

Online Learning

Lecture 2

Nicolò Cesa-Bianchi

Università degli Studi di Milano

Follow the Regularized Leader

▶ If losses lack curvature, FTL is unstable

▶ We can introduce curvature using a real-valued regularizer ψ

▶ $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \ell_s(\mathbf{w})$

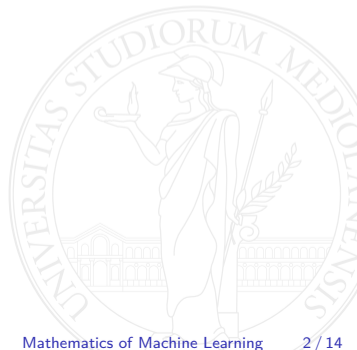
▶ Example: the SVM objective function is $\operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \frac{1}{m} \sum_{t=1}^m \ell_t(\mathbf{w})$

▶ If ℓ_t are all convex, this is equivalent to FTL over λ -strongly convex losses $\frac{\lambda}{2} \|\cdot\|_2^2 + \ell_t$

▶ How can we solve this constrained optimization problem?

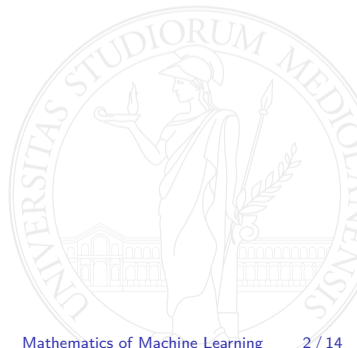
FTRL with linearized losses

► Linearized regret:
$$R_T(\mathbf{u}) = \sum_{t=1}^T \left(\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u}) \right) \leq \sum_{t=1}^T \nabla \ell_t(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{u})$$



FTRL with linearized losses

- ▶ Linearized regret: $R_T(\mathbf{u}) = \sum_{t=1}^T (\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u})) \leq \sum_{t=1}^T \nabla \ell_t(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{u})$
- ▶ Pretend all losses are linear: $\ell'_t(\mathbf{w}) = \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$



FTRL with linearized losses

- ▶ Linearized regret: $R_T(\mathbf{u}) = \sum_{t=1}^T (\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u})) \leq \sum_{t=1}^T \nabla \ell_t(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{u})$
- ▶ Pretend all losses are linear: $\ell'_t(\mathbf{w}) = \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$
- ▶ FTRL with linearized losses:

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \ell'_s(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s)$$

FTRL with linearized losses

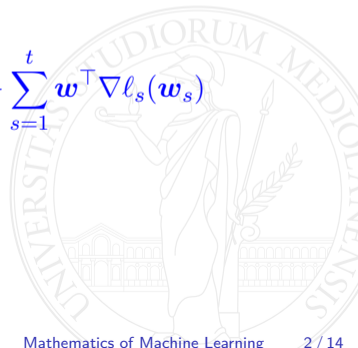
▶ Linearized regret: $R_T(\mathbf{u}) = \sum_{t=1}^T (\ell_t(\mathbf{w}_t) - \ell_t(\mathbf{u})) \leq \sum_{t=1}^T \nabla \ell_t(\mathbf{w}_t)^\top (\mathbf{w}_t - \mathbf{u})$

▶ Pretend all losses are linear: $\ell'_t(\mathbf{w}) = \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$

▶ FTRL with linearized losses:

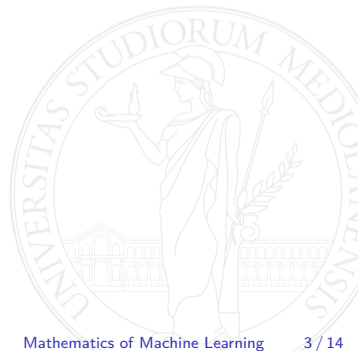
$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \ell'_s(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s)$$

▶ We still have a constrained optimization problem to solve



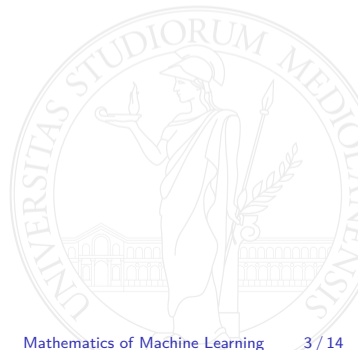
Bregman divergences

- ▶ Assume ψ is convex



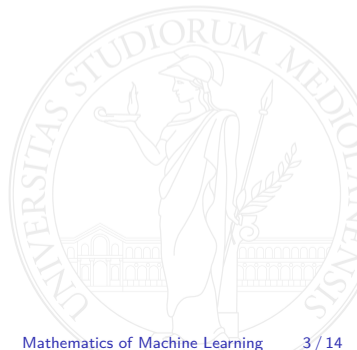
Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$



