

# Machine Learning Theory 2022

## Lecture 10

**Wouter M. Koolen**

Download these slides now from [elo.mastermath.nl](http://elo.mastermath.nl)!

### Online Convex Optimisation

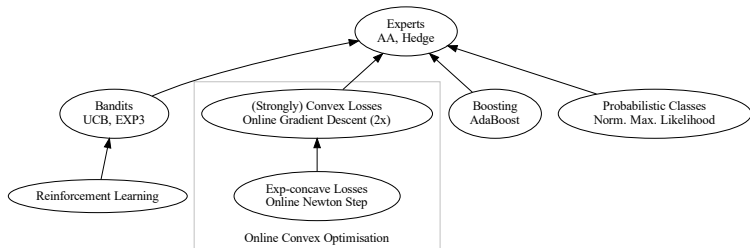
- ▶ Gradient Descent for Convex Losses
- ▶ Online to Batch Conversion
- ▶ Gradient Descent for Strongly Convex Losses



homework roulette  
in the break

## Recap

# Overview of Second Half of Course



Material: course notes and selection of sources on MLT website.

## Recap: Finite Classes

So far we have seen learning “finite sets”:

Our learning algorithms behave like the **best** among  $K$  strategies.

- ▶  $K$ -Experts setting
  - ▶ Mix loss : Aggregating Algorithm
  - ▶ Dot loss : Hedge algorithm
- ▶  $K$ -armed bandit settings
  - ▶ Adversarial bandit : EXP3
  - ▶ Stochastic bandit : UCB

# Outlook: Beyond the Finite

What if we want to compete with **infinite** sets?

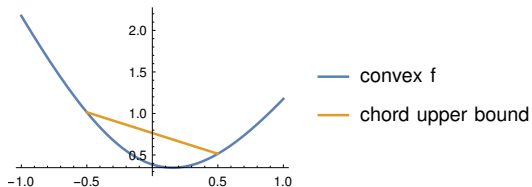
- ▶ Can we?
- ▶ How?

In each case, **lower bounds** grow with  $K$ :  $\ln K$ ,  $\sqrt{T \ln K}$ ,  $\sqrt{TK \ln K}$ ,  $K/\Delta \ln T$ . So hopeless in the **unstructured**  $K \rightarrow \infty$  case.

Today: compete with **continuous** sets of actions, parameterised such that the loss is a **convex** function of the action.

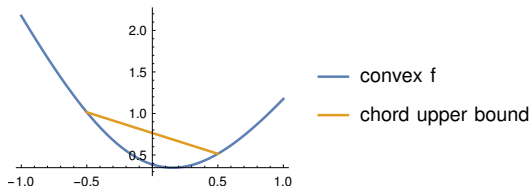
# Convexity Review

# Convex Functions I : definition



Fix a convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ .

# Convex Functions I : definition



Fix a convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ .

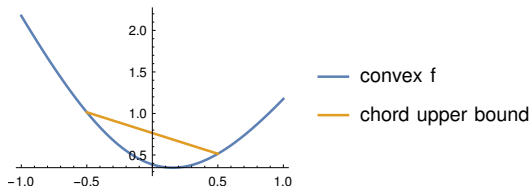
## Definition

A function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is convex if for all  $\mathbf{x}, \mathbf{y} \in \mathcal{U}$  and weights  $\theta \in [0, 1]$ ,

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y}).$$



# Convex Functions I : definition



Fix a convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ .

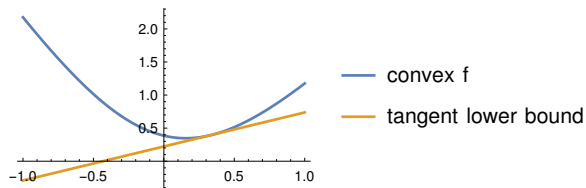
## Definition

A function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is convex if for all  $\mathbf{x}, \mathbf{y} \in \mathcal{U}$  and weights  $\theta \in [0, 1]$ ,

$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y}).$$

Extends to arbitrary mixtures:  $f(\mathbb{E}[X]) \leq \mathbb{E}[f(X)]$  (Jensen).

