

The Online Approach to Machine Learning

Nicolò Cesa-Bianchi

Università degli Studi di Milano



Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel
- 4 The joy of convex
- 5 The joy of convex (without the gradient)



Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel
- 4 The joy of convex
- 5 The joy of convex (without the gradient)





Classification/regression tasks

- Predictive models h mapping data instances X to labels Y (e.g., binary classifier)
- Training data $S_T = ((X_1, Y_1), \dots, (X_T, Y_T))$ (e.g., email messages with spam vs. nonspam annotations)
- Learning algorithm A (e.g., Support Vector Machine) maps training data S_T to model $h = A(S_T)$

Evaluate the **risk** of the trained model h with respect to a given **loss function**



Two notions of risk

View data as a statistical sample: **statistical risk**

$$\mathbb{E} \left[\text{loss} \left(\underbrace{A(S_T)}_{\text{trained model}}, \underbrace{(X, Y)}_{\text{test example}} \right) \right]$$

Training set $S_T = ((X_1, Y_1), \dots, (X_T, Y_T))$ and test example (X, Y) drawn i.i.d. from the same unknown and fixed distribution

View data as an arbitrary sequence: **sequential risk**

$$\sum_{t=1}^T \text{loss} \left(\underbrace{A(S_{t-1})}_{\text{trained model}}, \underbrace{(X_t, Y_t)}_{\text{test example}} \right)$$

Sequence of models trained on growing prefixes $S_t = ((X_1, Y_1), \dots, (X_t, Y_t))$ of the data sequence

Regrets, I had a few

Learning algorithm A maps datasets to models in a given class \mathcal{H}

Variance error in statistical learning

$$\mathbb{E} \left[\text{loss}(A(S_T), (X, Y)) \right] - \inf_{h \in \mathcal{H}} \mathbb{E} \left[\text{loss}(h, (X, Y)) \right]$$

compare to expected loss of best model in the class

Regret in online learning

$$\sum_{t=1}^T \text{loss}(A(S_{t-1}), (X_t, Y_t)) - \inf_{h \in \mathcal{H}} \sum_{t=1}^T \text{loss}(h, (X_t, Y_t))$$

compare to cumulative loss of best model in the class



Incremental model update

A natural blueprint for online learning algorithms

For $t = 1, 2, \dots$

- 1 Apply current model h_{t-1} to next data element (X_t, Y_t)
- 2 Update current model: $h_{t-1} \rightarrow h_t \in \mathcal{H}$

Goal: control regret

$$\sum_{t=1}^T \text{loss}(h_{t-1}, (X_t, Y_t)) - \inf_{h \in \mathcal{H}} \sum_{t=1}^T \text{loss}(h, (X_t, Y_t))$$

View this as a **repeated game** between a player generating predictors $h_t \in \mathcal{H}$ and an opponent generating data (X_t, Y_t)



Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel
- 4 The joy of convex
- 5 The joy of convex (without the gradient)



Theory of repeated games



James Hannan
(1922–2010)



David Blackwell
(1919–2010)

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)

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



Prediction with expert advice



Volodya Vovk



Manfred Warmuth

	$t = 1$	$t = 2$	\dots
1	$l_1(1)$	$l_2(1)$	\dots
2	$l_1(2)$	$l_2(2)$	\dots
\vdots	\vdots	\vdots	\ddots
N	$l_1(N)$	$l_2(N)$	

Opponent's moves y_1, y_2, \dots define a **sequential prediction problem** with a **time-varying loss function** $l(i_t, y_t) = l_t(i_t)$



Playing the experts game

N actions



For $t = 1, 2, \dots$

- Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)



Playing the experts game

N actions



For $t = 1, 2, \dots$

- 1 Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)
- 2 Player picks an action I_t (possibly using randomization) and incurs loss $\ell_t(I_t)$



Playing the experts game

N actions



For $t = 1, 2, \dots$

- 1 Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)
- 2 Player picks an action I_t (possibly using randomization) and incurs loss $\ell_t(I_t)$
- 3 Player gets **feedback information**: $\ell_t(1), \dots, \ell_t(N)$



Oblivious opponents

Losses $\ell_t(1), \dots, \ell_t(N)$ for all $t = 1, 2, \dots$ are fixed beforehand, and unknown to the (randomized) player

Oblivious regret minimization

$$R_T \stackrel{\text{def}}{=} \mathbb{E} \left[\sum_{t=1}^T \ell_t(I_t) \right] - \min_{i=1, \dots, N} \sum_{t=1}^T \ell_t(i) \stackrel{\text{want}}{=} o(T)$$



Lower bound using random losses

- $\ell_t(i) \rightarrow L_t(i) \in \{0, 1\}$ independent random coin flip

- For any player strategy $\mathbb{E} \left[\sum_{t=1}^T L_t(I_t) \right] = \frac{T}{2}$