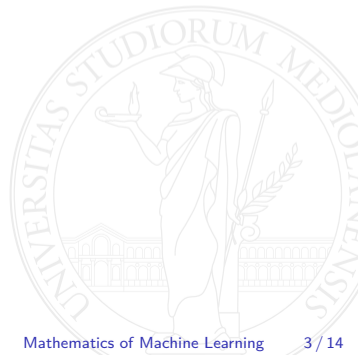
Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}



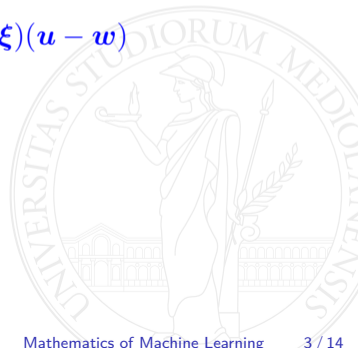
Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}
- ▶ If ψ is μ -strongly convex w.r.t. $\|\cdot\|$, then $B_\psi(\mathbf{u}, \mathbf{w}) \geq \frac{\mu}{2} \|\mathbf{u} - \mathbf{w}\|^2$



Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}
- ▶ If ψ is μ -strongly convex w.r.t. $\|\cdot\|$, then $B_\psi(\mathbf{u}, \mathbf{w}) \geq \frac{\mu}{2} \|\mathbf{u} - \mathbf{w}\|^2$
- ▶ If ψ is twice differentiable, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2}(\mathbf{u} - \mathbf{w})^\top \nabla^2\psi(\boldsymbol{\xi})(\mathbf{u} - \mathbf{w})$
for some $\boldsymbol{\xi}$ on the line segment joining \mathbf{u} and \mathbf{w}



Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}
- ▶ If ψ is μ -strongly convex w.r.t. $\|\cdot\|$, then $B_\psi(\mathbf{u}, \mathbf{w}) \geq \frac{\mu}{2} \|\mathbf{u} - \mathbf{w}\|^2$
- ▶ If ψ is twice differentiable, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2}(\mathbf{u} - \mathbf{w})^\top \nabla^2\psi(\boldsymbol{\xi})(\mathbf{u} - \mathbf{w})$
for some $\boldsymbol{\xi}$ on the line segment joining \mathbf{u} and \mathbf{w}
- ▶ If $\psi = \frac{1}{2} \|\cdot\|_2^2$, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u} - \mathbf{w}\|_2^2$

squared Euclidean distance

Bregman divergences

- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}
- ▶ If ψ is μ -strongly convex w.r.t. $\|\cdot\|$, then $B_\psi(\mathbf{u}, \mathbf{w}) \geq \frac{\mu}{2} \|\mathbf{u} - \mathbf{w}\|^2$
- ▶ If ψ is twice differentiable, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2}(\mathbf{u} - \mathbf{w})^\top \nabla^2\psi(\boldsymbol{\xi})(\mathbf{u} - \mathbf{w})$
for some $\boldsymbol{\xi}$ on the line segment joining \mathbf{u} and \mathbf{w}
- ▶ If $\psi = \frac{1}{2} \|\cdot\|_2^2$, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u} - \mathbf{w}\|_2^2$
- ▶ Δ_d is the d -dimensional probability simplex

squared Euclidean distance

Bregman divergences

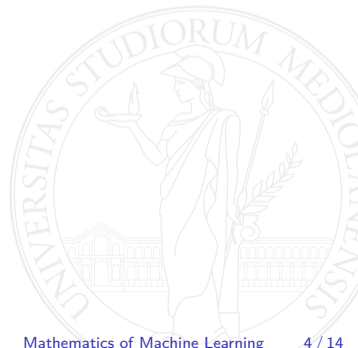
- ▶ Assume ψ is convex
- ▶ Bregman divergence: $B_\psi(\mathbf{u}, \mathbf{w}) = \psi(\mathbf{u}) - \psi(\mathbf{w}) - \nabla\psi(\mathbf{w})^\top(\mathbf{u} - \mathbf{w})$
- ▶ Error in first-order Taylor expansion of ψ around \mathbf{w}
- ▶ If ψ is μ -strongly convex w.r.t. $\|\cdot\|$, then $B_\psi(\mathbf{u}, \mathbf{w}) \geq \frac{\mu}{2} \|\mathbf{u} - \mathbf{w}\|^2$
- ▶ If ψ is twice differentiable, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2}(\mathbf{u} - \mathbf{w})^\top \nabla^2\psi(\boldsymbol{\xi})(\mathbf{u} - \mathbf{w})$
for some $\boldsymbol{\xi}$ on the line segment joining \mathbf{u} and \mathbf{w}
- ▶ If $\psi = \frac{1}{2} \|\cdot\|_2^2$, then $B_\psi(\mathbf{u}, \mathbf{w}) = \frac{1}{2} \|\mathbf{u} - \mathbf{w}\|_2^2$ squared Euclidean distance
- ▶ Δ_d is the d -dimensional probability simplex
- ▶ If $\mathbf{p} \in \Delta_d$ and $\psi(\mathbf{p}) = \sum_i p(i) \ln p(i)$, then $B_\psi(\mathbf{p}, \mathbf{q}) = \sum_i p(i) \ln \frac{p(i)}{q(i)}$ KL-divergence

Bregman projections

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be strictly convex and differentiable in \mathbb{V} .

