary/non-stationary
 - ▶ Adversarial, Contextual
 - ▶ Graph-based
- ▶ Applications
 - ▶ Web-searches (contextual, max adv. income of Google, Bing etc.)
 - ▶ Clinical trials (minimize patient losses)
 - ▶ Adaptive routing (minimize delays)

Stochastic, stationary Bandits

- ▶ Rewards for actions are generated by stationary distributions
- ▶ $R(T) = T\hat{\mu} - \sum_{j=1}^K \mu_j * T_j$, T_j is the number action j was played
- ▶ Estimate μ_k accurately to minimize the regret

Upper Confidence Bound (UCB) algorithm

- ▶ Calculate confidence bounds for each action
- ▶ Chernoff-Hoeffding bound:
 - ▶ Let X_1, X_2, \dots, X_K independent random variables in the range $[0,1]$ with $\mu = \text{Exp}(X) = \frac{1}{K} \sum_{j=1}^K \mu_j \rightarrow P\left(\frac{1}{T} \sum_{t=1}^T X_i \geq \mu + a\right) \leq e^{-2a^2 T}$
- ▶ $a = \sqrt{\frac{2 \log(T)}{T_j}} \rightarrow P\left(\frac{1}{T} \sum_{t=1}^T X_i \geq \mu + a\right) \leq T^{-4}$, converges quickly to 0
- ▶ Choose the action with highest upper confidence bound
$$\max \mu_j + \sqrt{\frac{2 \log(T)}{T_j}}$$
- ▶ Balances exploitation vs. exploration
- ▶ $R_{\text{UCB}} \leq \sum_{j=1}^K \frac{2 \log(T)}{\Delta_j} + \left(1 + \frac{\pi^2}{3}\right) \Delta_j, \quad \Delta_j = \hat{\mu} - \mu_j$

Further readings

- ▶ Got interested ?
 - ▶ Shai Shalev-Shwartz “Online Learning and Online Convex Optimization”
 - ▶ Sebastien Bubeck, Nicolò Cesa-Bianchi „Regret Analysis of Stochastic and Non-stochastic Multi-armed Bandit Problems“
 - ▶ Nicolò Cesa-Bianchi “Prediction, learning, and games”