

# Is Mirror Descent a special case of Exponential Weights?

**Dirk van der Hoeven** and Tim van Erven

08-11-2017



Universiteit  
Leiden  
The Netherlands

# Setting: Online Linear Optimization

We consider the Online Linear Optimization setting, which proceeds in rounds  $t = 1, \dots, T$ . In each round  $t$  we

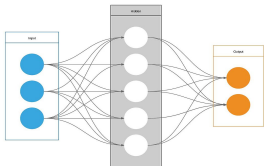
- 1 Choose a point  $\mathbf{w}_t \in \mathcal{K} \in \text{reals}^d$ , where  $\mathcal{K}$  is a convex set.
- 2 Receive gradient of convex loss function  $\mathbf{g}_t = \nabla f_t(\mathbf{w}_t)$
- 3 suffer loss  $\langle \mathbf{w}_t, \mathbf{g}_t \rangle$

Goal: keep regret  $\mathcal{R}_T(\mathbf{u})$  small

$$\mathcal{R}_T(\mathbf{u}) = \sum_{t=1}^T \langle \mathbf{w}_t, \mathbf{g}_t \rangle - \min_{\mathbf{u} \in \mathcal{K}} \sum_{t=1}^T \langle \mathbf{u}, \mathbf{g}_t \rangle$$

# Motivation

The Online Linear Optimization setting has many uses, among them are training neural networks <sup>1</sup>, gambling <sup>2</sup>, spam filtering <sup>3</sup>, and portfolio selection<sup>4</sup>.



<sup>1</sup> Neural network image by LearnDataSci from [www.learn datasci.com](http://www.learn datasci.com)

<sup>2</sup> Gambling image by History Channel from <http://www.history.com/news/ask-history/where-did-poker-originate>

<sup>3</sup> Spam image by Qwertyxp2000 from [https://commons.wikimedia.org/wiki/File:Spam\\_can.png](https://commons.wikimedia.org/wiki/File:Spam_can.png)

<sup>4</sup> stock market image by James Smith from <https://pixabay.com/en/business-stock-finance-market-1730089/>

# Algorithms

Under appropriate conditions  $\mathcal{R}_T(\mathbf{u}) = O(\sqrt{T})$

Multiple algorithms:

1. Online Gradient Descent.
2. Mirror Descent.
3. Exponential Weights.

# Problem

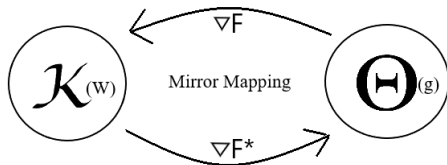
- ▶ Usually Gradient Descent and Exponential Weights are seen as special cases of Mirror Descent.
- ▶ Koolen (2016) found that Gradient Descent can be seen as a special case of Exponential Weights.
- ▶ Some interesting implications, but since Gradient Descent is a special case of Mirror Descent this also raises the following question:

**Is Mirror Descent a special case of Exponential Weights?**

# Mirror Descent

Choose suitable Legendre function  $F^*$ . Initialize  $\mathbf{w}_1 = \arg \min_{\mathbf{w}} F^*(\mathbf{w})$ , then update with:

$$\mathbf{w}_{t+1} = (\nabla F^*)^{-1} \left( \nabla F^*(\mathbf{w}_t) - \eta \mathbf{g}_t \right)$$



When  $F^*(\mathbf{w}_t) = \frac{1}{2} \|\mathbf{w}_t\|_2^2$  we obtain Gradient Descent:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t,$$

# Prediction with Expert Advice

- ▶ A special case of the Online Linear Optimization setting: The weight vector  $w_t$  as a probability distribution  $p_t$  on experts  $k = 1, \dots, K$ .
- ▶ We choose a specific loss function:  $\mathbb{E}_{k \sim p_t} [\langle e_k, g_t \rangle]$ , where  $e_k$  is the basis vector in direction  $k$  and  $g_t^k$  is the loss of expert  $k$  at time  $t$ .
- ▶ Same goal as in Online Linear Optimization, find an algorithm that has regret that grows sub-linear with  $T$ .

# Exponential Weights

Initialize  $p_t$  with *prior distribution*  $\pi$ , then update with:

$$p_{t+1}(k) = \frac{\pi(k) \exp(-\eta \sum_{i=1}^t g_i^k)}{\sum_{k=1}^K \pi(k) \exp(-\eta \sum_{i=1}^t g_i^k)},$$

where  $\eta$  is the learning rate. Usually,  $\pi$  is chosen as the uniform distribution over the experts.