Then $\operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} f(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_f(\mathbf{w}, \mathbf{w}')$ where $\mathbf{w}' = \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$



Bregman projections

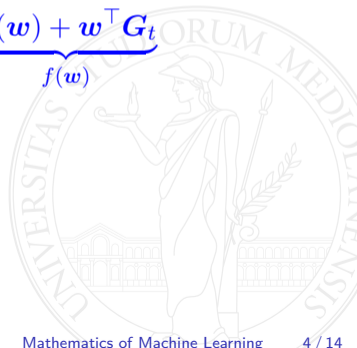
Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be strictly convex and differentiable in \mathbb{V} .

Then $\operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} f(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_f(\mathbf{w}, \mathbf{w}')$ where $\mathbf{w}' = \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \underbrace{\psi(\mathbf{w}) + \mathbf{w}^\top \mathbf{G}_t}_{f(\mathbf{w})}$$

where $\mathbf{G}_t = \sum_{s=1}^t \nabla \ell_s(\mathbf{w}_s)$



Bregman projections

Theorem

Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be strictly convex and differentiable in \mathbb{V} .

Then $\operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} f(\mathbf{w}) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_f(\mathbf{w}, \mathbf{w}')$ where $\mathbf{w}' = \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w})$

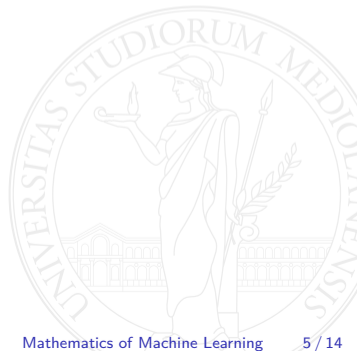
$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s) = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \underbrace{\psi(\mathbf{w}) + \mathbf{w}^\top \mathbf{G}_t}_{f(\mathbf{w})}$$

where $\mathbf{G}_t = \sum_{s=1}^t \nabla \ell_s(\mathbf{w}_s)$

Fact: $B_f \equiv B_\psi$ (because B_ψ is invariant with respect to addition of linear functions)

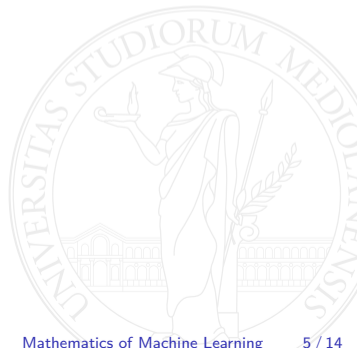
FTRL update in two steps:

$$1. \mathbf{w}'_{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s)$$



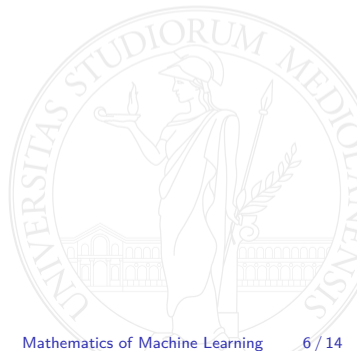
FTRL update in two steps:

1. $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} \psi(\mathbf{w}) + \sum_{s=1}^t \mathbf{w}^\top \nabla \ell_s(\mathbf{w}_s)$
2. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$ (Bregman projection)



Convex duality

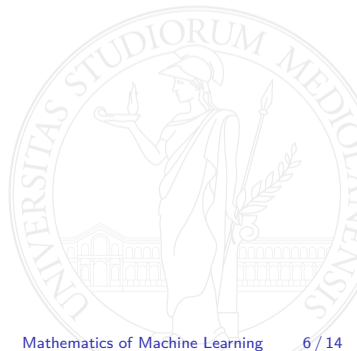
- ▶ Assume ψ is convex



Convex duality

► Assume ψ is convex

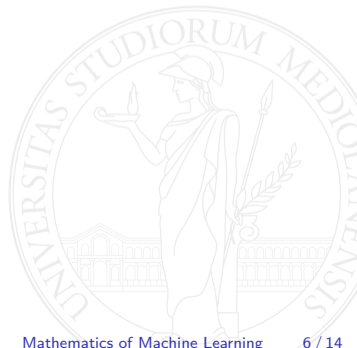
► $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\psi(\mathbf{w}) + \mathbf{w}^\top \mathbf{G}_t \right) = \operatorname{argmax}_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$ for $\boldsymbol{\theta}_{t+1} = -\mathbf{G}_t$



Convex duality

- ▶ Assume ψ is convex
- ▶ $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} \left(\psi(\mathbf{w}) + \mathbf{w}^\top \mathbf{G}_t \right) = \operatorname{argmax}_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$ for $\boldsymbol{\theta}_{t+1} = -\mathbf{G}_t$
- ▶ The convex function ψ^* is the **convex conjugate** of ψ

$$\psi^*(\boldsymbol{\theta}) = \max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta} - \psi(\mathbf{w}) \right)$$



