

# Computational Learning Theory

## Lecture 13: Online Convex Optimisation

Lecturer: James Worrell

In this lecture we briefly introduce some basic concepts behind online convex optimisation. We present the online subgradient descent procedure and show how the Perceptron algorithm can be seen as an instance of this procedure. More generally we show that both the Perceptron and Winnow algorithms are instances of follow-the-leader algorithms, but with different regularisation functions.

### 1 Subgradients

Let  $S \subseteq \mathbb{R}^n$  be a convex set. One characterisation of when a function  $f : S \rightarrow \mathbb{R}$  is convex is that for all  $w \in S$  there exists a vector  $z \in \mathbb{R}^n$  such that for  $u \in S$  we have

$$f(u) - f(w) \geq z \cdot (u - w). \quad (1)$$

Inequality (1) says that the linear function  $g(u) = f(w) + z \cdot (u - w)$  is a lower bound of  $f$  which is tangent to  $f$  at  $w$ . (See the book of Boyd and Vandenberghe for details.)

The set of vectors  $z$  that satisfy (1) defines the collection

$$\partial f(w) = \{z \in \mathbb{R}^n : f(u) - f(w) \geq z \cdot (u - w) \text{ for all } u \in S\}$$

of *subgradients* of  $f$  at  $w$ .

If  $f$  is convex and differentiable at  $w$  then there is exactly one subgradient at  $w$ , which is equal to the gradient:  $\partial f(w) = \{\nabla f(w)\}$ .

**Example 1.** Consider the hinge function  $f : \mathbb{R} \rightarrow \mathbb{R}$ , given by  $f(w) = \max(0, -w)$ . Then  $\partial f(w) = \{0\}$  if  $w > 0$ ,  $\partial f(w) = \{-1\}$  if  $w < 0$ , and  $\partial f(w) = [-1, 0]$  if  $w = 0$ . More generally, consider a labelled feature vector  $(x, y)$ , with  $x \in \mathbb{R}^n$  and  $y \in \{-1, +1\}$ . Let  $f : \mathbb{R} \rightarrow \mathbb{R}$  be the loss function

$$f(w) = \max\{0, -y(w \cdot x)\}$$

of the linear classifier  $w$  at the point  $(x, y)$ . Then

$$\partial f(w) = \begin{cases} \{0\} & \text{if } y(w \cdot x) > 0 \\ \{-y x\} & \text{if } y(w \cdot x) < 0 \\ \{-y \alpha x : \alpha \in [0, 1]\} & \text{if } y(w \cdot x) = 0 \end{cases}$$

Let  $\rho > 0$ . We say that  $f : S \rightarrow \mathbb{R}$  is  $\rho$ -Lipschitz if for all  $u, w \in S$ ,

$$|f(u) - f(w)| \leq \rho \|u - w\|.$$

Suppose that  $f : S \rightarrow \mathbb{R}$  is  $\rho$ -Lipschitz and that  $S$  is an open convex subset of  $\mathbb{R}^n$ . Then for any  $w \in S$  and subgradient  $z \in \partial f(w)$  we have  $\|z\| \leq \rho$ . Indeed pick  $\varepsilon > 0$  sufficiently small that  $u := w + \varepsilon z \in S$ . Then

$$\rho \varepsilon \|z\| = \rho \|u - w\| \geq |f(u) - f(w)| \geq z \cdot (u - w) = \varepsilon \|z\|^2.$$

Hence  $\|z\| \leq \rho$ .

## 2 Online Convex Optimisation

The setup of online convex optimisation is as follows. Let  $S \subseteq \mathbb{R}^n$  be a convex set and  $G$  a set of convex loss functions  $S \rightarrow \mathbb{R}$ .

for  $t = 1, 2, \dots$

the learner chooses prediction  $\mathbf{w}_t \in S$

the adversary chooses convex loss function  $f_t \in G$

the learner loses  $f_t(\mathbf{w}_t)$

The regret of the learner over the first  $T$  rounds is then the difference between their total loss and the total loss of the best comparator:

$$\text{Regret}(T) = \sup_{\mathbf{u} \in S} \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})).$$

Our goal is to minimize regret. In particular, we would like the regret to be sublinear in the number of rounds  $T$  so that our *average* regret goes to 0 as  $T$  goes to infinity.