## Convex Functions II : tangent bound

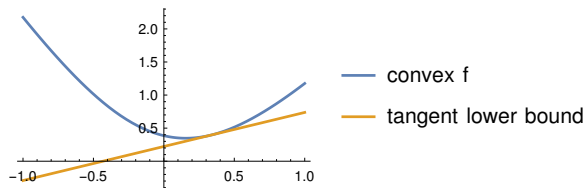


### Fact

A differentiable function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is convex iff for all  $x, y \in \mathcal{U}$

$$f(y) - f(x) \geq \langle y - x, \nabla f(x) \rangle$$

## Convex Functions II : tangent bound



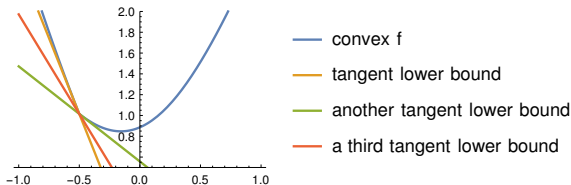
### Fact

A differentiable function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is convex iff for all  $\mathbf{x}, \mathbf{y} \in \mathcal{U}$

$$f(\mathbf{y}) - f(\mathbf{x}) \geq \langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{x}) \rangle$$

Symmetrically,  $\langle \mathbf{y} - \mathbf{x}, \nabla f(\mathbf{y}) \rangle \geq f(\mathbf{y}) - f(\mathbf{x})$ .

# Convex Functions III : sub-gradient



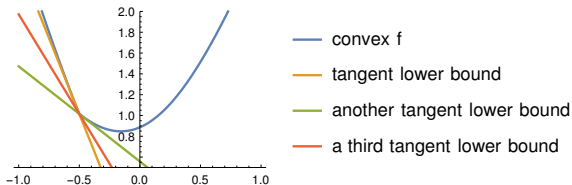
## Fact (Sub-gradient)

For any convex  $f : \mathcal{U} \rightarrow \mathbb{R}$ , possibly non-differentiable, and point  $x \in \mathcal{U}$ , there always exists **some** vector  $g \in \mathbb{R}^d$  such that for all  $y \in \mathcal{U}$

$$f(y) - f(x) \geq \langle y - x, g \rangle$$

Any such vector  $g$  is called a **sub-gradient** (of  $f$  at  $x$ ).

# Convex Functions III : sub-gradient



## Fact (Sub-gradient)

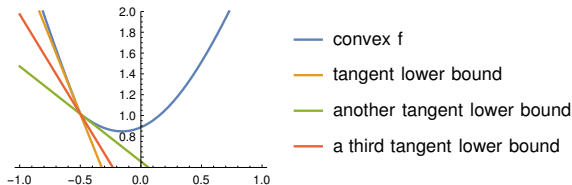
For any convex  $f : \mathcal{U} \rightarrow \mathbb{R}$ , possibly non-differentiable, and point  $x \in \mathcal{U}$ , there always exists **some** vector  $g \in \mathbb{R}^d$  such that for all  $y \in \mathcal{U}$

$$f(y) - f(x) \geq \langle y - x, g \rangle$$

Any such vector  $g$  is called a **sub-gradient** (of  $f$  at  $x$ ).

The gradient of a differentiable function is a sub-gradient.

# Convex Functions III : sub-gradient



## Fact (Sub-gradient)

For any convex  $f : \mathcal{U} \rightarrow \mathbb{R}$ , possibly non-differentiable, and point  $x \in \mathcal{U}$ , there always exists **some** vector  $g \in \mathbb{R}^d$  such that for all  $y \in \mathcal{U}$

$$f(y) - f(x) \geq \langle y - x, g \rangle$$

Any such vector  $g$  is called a **sub-gradient** (of  $f$  at  $x$ ).

The gradient of a differentiable function is a sub-gradient.

We will abuse notation and denote **any** sub-gradient by  $\nabla f(x)$ .

# Online Convex Optimisation

# Online Convex Optimisation

General yet simple sequential decision problem.

Fix a convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ .

## Protocol

For  $t = 1, 2, \dots$

- ▶ Learner chooses a point  $\mathbf{w}_t \in \mathcal{U}$ .
- ▶ Adversary reveals convex loss function  $f_t : \mathcal{U} \rightarrow \mathbb{R}$ .
- ▶ Learner's loss is  $f_t(\mathbf{w}_t)$



# Online Convex Optimisation

General yet simple sequential decision problem.

Fix a convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ .