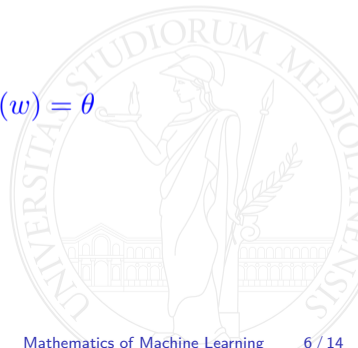
Convex duality

- ▶ Assume ψ is convex
- ▶ $\mathbf{w}'_{t+1} = \underset{\mathbf{w}}{\operatorname{argmin}} \left(\psi(\mathbf{w}) + \mathbf{w}^\top \mathbf{G}_t \right) = \underset{\mathbf{w}}{\operatorname{argmax}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$ for $\boldsymbol{\theta}_{t+1} = -\mathbf{G}_t$
- ▶ The convex function ψ^* is the **convex conjugate** of ψ

$$\psi^*(\boldsymbol{\theta}) = \max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta} - \psi(\mathbf{w}) \right)$$

- ▶ If ψ is differentiable

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^\top \boldsymbol{\theta} - \psi(\mathbf{w})) = \boldsymbol{\theta} - \psi'(\mathbf{w}) = 0 \quad \text{iff} \quad \psi'(\mathbf{w}) = \boldsymbol{\theta}$$



Convex duality

- ▶ Assume ψ is convex

- ▶ $w'_{t+1} = \operatorname{argmin}_w (\psi(w) + w^\top G_t) = \operatorname{argmax}_w (w^\top \theta_{t+1} - \psi(w))$ for $\theta_{t+1} = -G_t$

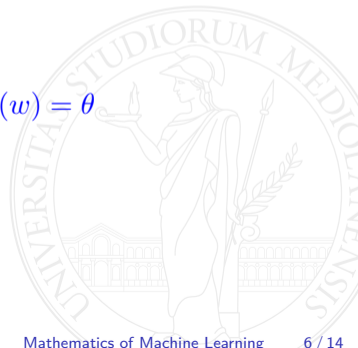
- ▶ The convex function ψ^* is the **convex conjugate** of ψ

$$\psi^*(\theta) = \max_w (w^\top \theta - \psi(w))$$

- ▶ If ψ is differentiable

$$\frac{\partial}{\partial w} (w\theta - \psi(w)) = \theta - \psi'(w) = 0 \quad \text{iff} \quad \psi'(w) = \theta$$

- ▶ $w^* = \operatorname{argmax}_{w \in \mathbb{R}} (w\theta - \psi(w))$



Convex duality

- ▶ Assume ψ is convex

- ▶ $w'_{t+1} = \operatorname{argmin}_w (\psi(w) + w^\top G_t) = \operatorname{argmax}_w (w^\top \theta_{t+1} - \psi(w))$ for $\theta_{t+1} = -G_t$

- ▶ The convex function ψ^* is the **convex conjugate** of ψ

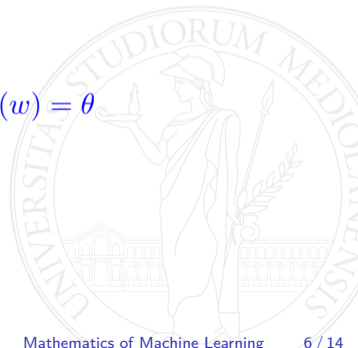
$$\psi^*(\theta) = \max_w (w^\top \theta - \psi(w))$$

- ▶ If ψ is differentiable

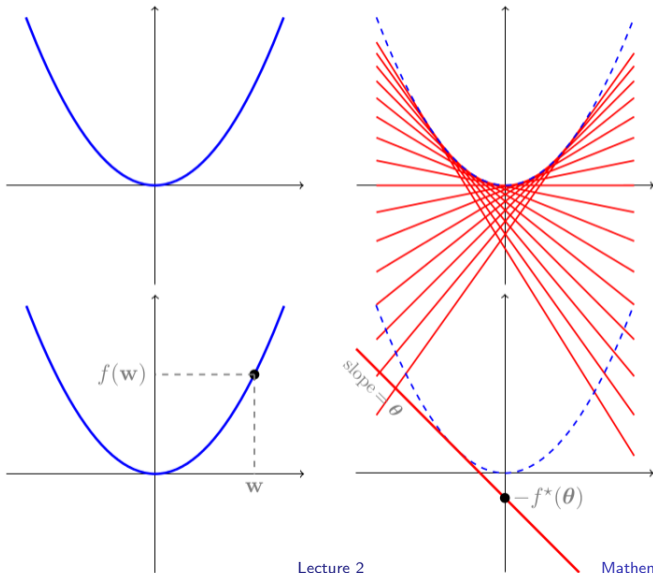
$$\frac{\partial}{\partial w} (w\theta - \psi(w)) = \theta - \psi'(w) = 0 \quad \text{iff} \quad \psi'(w) = \theta$$