### 2.1 Follow the (Regularised) Leader

A natural policy for the learner in the situation of online convex optimisation is called *follow the leader*. The idea is to select in round  $t$  a vector  $\mathbf{w}_t \in S$  that minimises the total loss on the previous  $t-1$  rounds:

$$\mathbf{w}_t \in \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} f_i(\mathbf{w}). \quad (2)$$

Follow the leader is similar to the idea of minimising empirical error in the setting of batch learning.

The following is an example of a series of loss functions for which follow the leader generates regret linear in the number of rounds.

**Example 2.** Let  $S = [-1, 1]$  be the set of possible predictions of the learner. Assume that the adversary provides the following alternating sequence of (linear) loss functions:

$$f_1(w) = -0.5w, \quad f_2(w) = w, \quad f_3(w) = -w, \quad f_4(w) = w, \dots$$

In response the learner will generate the following sequence of predictions:

$$w_1 = 0, \quad w_2 = 1, \quad w_3 = -1, \quad w_4 = 1, \dots$$

The total loss over the first  $T$  rounds is  $T - 1$ , but the comparator  $w = 0$  obtains a loss of 0 over the first  $T$  rounds, so the regret is at least  $T - 1$ . This falls short of our goal sublinear regret.

Intuitively, the bad behaviour of follow-the-leader algorithm in Example 2 stems from from the instability of the predictions, which oscillated between the two vertices of the set  $S$ . We can improve the stability of the follow-the-leader algorithm by adding a *regulariser term*  $R(\mathbf{w})$  to the expression that the learner minimises in each step, for some convex function  $R : S \rightarrow \mathbb{R}$ . This leads to the *follow-the-regularised-leader* (FTRL) rule:

$$\mathbf{w}_t \in \arg \min_{\mathbf{w} \in S} \left( \sum_{i=1}^{t-1} f_i(\mathbf{w}) + R(\mathbf{w}) \right). \quad (3)$$

The following theorem gives an upper bound on the regret of follow the regularised leader. Notice that the first term in the upper bound corresponds to the stability of the algorithm.

**Theorem 1.** For any number of rounds  $T$  and any  $\mathbf{u} \in S$ :

$$\sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{w}_{t+1})) + R(\mathbf{u}) - R(\mathbf{w}_1).$$

*Proof.* Subtracting  $\sum_{t=1}^T f_t(\mathbf{w}_t)$  from both sides, it is sufficient to show that

$$\sum_{t=1}^T f_t(\mathbf{w}_{t+1}) + R(\mathbf{w}_1) \leq \sum_{t=1}^T f_t(\mathbf{u}) + R(\mathbf{u}) \quad \text{for all } \mathbf{u} \in S \quad (4)$$

We prove (4) by induction on  $T$ . Assume by induction that

$$\sum_{t=1}^{T-1} f_t(\mathbf{w}_{t+1}) + R(\mathbf{w}_1) \leq \sum_{t=1}^{T-1} f_t(\mathbf{u}) + R(\mathbf{u}) \quad \text{for all } \mathbf{u} \in S$$

Adding  $f_T(\mathbf{w}_{T+1})$  to both sides we have

$$\sum_{t=1}^T f_t(\mathbf{w}_{t+1}) + R(\mathbf{w}_1) \leq \sum_{t=1}^{T-1} f_t(\mathbf{u}) + f_T(\mathbf{w}_{T+1}) + R(\mathbf{u}) \quad \text{for all } \mathbf{u} \in S$$

In particular this holds for  $\mathbf{u} = \mathbf{w}_{T+1}$ , that is,

$$\sum_{t=1}^T f_t(\mathbf{w}_{t+1}) + R(\mathbf{w}_1) \leq \sum_{t=1}^T f_t(\mathbf{w}_{T+1}) + R(\mathbf{w}_{T+1}) \quad (5)$$

But since  $\mathbf{w}_{T+1} \in \arg \min_{\mathbf{u} \in S} \left( \sum_{t=1}^T f_t(\mathbf{u}) + R(\mathbf{u}) \right)$  the inequality (5) directly implies (4).  $\square$

### 3 Online Subgradient Descent

In this section we present an algorithm for online convex optimisation that can be seen as an instance of FTRL with regulariser  $R(\mathbf{w}) = \frac{1}{2\eta} \|\mathbf{w}\|^2$ , where  $\eta > 0$  is a parameter of the algorithm.