- Then the expected regret is

$$\mathbb{E} \left[\max_{i=1, \dots, N} \sum_{t=1}^T \left(\frac{1}{2} - L_t(i) \right) \right] = (1 - o(1)) \sqrt{\frac{T \ln N}{2}}$$



Exponentially weighted forecaster

At time t pick action $I_t = i$ with probability proportional to

$$\exp\left(-\eta \sum_{s=1}^{t-1} \ell_s(i)\right)$$

the sum at the exponent is the **total loss** of action i up to now

Regret bound

[Experts' paper, 1997]

- If $\eta = \sqrt{(\ln N)/(8T)}$ then $R_T \leq \sqrt{\frac{T \ln N}{2}}$
- Matching lower bound including constants
- Dynamic choice $\eta_t = \sqrt{(\ln N)/(8t)}$ only loses small constants

The bandit problem: playing an unknown game

N actions



For $t = 1, 2, \dots$

- ① Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)



The bandit problem: playing an unknown game

N actions



For $t = 1, 2, \dots$

- 1 Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)
- 2 Player picks an action I_t (possibly using randomization) and incurs loss $\ell_t(I_t)$



The bandit problem: playing an unknown game

N actions



For $t = 1, 2, \dots$

- 1 Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)
- 2 Player picks an action I_t (possibly using randomization) and incurs loss $\ell_t(I_t)$
- 3 Player gets **feedback information**: Only $\ell_t(I_t)$ is revealed



The bandit problem: playing an unknown game

N actions



For $t = 1, 2, \dots$

- 1 Loss $\ell_t(i) \in [0, 1]$ is assigned to every action $i = 1, \dots, N$ (hidden from the player)
- 2 Player picks an action I_t (possibly using randomization) and incurs loss $\ell_t(I_t)$
- 3 Player gets **feedback information**: Only $\ell_t(I_t)$ is revealed

Many applications

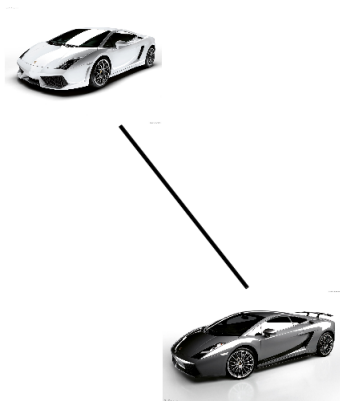
Ad placement, dynamic content adaptation, routing, online auctions

Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel**
- 4 The joy of convex
- 5 The joy of convex (without the gradient)