- ▶ $w^* = \operatorname{argmax}_{w \in \mathbb{R}} (w\theta - \psi(w))$

- ▶ $\psi(w^*) = w^* \theta - \psi^*(\theta)$

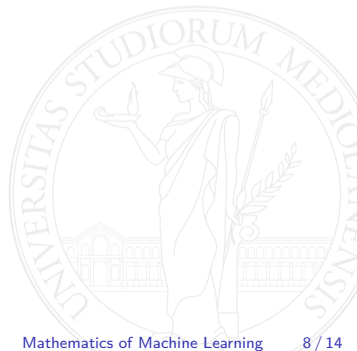


Convex duality



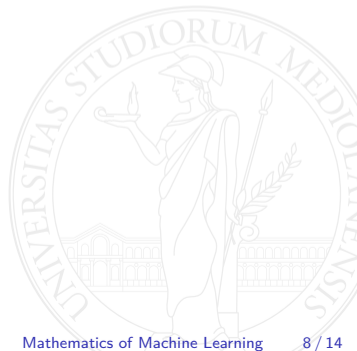
Strongly convex regularizers

$$\blacktriangleright \mathbf{w}'_{t+1} = \operatorname{argmax}_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$$



Strongly convex regularizers

- ▶ $\mathbf{w}'_{t+1} = \underset{\mathbf{w}}{\operatorname{argmax}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$
- ▶ $\psi^*(\boldsymbol{\theta}_{t+1}) = \max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$

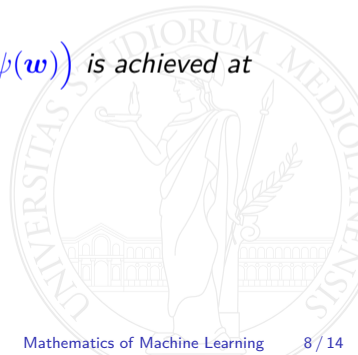


Strongly convex regularizers

- ▶ $\mathbf{w}'_{t+1} = \operatorname{argmax}_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$
- ▶ $\psi^*(\boldsymbol{\theta}_{t+1}) = \max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$

Theorem

If ψ is strongly convex, then ψ^* is differentiable and $\max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta} - \psi(\mathbf{w}) \right)$ is achieved at $\mathbf{w}^* = \nabla \psi^*(\boldsymbol{\theta})$



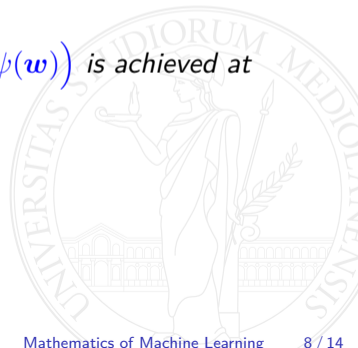
Strongly convex regularizers

- ▶ $\mathbf{w}'_{t+1} = \operatorname{argmax}_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$
- ▶ $\psi^*(\boldsymbol{\theta}_{t+1}) = \max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}) \right)$

Theorem

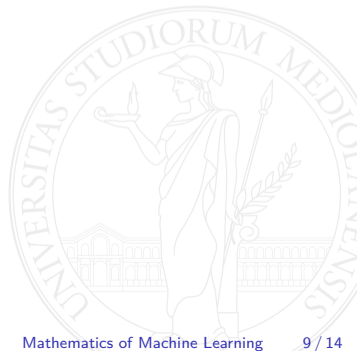
If ψ is strongly convex, then ψ^* is differentiable and $\max_{\mathbf{w}} \left(\mathbf{w}^\top \boldsymbol{\theta} - \psi(\mathbf{w}) \right)$ is achieved at $\mathbf{w}^* = \nabla \psi^*(\boldsymbol{\theta})$

This implies $\mathbf{w}'_{t+1} = \nabla \psi^*(\boldsymbol{\theta}_{t+1})$ for ψ strongly convex



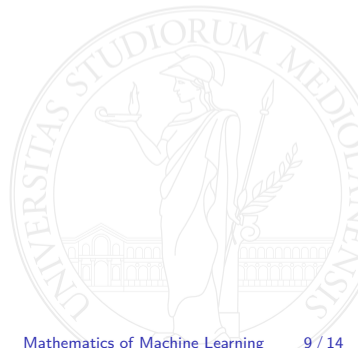
FTRL update (for strongly convex regularizers)

1. $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \nabla \ell_t(\mathbf{w}_t)$ (gradient update)



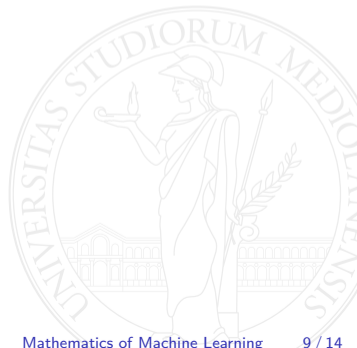
FTRL update (for strongly convex regularizers)