## Protocol

For  $t = 1, 2, \dots$

- ▶ Learner chooses a point  $\mathbf{w}_t \in \mathcal{U}$ .
- ▶ Adversary reveals convex loss function  $f_t : \mathcal{U} \rightarrow \mathbb{R}$ .
- ▶ Learner's loss is  $f_t(\mathbf{w}_t)$

## Objective:

Regret w.r.t. best point after  $T$  rounds:

$$R_T = \max_{\mathbf{u} \in \mathcal{U}} \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u}))$$

## Example loss functions

| Setting              | loss function $f_t(\mathbf{u})$                   |
|----------------------|---|
| Hedge setting        | $\mathbf{u}^\top \ell_t$                          |
| Point prediction     | $\ \mathbf{u} - \mathbf{x}_t\ ^2$                 |
| Regression           | $(\mathbf{u}^\top \mathbf{x}_t - y_t)^2$          |
| Logistic regression  | $\ln(1 + e^{-y_t \mathbf{u}^\top \mathbf{x}_t})$  |
| Hinge loss           | $\max\{0, 1 - y_t \mathbf{u}^\top \mathbf{x}_t\}$ |
| Investment           | $-\ln(\mathbf{u}^\top \mathbf{x}_t)$              |
| Offline optimisation | $f(\mathbf{u})$                                   |

# Online Gradient Descent (OGD)

Let  $\mathcal{U}$  be a closed convex set containing  $\mathbf{0}$ .

## Definition

Online Gradient Descent with learning rate  $\eta > 0$  plays

$$\mathbf{w}_1 = \mathbf{0} \quad \text{and} \quad \mathbf{w}_{t+1} = \Pi_{\mathcal{U}}(\mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t))$$

where  $\Pi_{\mathcal{U}}(\mathbf{w}) = \arg \min_{\mathbf{u} \in \mathcal{U}} \|\mathbf{u} - \mathbf{w}\|$  is the projection onto  $\mathcal{U}$ .

# Online Gradient Descent (OGD)

## Theorem

Let  $\|\nabla f_t(\mathbf{u})\| \leq G$  and  $\|\mathbf{u}\| \leq D$  for all  $\mathbf{u} \in \mathcal{U}$ . Then

$$R_T = \max_{\mathbf{u} \in \mathcal{U}} \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

# Online Gradient Descent (OGD)

## Theorem

Let  $\|\nabla f_t(\mathbf{u})\| \leq G$  and  $\|\mathbf{u}\| \leq D$  for all  $\mathbf{u} \in \mathcal{U}$ . Then

$$R_T = \max_{\mathbf{u} \in \mathcal{U}} \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) \leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} T G^2$$

## Corollary

Tuning  $\eta = \frac{D}{G\sqrt{T}}$  results in

$$R_T \leq D G \sqrt{T}$$

# Pythagorean Inequality

## Lemma (Pythagorean Inequality)

Fix a closed convex set  $\mathcal{U} \subseteq \mathbb{R}^d$ . Let  $\mathbf{x} \in \mathcal{U}$ ,  $\mathbf{y} \in \mathbb{R}^d$  and

$$\hat{\mathbf{y}} = \Pi_{\mathcal{U}}(\mathbf{y}) = \arg \min_{\mathbf{u} \in \mathcal{U}} \|\mathbf{u} - \mathbf{y}\|^2.$$

Then

$$\|\mathbf{x} - \hat{\mathbf{y}}\|^2 + \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \leq \|\mathbf{x} - \mathbf{y}\|^2$$

NB: not to be confused with **triangle inequality**

$$\|\mathbf{x} - \mathbf{y}\| \leq \|\mathbf{x} - \hat{\mathbf{y}}\| + \|\hat{\mathbf{y}} - \mathbf{y}\|.$$

# Proof of GD regret bound I

Fix any  $\mathbf{u} \in \mathcal{U}$ . We have

$$f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle$$

Moreover,

$$\begin{aligned} \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 &= \|\Pi_{\mathcal{U}}(\mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t)) - \mathbf{u}\|^2 \\ &\stackrel{\text{Pyth. Ineq.}}{\leq} \|\mathbf{w}_t - \eta \nabla f_t(\mathbf{w}_t) - \mathbf{u}\|^2 \\ &= \|\mathbf{w}_t - \mathbf{u}\|^2 - 2\eta \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \eta^2 \|\nabla f_t(\mathbf{w}_t)\|^2 \end{aligned}$$

Hence

$$\langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle \leq \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta} + \frac{\eta}{2} \|\nabla f_t(\mathbf{w}_t)\|^2$$