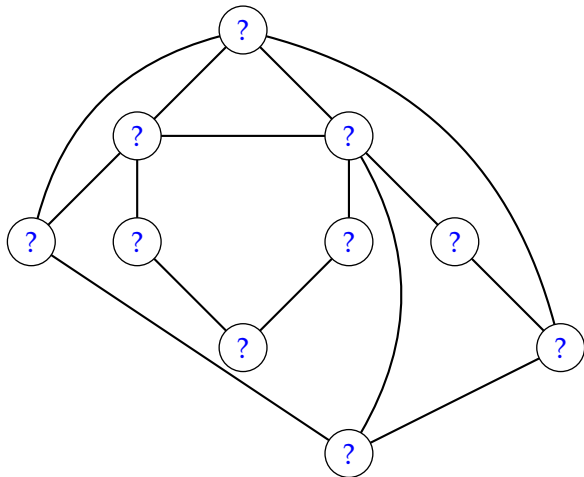
Undirected



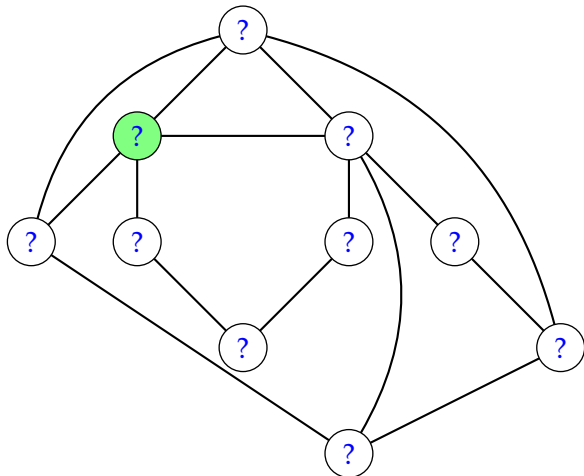
Directed



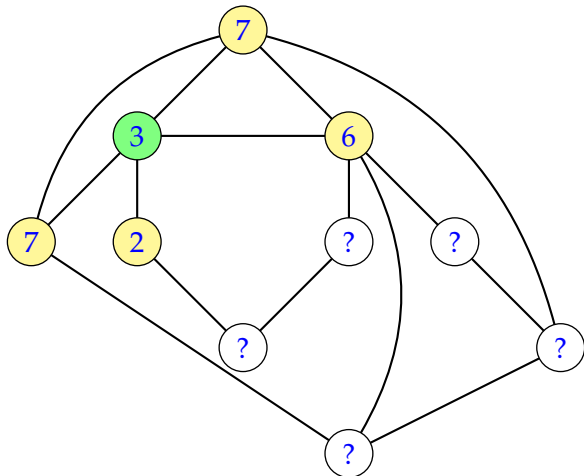
A graph of relationships over actions



A graph of relationships over actions

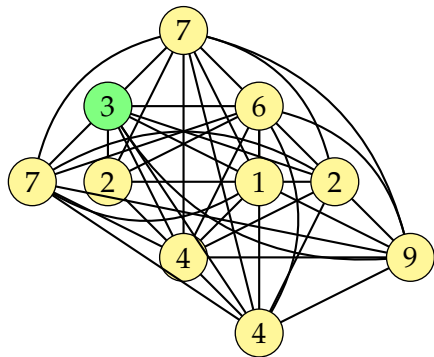


A graph of relationships over actions

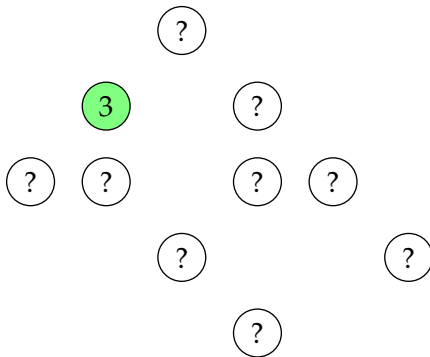


Recovering expert and bandit settings

Experts: clique



Bandits: empty graph



Exponentially weighted forecaster — Reprise

Player's strategy [Alon, C-B, Gentile, Mannor, Mansour and Shamir, 2013]

$$\bullet \mathbb{P}_t(I_t = i) \propto \exp\left(-\eta \sum_{s=1}^{t-1} \widehat{\ell}_s(i)\right) \quad i = 1, \dots, N$$

$$\bullet \widehat{\ell}_t(i) = \begin{cases} \frac{\ell_t(i)}{\mathbb{P}_t(\ell_t(i) \text{ observed})} & \text{if } \ell_t(i) \text{ is observed} \\ 0 & \text{otherwise} \end{cases}$$

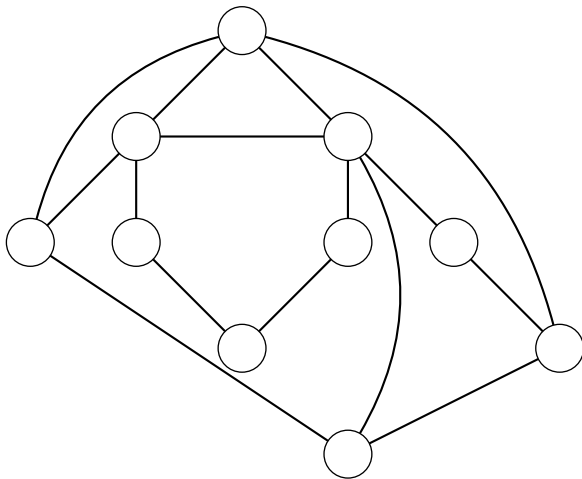
Importance sampling estimator

$$\mathbb{E}_t[\widehat{\ell}_t(i)] = \ell_t(i) \quad \text{unbiased}$$

$$\mathbb{E}_t[\widehat{\ell}_t(i)^2] \leq \frac{1}{\mathbb{P}_t(\ell_t(i) \text{ observed})} \quad \text{variance control}$$

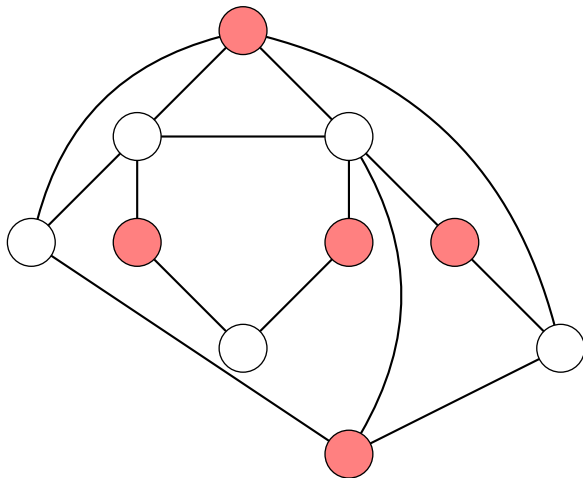
Independence number $\alpha(G)$

The size of the largest **independent set**



Independence number $\alpha(G)$

The size of the largest **independent set**



Regret bounds

Analysis (undirected graphs)

$$\begin{aligned} R_T &\leq \frac{\ln N}{\eta} + \frac{\eta}{2} \sum_{t=1}^T \sum_{i=1}^N \underbrace{\mathbb{P}(I_t = i \mid \ell_t(i) \text{ observed})}_{\leq \alpha(G)} \\ &= \sqrt{\alpha(G) T \ln N} \quad \text{by tuning } \eta \end{aligned}$$