1. $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \nabla \ell_t(\mathbf{w}_t)$ (gradient update)
2. $\mathbf{w}'_{t+1} = \nabla \psi^*(\boldsymbol{\theta}_{t+1})$ (mirror mapping)



FTRL update (for strongly convex regularizers)

1. $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \nabla \ell_t(\mathbf{w}_t)$ (gradient update)
2. $\mathbf{w}'_{t+1} = \nabla \psi^*(\boldsymbol{\theta}_{t+1})$ (mirror mapping)
3. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$ (Bregman projection)



The Euclidean case: Lazy Online Gradient Descent

► $\psi = \frac{1}{2\eta} \|\cdot\|_2^2$

strongly convex w.r.t. $\|\cdot\|_2$ ($\eta > 0$ is a learning rate)



The Euclidean case: Lazy Online Gradient Descent

- ▶ $\psi = \frac{1}{2\eta} \|\cdot\|_2^2$
- ▶ $\psi^* = \frac{\eta}{2} \|\cdot\|_2^2$

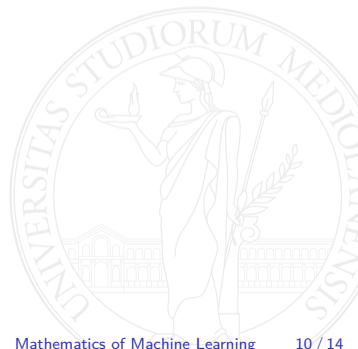
strongly convex w.r.t. $\|\cdot\|_2$ ($\eta > 0$ is a learning rate)



The Euclidean case: Lazy Online Gradient Descent

- ▶ $\psi = \frac{1}{2\eta} \|\cdot\|_2^2$
- ▶ $\psi^* = \frac{\eta}{2} \|\cdot\|_2^2$
- ▶ $\nabla\psi^*(\boldsymbol{\theta}) = \eta\boldsymbol{\theta}$

strongly convex w.r.t. $\|\cdot\|_2$ ($\eta > 0$ is a learning rate)



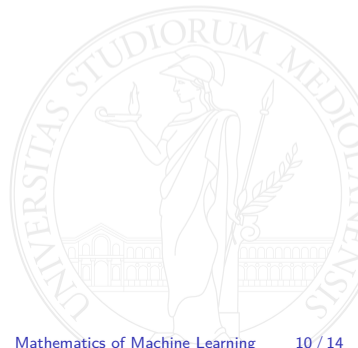
The Euclidean case: Lazy Online Gradient Descent

- ▶ $\psi = \frac{1}{2\eta} \|\cdot\|_2^2$
- ▶ $\psi^* = \frac{\eta}{2} \|\cdot\|_2^2$
- ▶ $\nabla\psi^*(\boldsymbol{\theta}) = \eta\boldsymbol{\theta}$

strongly convex w.r.t. $\|\cdot\|_2$ ($\eta > 0$ is a learning rate)

FTRL update (Projected Lazy OGD):

1. $\mathbf{w}'_{t+1} = -\eta \sum_{s=1}^t \nabla \ell_s(\mathbf{w}_s)$



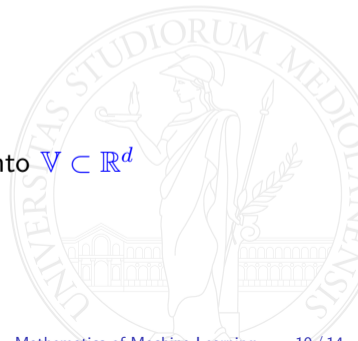
The Euclidean case: Lazy Online Gradient Descent

- ▶ $\psi = \frac{1}{2\eta} \|\cdot\|_2^2$
- ▶ $\psi^* = \frac{\eta}{2} \|\cdot\|_2^2$
- ▶ $\nabla\psi^*(\boldsymbol{\theta}) = \eta\boldsymbol{\theta}$

strongly convex w.r.t. $\|\cdot\|_2$ ($\eta > 0$ is a learning rate)

FTRL update (Projected Lazy OGD):

1. $\mathbf{w}'_{t+1} = -\eta \sum_{s=1}^t \nabla \ell_s(\mathbf{w}_s)$
2. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} \|\mathbf{w} - \mathbf{w}'_{t+1}\|_2$ standard Euclidean projection onto $\mathbb{V} \subset \mathbb{R}^d$



The Entropic case: Exponentiated Gradient

- ▶ $\mathbb{V} = \Delta_d$ (probability simplex)