## Proof of GD regret bound II

Summing over  $T$  rounds, we find

$$\begin{aligned}\sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u})) &\leq \sum_{t=1}^T \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle \\ &\leq \underbrace{\sum_{t=1}^T \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta}}_{\text{telescopes}} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{w}_t)\|^2 \\ &\leq \frac{\|\mathbf{u}\|^2 - \|\mathbf{w}_{T+1} - \mathbf{u}\|^2}{2\eta} + \frac{\eta}{2} \sum_{t=1}^T \|\nabla f_t(\mathbf{w}_t)\|^2 \\ &\leq \frac{D^2}{2\eta} + \frac{\eta}{2} TG^2\end{aligned}$$



## Online to Batch Conversion

## Online to Batch Conversion

Goal: obtain an estimator  $\hat{w}_T$  with small expected excess risk.

$$\mathbb{E}_{f_1, \dots, f_T} \left[ \mathbb{E}_f [f(\hat{w}_T) - f(u^*)] \right] \leq \text{small}$$

where the training set  $f_1, \dots, f_T$  and the test sample  $f$  are drawn i.i.d. and  $u^*$  optimises the risk  $u \mapsto \mathbb{E}_f[f(u)]$ .

## Online to Batch Conversion

Goal: obtain an estimator  $\hat{w}_T$  with small expected excess risk.

$$\mathbb{E}_{f_1, \dots, f_T} \left[ \mathbb{E}_f [f(\hat{w}_T) - f(u^*)] \right] \leq \text{small}$$

where the training set  $f_1, \dots, f_T$  and the test sample  $f$  are drawn i.i.d. and  $u^*$  optimises the risk  $u \mapsto \mathbb{E}_f[f(u)]$ .

Idea: use online learning algorithm. Given training sample  $f_1, \dots, f_T$ , the algorithm picks  $w_1, \dots, w_T$ . Let us define the *average iterate estimator*

$$\hat{w}_T = \frac{1}{T} \sum_{t=1}^T w_t.$$

# Online to Batch Conversion

Goal: obtain an estimator  $\hat{w}_T$  with small expected excess risk.

$$\mathbb{E}_{f_1, \dots, f_T} \left[ \mathbb{E}_f [f(\hat{w}_T) - f(u^*)] \right] \leq \text{small}$$

where the training set  $f_1, \dots, f_T$  and the test sample  $f$  are drawn i.i.d. and  $u^*$  optimises the risk  $u \mapsto \mathbb{E}_f[f(u)]$ .

Idea: use online learning algorithm. Given training sample  $f_1, \dots, f_T$ , the algorithm picks  $w_1, \dots, w_T$ . Let us define the *average iterate estimator*

$$\hat{w}_T = \frac{1}{T} \sum_{t=1}^T w_t.$$

## Theorem

*An online regret bound  $R_T \leq B(T)$  implies*

$$\mathbb{E}_{\text{iid } f_1, \dots, f_T, f} [f(\hat{w}_T) - f(u^*)] \leq \frac{B(T)}{T}$$

# Online to Batch Proof

$$\begin{aligned} & \mathbb{E}_{\text{iid } f_1, \dots, f_T, f} [f(\hat{\mathbf{w}}_T) - f(\mathbf{u}^*)] \\ & \leq \mathbb{E}_{\text{iid } f_1, \dots, f_T, f} \left[ \frac{1}{T} \sum_{t=1}^T (f(\mathbf{w}_t) - f(\mathbf{u}^*)) \right] \\ & = \mathbb{E}_{\text{iid } f_1, \dots, f_T, f} \left[ \frac{1}{T} \sum_{t=1}^T (f_t(\mathbf{w}_t) - f_t(\mathbf{u}^*)) \right] \leq \frac{B(T)}{T} \end{aligned}$$

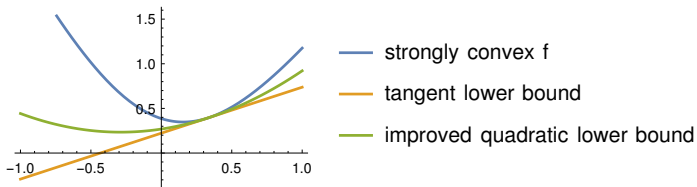
The first step is convexity of  $f$ . The last step uses that  $f$  and  $f_t$  have the same distribution (and  $\mathbf{w}_t$  is not a function of  $f_t$ ).

# Online Strongly Convex Optimisation

# Structure

What if I **know more** about my setting than **convexity of the loss function**? Can I learn faster?