If graph is **directed**, then bound worsens only by log factors

Special cases

- **Experts:** $\alpha(G) = 1$ $R_T \leq \sqrt{T \ln N}$
- **Bandits:** $\alpha(G) = N$ $R_T \leq \sqrt{TN \ln N}$



The loss of action i at time t depends on the player's past m actions
 $\ell_t(i) \rightarrow L_t(I_{t-m}, \dots, I_{t-1}, i)$

Adaptive regret

$$R_T^{\text{ada}} = \mathbb{E} \left[\sum_{t=1}^T L_t(I_{t-m}, \dots, I_{t-1}, I_t) - \min_{i=1, \dots, K} \sum_{t=1}^T L_t(\underbrace{i, \dots, i}_m, i) \right]$$

Minimax adaptive regret (for any constant $m > 1$)

$$R_T^{\text{ada}} = \Theta(T^{2/3})$$

Partial monitoring: not observing any loss

Dynamic pricing: Perform as the best fixed price

- 1 Post a T-shirt price
- 2 Observe if next customer buys or not
- 3 Adjust price

Feedback does not reveal the player's loss



	1	2	3	4	5
1	0	1	2	3	4
2	c	0	1	2	3
3	c	c	0	1	2
4	c	c	c	0	1
5	c	c	c	c	0

Loss matrix

	1	2	3	4	5
1	1	1	1	1	1
2	0	1	1	1	1
3	0	0	1	1	1
4	0	0	0	1	1
5	0	0	0	0	1

Feedback matrix



A characterization of minimax regret

Special case

Multiarmed bandits: loss and feedback matrix are the same

A general gap theorem [Bartok, Foster, Pál, Rakhlin and Szepesvári, 2013]

- A constructive characterization of the minimax regret for any pair of loss/feedback matrix
- Only three possible rates for nontrivial games:
 - 1 Easy games (e.g., bandits): $\Theta(\sqrt{T})$
 - 2 Hard games (e.g., revealing action): $\Theta(T^{2/3})$
 - 3 Impossible games: $\Theta(T)$



Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel
- 4 The joy of convex**
- 5 The joy of convex (without the gradient)



A game equivalent to prediction with expert advice

Online linear optimization in the simplex

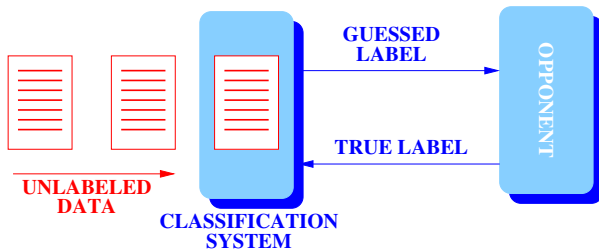
- 1 Play point \mathbf{p}_t from the N -dimensional simplex Δ_N
- 2 Incur linear loss $\mathbb{E}[\ell_t(\mathbf{I}_t)] = \mathbf{p}_t^\top \tilde{\ell}_t$
- 3 Observe loss gradient $\tilde{\ell}_t$

Regret: compete against the best point in the simplex

$$\begin{aligned} \sum_{t=1}^T \mathbf{p}_t^\top \tilde{\ell}_t - \underbrace{\min_{\mathbf{q} \in \Delta_N} \sum_{t=1}^T \mathbf{q}^\top \tilde{\ell}_t}_{\text{Best point in the simplex}} \\ = \min_{i=1, \dots, N} \frac{1}{T} \sum_{t=1}^T \tilde{\ell}_t(i) \end{aligned}$$



From game theory to machine learning



- Opponent's moves y_t are viewed as **values or labels** assigned to observations $x_t \in \mathbb{R}^d$ (e.g., categories of documents)
- A repeated game between the player choosing an element w_t of a **linear space** and the opponent choosing a label y_t for x_t
- Regret with respect to **best element** in the linear space

- 1 Play point \mathbf{w}_t from a **convex linear space** S
- 2 Incur convex loss $l_t(\mathbf{w}_t)$
- 3 Observe loss gradient $\nabla l_t(\mathbf{w}_t)$
- 4 Update point: $\mathbf{w}_t \rightarrow \mathbf{w}_{t+1} \in S$