The Entropic case: Exponentiated Gradient

- ▶ $\mathbb{V} = \Delta_d$ (probability simplex)
- ▶ $\psi(\mathbf{p}) = \frac{1}{\eta} \sum_i p(i) \ln p(i)$ for $\mathbf{p} \in \Delta_d$

strongly convex w.r.t. $\|\cdot\|_1$



The Entropic case: Exponentiated Gradient

- ▶ $\mathbb{V} = \Delta_d$ (probability simplex)
- ▶ $\psi(\mathbf{p}) = \frac{1}{\eta} \sum_i p(i) \ln p(i)$ for $\mathbf{p} \in \Delta_d$
- ▶ $\psi^*(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta_d} (\mathbf{p}^\top \boldsymbol{\theta} - \psi(\mathbf{p})) = \frac{1}{\eta} \ln \left(\sum_i e^{\eta \theta(i)} \right)$

strongly convex w.r.t. $\|\cdot\|_1$

we solve the constrained problem



The Entropic case: Exponentiated Gradient

▶ $\mathbb{V} = \Delta_d$ (probability simplex)

▶ $\psi(\mathbf{p}) = \frac{1}{\eta} \sum_i p(i) \ln p(i)$ for $\mathbf{p} \in \Delta_d$

▶ $\psi^*(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta_d} (\mathbf{p}^\top \boldsymbol{\theta} - \psi(\mathbf{p})) = \frac{1}{\eta} \ln \left(\sum_i e^{\eta \theta(i)} \right)$

▶ $\nabla \psi^*(\boldsymbol{\theta})_i = \frac{e^{\eta \theta(i)}}{\sum_{j=1}^d e^{\eta \theta(j)}}$

strongly convex w.r.t. $\|\cdot\|_1$

we solve the constrained problem



The Entropic case: Exponentiated Gradient

- ▶ $\mathbb{V} = \Delta_d$ (probability simplex)
- ▶ $\psi(\mathbf{p}) = \frac{1}{\eta} \sum_i p(i) \ln p(i)$ for $\mathbf{p} \in \Delta_d$
- ▶ $\psi^*(\boldsymbol{\theta}) = \max_{\mathbf{p} \in \Delta_d} (\mathbf{p}^\top \boldsymbol{\theta} - \psi(\mathbf{p})) = \frac{1}{\eta} \ln \left(\sum_i e^{\eta \theta(i)} \right)$
- ▶ $\nabla \psi^*(\boldsymbol{\theta})_i = \frac{e^{\eta \theta(i)}}{\sum_{j=1}^d e^{\eta \theta(j)}}$

FTRL update (EG):

$$p_{t+1}(i) = \frac{\exp \left(-\eta \sum_{s=1}^t \nabla \ell_s(\mathbf{p}_s)_i \right)}{\sum_{j=1}^d \exp \left(-\eta \sum_{s=1}^t \nabla \ell_s(\mathbf{p}_s)_j \right)}$$

strongly convex w.r.t. $\|\cdot\|_1$

we solve the constrained problem



Online Mirror Descent

Update minimizes trade-off between linearized loss and Bregman divergence from previous iterate

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_{\psi}(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^{\top} \nabla \ell_t(\mathbf{w}_t)$$



Online Mirror Descent

Update minimizes trade-off between linearized loss and Bregman divergence from previous iterate

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$$

Projected update:

1. $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$



Online Mirror Descent

Update minimizes trade-off between linearized loss and Bregman divergence from previous iterate

$$\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$$

Projected update:

1. $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \nabla \ell_t(\mathbf{w}_t)$
2. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$



OMD: computation of the update

Assume ψ is strongly convex



OMD: computation of the update

Assume ψ is strongly convex

$$\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t \quad (\mathbf{g}_t = \nabla l_t(\mathbf{w}_t))$$



OMD: computation of the update

Assume ψ is strongly convex

$$\begin{aligned}\mathbf{w}'_{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t && (\mathbf{g}_t = \nabla l_t(\mathbf{w}_t)) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) - \psi(\mathbf{w}_t) - \nabla \psi(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t\end{aligned}$$



OMD: computation of the update

Assume ψ is strongly convex

$$\begin{aligned}\mathbf{w}'_{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t && (\mathbf{g}_t = \nabla l_t(\mathbf{w}_t)) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) - \psi(\mathbf{w}_t) - \nabla \psi(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) + \mathbf{w}^\top (\eta_t \mathbf{g}_t - \nabla \psi(\mathbf{w}_t))\end{aligned}$$



OMD: computation of the update

Assume ψ is strongly convex

$$\begin{aligned}\mathbf{w}'_{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t && (\mathbf{g}_t = \nabla l_t(\mathbf{w}_t)) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) - \psi(\mathbf{w}_t) - \nabla \psi(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) + \mathbf{w}^\top (\eta_t \mathbf{g}_t - \nabla \psi(\mathbf{w}_t))\end{aligned}$$

We use $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} (\psi(\mathbf{w}) + \mathbf{w}^\top (\eta_t \mathbf{g}_t - \nabla \psi(\mathbf{w}_t))) = \operatorname{argmax}_{\mathbf{w}} (\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}))$
for $\boldsymbol{\theta}_{t+1} = \nabla \psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t$