**Online Subgradient Descent** (with parameter  $\eta > 0$ )

**for**  $t = 1, \dots, T$

Predict  $\mathbf{w}_t := \arg \min_{\mathbf{w} \in S} \left\| \mathbf{w} - \eta \sum_{i=1}^{t-1} \mathbf{z}_i \right\|$

Receive  $g_t : S \rightarrow \mathbb{R}$

Choose a subgradient  $\mathbf{z}_t \in \partial g_t(\mathbf{w}_t)$

**end**

The prediction  $\mathbf{w}_t$  is chosen to be the closest point in  $S$  to the sum of the subgradients  $-\eta \sum_{i=1}^{t-1} \mathbf{z}_i$ . In case  $S = \mathbb{R}^n$  we simply have  $\mathbf{w}_t = -\eta \sum_{i=1}^{t-1} \mathbf{z}_i$ .

**Remark 1.** Suppose that each loss function  $g_t$  is a linear function  $g_t(\mathbf{w}) = \mathbf{v}_t \cdot \mathbf{w}$ . Then the above algorithm makes exactly the same predictions as FTRL with regulariser  $R(\mathbf{w}) = \frac{1}{2\eta} \|\mathbf{w}\|^2$ . Let  $\mathbf{w}_1, \dots, \mathbf{w}_T$  be the predictions on the online subgradient descent algorithm. Observe that the only choice for  $\mathbf{z}_i \in \partial g_i(\mathbf{w}_i)$  is  $\mathbf{z}_i = \mathbf{v}_i$ , and so we have

$$\mathbf{w}_t = \arg \min_{\mathbf{w} \in S} \frac{1}{2\eta} \left\| \mathbf{w} - \eta \sum_{i=1}^{t-1} \mathbf{v}_i \right\|^2. \quad (6)$$

On the other hand, the FTRL update rule (3) with the given sequence of loss functions is

$$\mathbf{w}_t = \arg \min_{\mathbf{w} \in S} \left( \frac{1}{2\eta} \|\mathbf{w}\|^2 - \sum_{i=1}^{t-1} \mathbf{v}_i \cdot \mathbf{w} \right). \quad (7)$$

But, completing the square, we see that the objectives in the two optimisation problems (6) and (7) are equal up to a constant.

The main result of this section shows that the online subgradient descent procedure has total regret  $O(\sqrt{T})$  over  $T$  rounds, assuming that the loss functions  $f_t$  are convex,  $\rho$ -Lipschitz for some fixed  $\rho$ , and defined on a region  $S$  of bounded diameter. This achieves our goal of sublinear regret.

**Theorem 2.** Consider a run of the online subgradient descent algorithm in which the loss functions  $g_1, \dots, g_T : S \rightarrow \mathbb{R}$  are  $\rho$ -Lipschitz and such that  $\sup\{\|\mathbf{u} - \mathbf{v}\| : \mathbf{u}, \mathbf{v} \in S\} \leq D$  for some  $\rho, D > 0$ . If  $\eta = \frac{D}{\rho\sqrt{2T}}$  then the regret after  $T$  rounds is at most  $D\rho\sqrt{2T}$ .

*Proof.* For all  $\mathbf{u} \in S$  we have

$$\sum_{t=1}^T (g_t(\mathbf{w}_t) - g_t(\mathbf{u})) \leq \sum_{t=1}^T \mathbf{z}_t \cdot (\mathbf{w}_t - \mathbf{u}) \quad (8)$$

$$\leq \sum_{t=1}^T \mathbf{z}_t \cdot (\mathbf{w}_t - \mathbf{w}_{t+1}) + \frac{\|\mathbf{u}\|^2 - \|\mathbf{w}_1\|^2}{2\eta} \quad (9)$$

$$\leq \sum_{t=1}^T \|\mathbf{z}_t\| \cdot \eta \|\mathbf{z}_t\| + \frac{\|\mathbf{u}\|^2 - \|\mathbf{w}_1\|^2}{2\eta} \quad (10)$$

$$\leq T\rho^2\eta D + \frac{D^2}{2\eta}.$$

Line (8) holds because  $\mathbf{z}_t$  is a subgradient of  $g_t$  at  $\mathbf{w}_t$ . Line (9) follows from Theorem 9 applied to FTRL with linear loss functions  $f_1, \dots, f_T : S \rightarrow \mathbb{R}$ , with  $f_t(\mathbf{w}) = \mathbf{z}_t \cdot \mathbf{w}$  for  $t = 1, \dots, T$ .<sup>1</sup> Line (10) follows from the Cauchy-Schwartz inequality, using the fact that

$$\|\mathbf{w}_{t+1} - \mathbf{w}_t\| \leq \left\| \eta \sum_{i=1}^t \mathbf{z}_i - \eta \sum_{i=1}^{t-1} \mathbf{z}_i \right\| = \|\eta \mathbf{z}_t\|.$$