Example

- Regression with square loss: $l_t(\mathbf{w}) = (\mathbf{w}^\top \mathbf{x}_t - y_t)^2$ $y_t \in \mathbb{R}$
- Classification with hinge loss: $l_t(\mathbf{w}) = [1 - y_t \mathbf{w}^\top \mathbf{x}_t]_+$
 $y_t \in \{-1, +1\}$

Regret

$$\sum_{t=1}^T l_t(\mathbf{w}_t) - \inf_{\mathbf{u} \in S} \sum_{t=1}^T l_t(\mathbf{u})$$

Finding a good online algorithm

Follow the leader

$$\mathbf{w}_{t+1} = \operatorname{arginf}_{\mathbf{w} \in S} \sum_{s=1}^t \ell_s(\mathbf{w})$$

Regret can be linear due to **lack of stability**

Example

$$S = [-1, +1] \quad \ell_1(\mathbf{w}) = 1 + \frac{\mathbf{w}}{2} \quad \ell_t(\mathbf{w}) = \begin{cases} -\mathbf{w} & \text{if } t \text{ is even} \\ +\mathbf{w} & \text{if } t \text{ is odd} \end{cases}$$



Regularized online learning

Strong convexity

$\Phi : S \rightarrow \mathbb{R}$ is β -strongly convex w.r.t. a norm $\|\cdot\|$ if for all $\mathbf{u}, \mathbf{v} \in S$

$$\Phi(\mathbf{v}) \geq \Phi(\mathbf{u}) + \nabla\Phi(\mathbf{u})^\top(\mathbf{v} - \mathbf{u}) + \frac{\beta}{2} \|\mathbf{u} - \mathbf{v}\|^2$$

Example: $\Phi(\mathbf{v}) = \frac{1}{2} \|\mathbf{v}\|^2$

Follow the regularized leader

[Shalev-Shwartz, 2007; Abernethy, Hazan and Rakhlin, 2008]

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in S} \left[\eta \sum_{s=1}^t \ell_s(\mathbf{w}) + \Phi(\mathbf{w}) \right]$$

Φ is a strongly convex regularizer defined on S

Deriving an incremental update

Linearization of convex losses

$$l_t(\mathbf{w}_t) - l_t(\mathbf{u}) \leq \underbrace{\nabla l_t(\mathbf{w}_t)^\top}_{\tilde{\ell}_t} \mathbf{w}_t - \underbrace{\nabla l_t(\mathbf{w}_t)^\top}_{\tilde{\ell}_t} \mathbf{u}$$

Follow the regularized leader with linearized losses

$$\begin{aligned} \mathbf{w}_{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathcal{S}} \left(\eta \underbrace{\sum_{s=1}^t \tilde{\ell}_s^\top}_{\theta_t} \mathbf{w} + \Phi(\mathbf{w}) \right) = \operatorname{argmax}_{\mathbf{w} \in \mathcal{S}} \left(-\eta \theta_t^\top \mathbf{w} - \Phi(\mathbf{w}) \right) \\ &= \nabla \Phi^*(-\eta \theta_t) \end{aligned}$$

Φ^* is the **convex dual** of Φ



Recall:

$$\mathbf{w}_{t+1} = \nabla\Phi^*(-\eta\boldsymbol{\theta}_t) = \nabla\Phi^*\left(-\eta\sum_{s=1}^t\nabla\ell_s(\mathbf{w}_s)\right)$$

Online Mirror Descent

Parameters: Strongly convex regularizer Φ and learning rate $\eta > 0$

Initialize: $\boldsymbol{\theta}_1 = \mathbf{0}$ // primal parameter

For $t = 1, 2, \dots$

- 1 Use $\mathbf{w}_t = \nabla\Phi^*(\boldsymbol{\theta}_t)$ // dual parameter (via mirror step)
- 2 Suffer loss $\ell_t(\mathbf{w}_t)$
- 3 Observe loss gradient $\nabla\ell_t(\mathbf{w}_t)$
- 4 Update $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \eta\nabla\ell_t(\mathbf{w}_t)$ // gradient step

Some examples

- **Exponentiated gradient:** $S = \Delta_N$ and $\Phi(\mathbf{w}) = \sum_{i=1}^d w_i \ln w_i$
[Kivinen and Warmuth, 1997]
- **Online Gradient Descent:** $S = \mathbb{R}^d$ and $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$
[Zinkevich, 2003]
- **p-norm Gradient Descent:** $S = \mathbb{R}^d$ and $\Phi(\mathbf{w}) = \frac{1}{2(p-1)} \|\mathbf{w}\|_p^2$
[Gentile, 2003]
- **Matrix gradient descent**
[Cavallanti, C-B and Gentile, 2010; Kakade, Shalev-Shwartz and Tewari, 2012]