# Different interpretation

- ▶ A non-standard interpretation of Exponential Weights arises if we use a non-uniform prior over a continuous set of experts parametrized by  $z \in \mathcal{K}$ .
- ▶ We use the **mean** of  $p_{t+1}$  as weights  $w_{t+1}$ .
- ▶ With a multivariate normal distribution as a prior the mean of Exponential Weights is the Gradient Descent algorithm.

# Setup

In the Online Linear Optimization setting, in each round  $t$  expert  $z$  receives loss  $\langle z, g_t \rangle$ .

Our loss becomes:

$$\begin{aligned}\mathbb{E}_{z \sim p_t}[\langle z, g_t \rangle] &= \langle \mathbb{E}_{z \sim p_t}[z], g_t \rangle \\ &= \langle w_t, g_t \rangle.\end{aligned}$$

We update  $\pi$  with:

$$p_{t+1}(z) = \frac{\pi(z) \exp(-\eta \sum_{i=1}^t \langle z, g_i \rangle)}{\int_{\mathcal{K}} \pi(z) \exp(-\eta \sum_{i=1}^t \langle z, g_i \rangle) dz}.$$

# Prior from an exponential family

Many distributions such as the normal, poisson, exponential, gamma, multinomial and many more can be written in the exponential family form:

$$p(\mathbf{z}) = e^{\langle \boldsymbol{\theta}, T(\mathbf{z}) \rangle - F(\boldsymbol{\theta})} K(\mathbf{z}),$$

where  $\boldsymbol{\theta}$  is the natural parameter,  $T(\mathbf{z})$  is the sufficient statistic,  $F(\boldsymbol{\theta})$  is the cumulant generating function, and  $K(\mathbf{z})$  is the carrier measure.

Mean:  $\mathbb{E}_p[\mathbf{z}] = \nabla F(\boldsymbol{\theta})$

# Main Result

## Theorem

*Let  $p_{t+1}$  be the Exponential Weights distribution at time  $t + 1$  with a prior from an exponential family. Let Mirror Descent be used with  $F^*$ , the convex conjugate of cumulant generating function  $F$ . Then the Mirror Descent algorithm is the mean of the Exponential Weights algorithm:*

$$\mathbb{E}_{z \sim p_{t+1}} [z] = \mathbf{w}_{t+1} = \nabla F \left( \nabla F^*(\mathbf{w}_t) - \eta \mathbf{g}_t \right).$$

## Example

With a standard multivariate normal prior the Exponential Weights distribution at time  $t + 1$  is:  $p_{t+1}(z) = N(z|\mathbf{w}_{t+1}, I)$ . The cumulant generating function is:

$$F\left(\sum_{i=1}^t \mathbf{g}_t\right) = \frac{1}{2} \left\| \sum_{i=1}^t \mathbf{g}_t \right\|_2^2$$

This gives the following mean:

$$\begin{aligned} \mathbb{E}_{\mathbf{z} \sim p_{t+1}}[\mathbf{z}] &= (\nabla F^*)^{-1} \left( \nabla F^*(\mathbf{w}_t) - \eta \mathbf{g}_t \right) \\ &= \mathbf{w}_t - \eta \mathbf{g}_t. \end{aligned}$$

# Applications of main result

1. Efficient sampling in the Linear Bandit setting
2. Nice theoretical properties of cumulant generating functions (self-concordant barriers)
3. Prior on the learning rate (exploit easy cases!)
4. Scale free algorithms (scaling of the loss becomes irrelevant)

# Conclusion

A large class of Online optimization algorithms is like learning distributions.

# References I

Koolen, W. (2016). Gradient descent as exponential weights. *Blog February 21:*

*<http://blog.wouterkoolen.info/GDasEW/post.html/>.*

van der Hoeven, D. (2016). Is mirror descent a special case of exponential weights? *MSC Thesis. Available from:*

*<http://pub.math.leidenuniv.nl/~hoevendvander/>.*



# Interpretation Gradient Descent

Gradient Descent does not learn the variance!

Can we learn the variance? Yes, with the online Newton algorithm:

$$p_{t+1}(z) = N(z|\mathbf{w}_{t+1}, (\sum_{i=1}^t c \mathbf{g}_i \mathbf{g}_i^T)^{-1})$$