Putting  $\eta = \frac{D}{\rho\sqrt{2T}}$ , the right-hand side above evaluates to  $D\rho\sqrt{2T}$ .  $\square$

### 3.1 Perceptron

The Perceptron algorithm can be seen as an instance of the online subgradient descent algorithm with parameter  $\eta = 1$ . Specifically let the loss function  $g_t : \mathbb{R}^n \rightarrow \mathbb{R}$  in round  $t$  of the online subgradient descent algorithm be

$$g_t(\mathbf{w}) = \max(0, -y_t(\mathbf{w} \cdot \mathbf{x}_t)),$$

where  $(\mathbf{x}_t, y_t)$  is the labelled example in round  $t$  of the Perceptron algorithm. Notice that the subgradient  $\mathbf{z}_t$  at round  $t$  satisfies  $\mathbf{z}_t = \mathbf{0}$  if the prediction  $\mathbf{w}_t$  has loss 0 and  $\mathbf{z}_t = -y_t \mathbf{x}_t$  otherwise. Thus the prediction  $\mathbf{w}_t = -\sum_{i=1}^{t-1} \mathbf{z}_i$  exactly matches the update rule of the Perceptron algorithm as presented in Lecture 11.

<sup>1</sup>Note that by Remark 1 FTRL makes predictions  $\mathbf{w}_1, \dots, \mathbf{w}_T$  when given this input by the adversary.

## 4 Winnow

Let us briefly sketch how the Winnow algorithm is based on follow-the-regularised-leader using the entropy function as a regulariser.

Let  $\mathbf{p}$  be a probability distribution of  $\{1, \dots, n\}$ . The entropy function is defined by

$$H(\mathbf{p}) = \sum_{i=1}^n p_i \log \frac{1}{p_i}.$$

Here, by convention,  $0 \log(1/0)$  is defined to be 0. Note that when  $\mathbf{p}$  puts all of its weight on a single feature, we have  $H(\mathbf{p}) = 0$ . When  $\mathbf{p}$  is uniform, we have  $H(\mathbf{p}) = \log n$ . In general, higher entropy implies “more uniform” weights. Because of this, using the negative entropy as the regularization term makes sense. Encouraging high entropy leads to more uniform weights, which leads to more stability in the algorithm.

Consider the FTRL algorithm with regulariser  $-\frac{1}{\eta}H$  for some parameter  $\eta > 0$ . Assume that in round  $t$  the adversary supplies the hinge loss function  $f_t : \mathbb{R}^n \rightarrow \mathbb{R}$  defined by

$$f_t(\mathbf{w}) = \max(0, -y_t(\mathbf{w} \cdot \mathbf{x}_t))$$

where  $(\mathbf{x}_t, y_t)$  is the  $t$ -th input to the Winnow algorithm. As in the case of online subgradient descent we can soundly apply the FTRL rule instead to the “surrogate” linear loss functions

$$g_t(\mathbf{w}) = \mathbf{z}_t \cdot \mathbf{w}$$

where  $\mathbf{z}_t$  is a subgradient of  $f_t$  at  $\mathbf{w}_t$ . Again we have that  $\mathbf{z}_t = \mathbf{0}$  if the prediction  $\mathbf{w}_t$  was correct on  $(\mathbf{x}_t, y_t)$  and  $\mathbf{z}_t = -y_t \mathbf{x}_t$  otherwise.

The prediction  $\mathbf{w}_{t+1}$  of FTRL algorithm in round  $t + 1$  is obtained by the following rule:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w} \in S} \sum_{i=1}^t \mathbf{z}_i \cdot \mathbf{w} - \frac{1}{\eta} H(\mathbf{w}). \quad (11)$$

Solving for  $\mathbf{w}$  using calculus one obtains an update rule equivalent to that given in the Winnow algorithm in Lecture 12.