General regret bound

Analysis relies on smoothness of Φ^* in order to bound increments $\Phi^*(\theta_{t+1}) - \Phi^*(\theta_t)$ via $\|\nabla \ell_t(\mathbf{w}_t)\|_*^2$

Oracle bound

[Kakade, Shalev-Shwartz and Tewari, 2012]

$$\underbrace{\sum_{t=1}^T \ell_t(\mathbf{w}_t)}_{\text{cumulative loss}} \leq \inf_{\mathbf{u} \in \mathcal{S}} \left(\underbrace{\sum_{t=1}^T \ell_t(\mathbf{u})}_{\text{model fit}} + \underbrace{\frac{\Phi(\mathbf{u})}{\eta}}_{\text{model cost}} \right) + \frac{\eta}{2} \sum_{t=1}^T \frac{\|\nabla \ell_t(\mathbf{w}_t)\|_*^2}{\beta}$$

ℓ_1, ℓ_2, \dots are arbitrary convex losses

- If gradients are bounded, then $R_T = \mathcal{O}(\sqrt{T})$
- This is optimal for general convex losses ℓ_t
- If all ℓ_t are **strongly convex**, then $R_T = \mathcal{O}(\ln T)$

Regularization via stochastic smoothing

Follow the perturbed leader

[Kalai and Vempala, 2005]

$$\mathbf{w}_{t+1} = \mathbb{E} \left[\operatorname{argmin}_{\mathbf{w} \in S} \left(\eta \boldsymbol{\theta}_t^\top \mathbf{w} + \mathbf{Z}^\top \mathbf{w} \right) \right]$$

- The distribution of \mathbf{Z} must be “stable” (small variance and small average sensitivity)
- For some choices of \mathbf{Z} , FPL becomes equivalent to OMD
[Abernethy, Lee, Sinha and Tewari, 2014]



Adaptive regularization

Online Ridge Regression

[Vovk, 2001; Azoury and Warmuth, 2001]

$$\sum_{t=1}^T (\mathbf{w}_t^\top \mathbf{x}_t - y_t)^2 \leq \inf_{\mathbf{u} \in \mathbb{R}^d} \left(\sum_{t=1}^T (\mathbf{u}^\top \mathbf{x}_t - y_t)^2 + \|\mathbf{u}\|^2 \right) + d \ln \left(1 + \frac{T}{d} \right)$$

$$\Phi_t(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_{\Lambda_t}^2 \quad \Lambda_t = \mathbf{I} + \sum_{s=1}^t \mathbf{x}_s \mathbf{x}_s^\top$$

More examples

- **Online Newton Step** [Hazan, Agarwal and Kale, 2007]
Logarithmic regret for exp-concave loss functions
- **AdaGrad** [Duchi, Hazan and Singer, 2010]
Competitive with “optimal” fixed regularizer
- **Scale-invariant algorithms** [Ross, Mineiro and Langford, 2013]
Regret invariant w.r.t. rescaling of single features

Nonstationarity

- If data source is not fitted well by any model in the class, then comparing to the **best model** $\mathbf{u} \in S$ is trivial
- Compare instead to the best **sequence** $\mathbf{u}_1, \mathbf{u}_2, \dots \in S$ of models

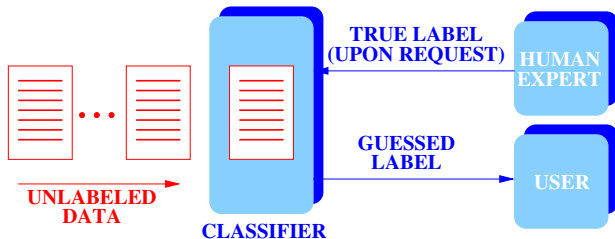
Shifting Regret for Online Mirror Descent

[Zinkevich, 2003]

$$\underbrace{\sum_{t=1}^T \ell_t(\mathbf{w}_t)}_{\text{cumulative loss}} \leq \inf_{\mathbf{u}_1, \dots, \mathbf{u}_T \in S} \underbrace{\sum_{t=1}^T \ell_t(\mathbf{u}_t)}_{\text{model fit}} + \underbrace{\sum_{t=1}^T \|\mathbf{u}_t - \mathbf{u}_{t-1}\|}_{\text{shifting model cost}} + \text{diam}(S) + \square$$



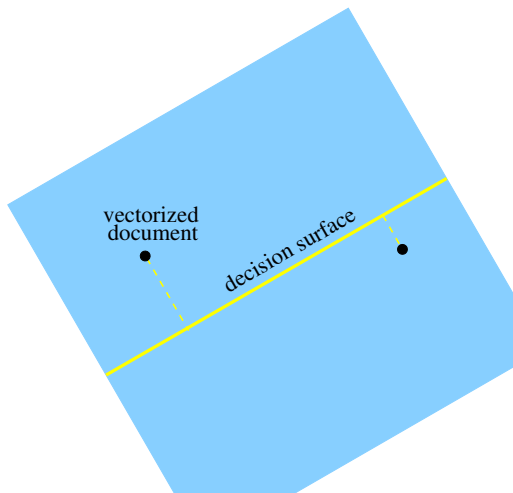
Online active learning



- 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?