OMD: computation of the update

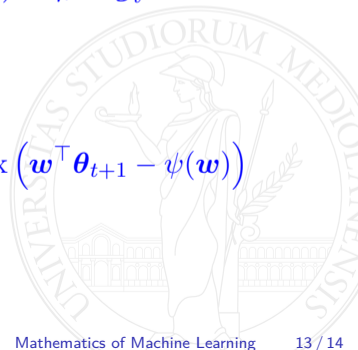
Assume ψ is strongly convex

$$\begin{aligned}\mathbf{w}'_{t+1} &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} B_\psi(\mathbf{w}, \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t && (\mathbf{g}_t = \nabla l_t(\mathbf{w}_t)) \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) - \psi(\mathbf{w}_t) - \nabla \psi(\mathbf{w}_t)^\top (\mathbf{w} - \mathbf{w}_t) + \eta_t \mathbf{w}^\top \mathbf{g}_t \\ &= \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \psi(\mathbf{w}) + \mathbf{w}^\top (\eta_t \mathbf{g}_t - \nabla \psi(\mathbf{w}_t))\end{aligned}$$

We use $\mathbf{w}'_{t+1} = \operatorname{argmin}_{\mathbf{w}} (\psi(\mathbf{w}) + \mathbf{w}^\top (\eta_t \mathbf{g}_t - \nabla \psi(\mathbf{w}_t))) = \operatorname{argmax}_{\mathbf{w}} (\mathbf{w}^\top \boldsymbol{\theta}_{t+1} - \psi(\mathbf{w}))$

for $\boldsymbol{\theta}_{t+1} = \nabla \psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t$

$$\mathbf{w}'_{t+1} = \nabla \psi^*(\boldsymbol{\theta}_{t+1}) = \nabla \psi^*(\nabla \psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t)$$



Some differences between FTRL and OMD

1. $\mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$
2. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$

1. $\mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$
2. $\mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

► $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

- ▶ $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t
- ▶ $\nabla\psi$ maps iterates to gradients

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

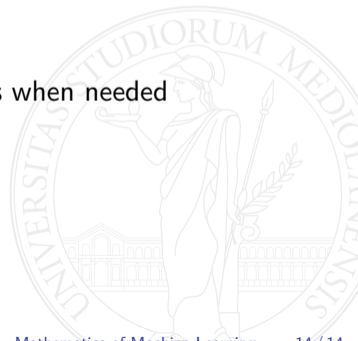
▶ $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t

▶ $\nabla\psi$ maps iterates to gradients

▶ FTRL updates a state variable $-\sum_{s=1}^t \mathbf{g}_s$ and maps it to iterates when needed

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

▶ $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t

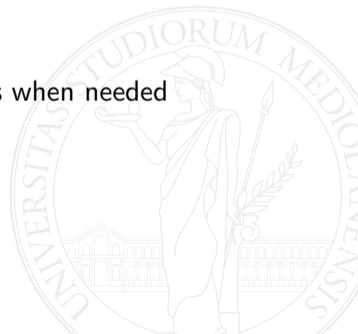
▶ $\nabla\psi$ maps iterates to gradients

▶ FTRL updates a state variable $-\sum_{s=1}^t \mathbf{g}_s$ and maps it to iterates when needed

▶ OMD maps iterates back to gradients before each update

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

▶ $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t

▶ $\nabla\psi$ maps iterates to gradients

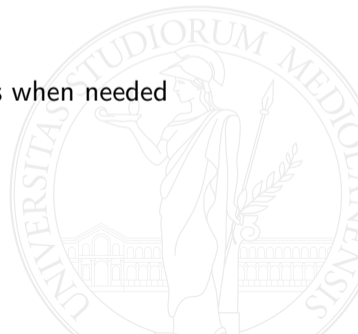
▶ FTRL updates a state variable $-\sum_{s=1}^t \mathbf{g}_s$ and maps it to iterates when needed

▶ OMD maps iterates back to gradients before each update

▶ OMD and FTRL have similar regret bounds in many cases

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathcal{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$



Some differences between FTRL and OMD

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(-\sum_{s=1}^t \mathbf{g}_s \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$

▶ $\nabla\psi^*$ maps gradients to iterates \mathbf{w}_t

▶ $\nabla\psi$ maps iterates to gradients

▶ FTRL updates a state variable $-\sum_{s=1}^t \mathbf{g}_s$ and maps it to iterates when needed

▶ OMD maps iterates back to gradients before each update

▶ OMD and FTRL have similar regret bounds in many cases

▶ In certain cases, FTRL works better than OMD when using dynamic learning rates

$$1. \mathbf{w}'_{t+1} = \nabla\psi^* \left(\nabla\psi(\mathbf{w}_t) - \eta_t \mathbf{g}_t \right)$$

$$2. \mathbf{w}_{t+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{V}} B_\psi(\mathbf{w}, \mathbf{w}'_{t+1})$$