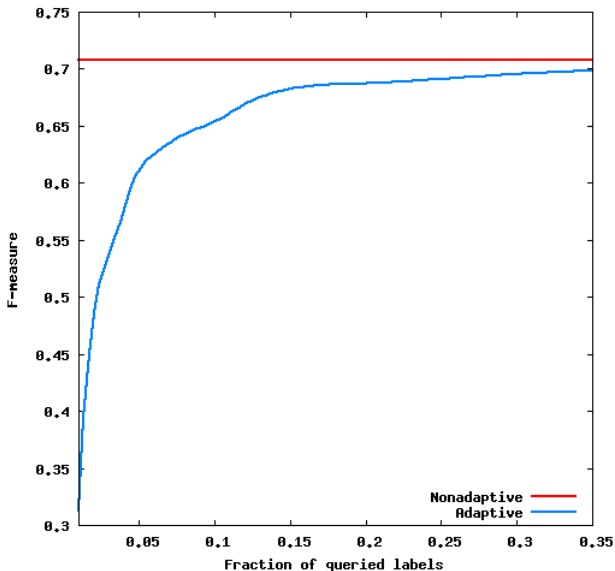


Opponent avoids causing mistakes on documents far away from decision surface

Probability of querying a document proportional to inverse distance to decision surface

Binary classification performance guarantee remains **identical** (in expectation) to the full sampling case

Experiments on document categorization



Stochastic Online Mirror Descent

Parameters: Strongly convex regularizer Φ and learning rate $\eta > 0$
Initialize: $\theta_1 = 0$ // primal parameter
For $t = 1, 2, \dots$

- 1 Use $\mathbf{w}_t = \nabla\Phi^*(\theta_t)$ // mirror step with projection on S
- 2 Suffer loss $\ell_t(\mathbf{w}_t)$
- 3 Compute estimate $\hat{\mathbf{g}}_t$ of loss gradient $\nabla\ell_t(\mathbf{w}_t)$
- 4 Update $\theta_{t+1} = \theta_t - \eta\hat{\mathbf{g}}_t$ // gradient step

Typically, $\Phi(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$ (stochastic OGD)





original



sampled

Obtain a few attributes
from each training example

- Use Stochastic OGD with square loss $\ell_t(\mathbf{w}) = \frac{1}{2}(\mathbf{w}^\top \mathbf{x}_t - y_t)^2$
- $\nabla \ell_t(\mathbf{w}) = (\mathbf{w}^\top \mathbf{x}_t - y_t) \mathbf{x}_t$

Unbiased estimate of square loss gradient using **two** attributes

- 1 Estimate of $\mathbf{w}^\top \mathbf{x}$: query \mathbf{x}_i according to $p(i) = \frac{|w_i|}{\|\mathbf{w}\|_1}$
- 2 Estimate of \mathbf{x} : query \mathbf{x}_j uniformly at random
- 3 Gradient estimate: $\hat{\mathbf{g}} = \left(\|\mathbf{w}\|_1 \operatorname{sgn}(w_i) \mathbf{x}_i - y \right) d \mathbf{x}_j \mathbf{e}_j$

Summary

- 1 My beautiful regret
- 2 A supposedly fun game I'll play again
- 3 A graphic novel
- 4 The joy of convex
- 5 The joy of convex (without the gradient)



Online convex optimization with bandit feedback

For $T = 1, 2, \dots$

- 1 Play point \mathbf{w}_t from a convex linear space S
- 2 Incur and observe convex loss $\ell_t(\mathbf{w}_t)$
- 3 Update point: $\mathbf{w}_t \rightarrow \mathbf{w}_{t+1} \in S$

Regret

$$R_T = \mathbb{E} \left[\sum_{t=1}^T \ell_t(\mathbf{w}_t) \right] - \inf_{\mathbf{u} \in S} \sum_{t=1}^T \ell_t(\mathbf{u})$$



Gradient descent without a gradient

[Flaxman, Kalai and McMahan, 2004]

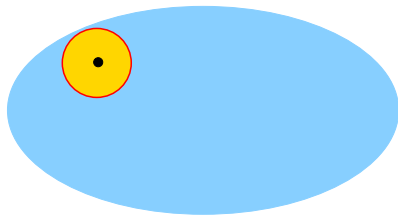
- Run stochastic OGD using a **perturbed version** of \mathbf{w}_t : $\mathbf{w}_t + \delta \mathbf{U}$ (\mathbf{U} is a random unit vector and $\delta > 0$)

- Gradient estimate $\hat{\mathbf{g}}_t = \frac{d}{\delta} l_t(\mathbf{w}_t + \delta \mathbf{U}) \mathbf{U}$

- **Fact (Stokes' Theorem)**: If l_t were differentiable, then

$$\mathbb{E}[\hat{\mathbf{g}}_t] = \nabla \mathbb{E}[l_t(\mathbf{w}_t + \delta \mathbf{B})]$$

where \mathbf{B} is a random vector in the unit ball



$\hat{\mathbf{g}}_t$ estimates the gradient of a **locally smoothed** version of l_t



- If ℓ_t is **Lipschitz**, then the smoothed version is a good approximation of ℓ_t
- Radius δ of perturbation controls bias/variance trade-off

Regret of stochastic OGD for convex and Lipschitz loss sequences

$$R_T = \mathcal{O}(T^{3/4})$$



Guarantees

- If ℓ_t is **Lipschitz**, then the smoothed version is a good approximation of ℓ_t
- Radius δ of perturbation controls bias/variance trade-off

Regret of stochastic OGD for convex and Lipschitz loss sequences

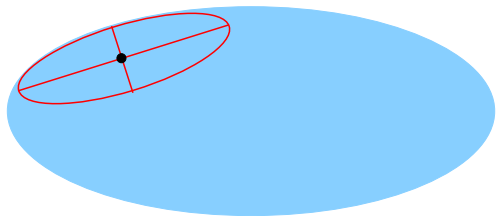
$$R_T = \mathcal{O}(T^{3/4})$$

The linear case

- Assume losses are **linear** functions on S , $\ell_t(\mathbf{w}) = \ell_t^\top \mathbf{w}$
- Can we achieve a better rate?



- Run stochastic OGD regularized with a **self-concordant function** for S
- Variance control through the associated **Dikin ellipsoid**
- Loss estimate $\hat{\ell}_t$ obtained via **perturbed point** $W_t \pm e_i \sqrt{\lambda_i}$ $\{e_i, \lambda_i\}$ is a randomly drawn eigenvector-eigenvalue pair of Dikin ellipsoid



Regret for linear functions

$$R_T = \mathcal{O}\left(d^{3/2} \sqrt{T \ln T}\right)$$



- Build an ε -cover $S_0 \subseteq S$ of size ε^{-d}
- Use experts algorithm (e.g., exponential weights) to draw actions $W_t \in S_0$ and use unbiased linear estimator for the loss

$$\hat{\ell}_t = P_t^{-1} W_t W_t^\top \ell_t \quad \text{where} \quad P_t = \mathbb{E}[W_t W_t^\top]$$

- Mix exponential weights with **exploration distribution** μ over the actions in the cover:

$$p_t(\mathbf{w}) = (1 - \gamma) \underbrace{q_t(\mathbf{w})}_{\text{exp. distrib.}} + \gamma \mu(\mathbf{w}) \quad (0 \leq \gamma \leq 1)$$

- μ controls the variance of the loss estimates by ensuring all directions are sampled often enough



Regret bound

$$R_T = \mathcal{O} \left(d \sqrt{\left(\frac{1}{d \lambda_{\min}} + 1 \right) T \ln T} \right)$$

λ_{\min} = smallest eigenvalue of $\mathbb{E}_{\mu} [W W^T]$

- λ_{\min}^{-1} is proportional to the variance of loss estimates
- When $\lambda_{\min} \approx \frac{1}{d}$ we get the **optimal bound** $\Theta \left(d \sqrt{T \ln T} \right)$
- If μ is **uniform over all actions**, the above happens when action space is approximately **isotropic**



Choose a basis under which the action set looks isotropic

- There are at most $\mathcal{O}(d^2)$ **contact points** between S_0 and Löwner ellipsoid (the min volume ellipsoid enclosing S_0)
- Put exploration distribution μ on these contact points
- This ensures that $\mathbb{E}_\mu[WW^\top]$ is isotropic: $\lambda_{\min} = \frac{1}{d}$

- Exploration on contact points of Löwner ellipsoid achieves optimal regret

$$R_T = \mathcal{O}\left(d \sqrt{T \ln T}\right)$$

- However, this construction is not efficient in general
- An efficient construction uses **volumetric ellipsoids** [Hazan, Gerber and Meka, 2014]

Conclusions

More applications

- Portfolio management
- Matrix completion
- Competitive analysis of algorithms
- Recommendation systems

Some open problems

- Exact rates for bandit convex optimization
- Trade-offs between regret bounds and running times
- Online tensor and spectral learning
- Problems with states