# Strongly Convex Case



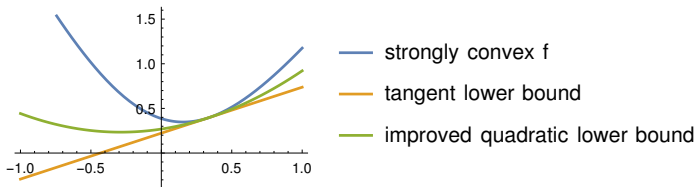
## Definition

A function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is *strongly convex* to degree  $\alpha \geq 0$  if

$$f(\mathbf{u}) - f(\mathbf{w}) \geq \langle \mathbf{u} - \mathbf{w}, \nabla f(\mathbf{w}) \rangle + \frac{\alpha}{2} \|\mathbf{u} - \mathbf{w}\|^2$$



# Strongly Convex Case



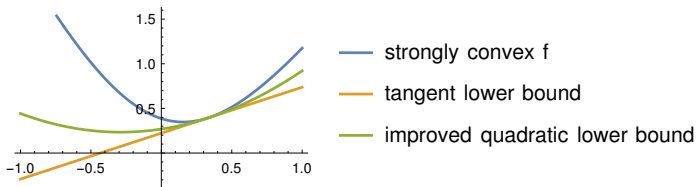
## Definition

A function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is *strongly convex* to degree  $\alpha \geq 0$  if

$$f(\mathbf{u}) - f(\mathbf{w}) \geq \langle \mathbf{u} - \mathbf{w}, \nabla f(\mathbf{w}) \rangle + \frac{\alpha}{2} \|\mathbf{u} - \mathbf{w}\|^2$$

Example:  $f(\mathbf{w}) = \|\mathbf{w} - \mathbf{x}_t\|^2$ .

# Strongly Convex Case



## Definition

A function  $f : \mathcal{U} \rightarrow \mathbb{R}$  is *strongly convex* to degree  $\alpha \geq 0$  if

$$f(\mathbf{u}) - f(\mathbf{w}) \geq \langle \mathbf{u} - \mathbf{w}, \nabla f(\mathbf{w}) \rangle + \frac{\alpha}{2} \|\mathbf{u} - \mathbf{w}\|^2$$

Example:  $f(\mathbf{w}) = \|\mathbf{w} - \mathbf{x}_t\|^2$ .

Idea: could this extra knowledge help in the regret rate?

# Online Gradient Descent with time-varying learning rate

Definition (OGD with time-varying learning rate)

$$\mathbf{w}_1 = \mathbf{0} \quad \text{and} \quad \mathbf{w}_{t+1} = \Pi_{\mathcal{U}}(\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t))$$

# Online Gradient Descent with time-varying learning rate

## Definition (OGD with time-varying learning rate)

$$\mathbf{w}_1 = \mathbf{0} \quad \text{and} \quad \mathbf{w}_{t+1} = \Pi_{\mathcal{U}}(\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t))$$

## Theorem

*For  $\alpha$ -strongly convex loss functions, OGD with learning rate  $\eta_t = \frac{1}{\alpha t}$  ensures*

$$R_T \leq \frac{G^2}{2\alpha} (1 + \ln T).$$

# Proof I

We start with

$$\begin{aligned}\|\mathbf{w}_{t+1} - \mathbf{u}\|^2 &= \|\Pi_{\mathcal{U}}(\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t)) - \mathbf{u}\|^2 \\ &\stackrel{\text{Pyth. Ineq.}}{\leq} \|\mathbf{w}_t - \eta_t \nabla f_t(\mathbf{w}_t) - \mathbf{u}\|^2 \\ &= \|\mathbf{w}_t - \mathbf{u}\|^2 - 2\eta_t \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle + \eta_t^2 \|\nabla f_t(\mathbf{w}_t)\|^2\end{aligned}$$

So that

$$\begin{aligned}&f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \\ &\leq \langle \mathbf{w}_t - \mathbf{u}, \nabla f_t(\mathbf{w}_t) \rangle - \frac{\alpha}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \\ &\leq \frac{\|\mathbf{w}_t - \mathbf{u}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{u}\|^2 + \eta_t^2 \|\nabla f_t(\mathbf{w}_t)\|^2}{2\eta_t} - \frac{\alpha}{2} \|\mathbf{w}_t - \mathbf{u}\|^2 \\ &= \|\mathbf{w}_t - \mathbf{u}\|^2 \left( \frac{1}{2\eta_t} - \frac{\alpha}{2} \right) - \frac{\|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta_t} + \frac{\eta_t \|\nabla f_t(\mathbf{w}_t)\|^2}{2}\end{aligned}$$

## Proof II

Summing over rounds gives

$$\begin{aligned} & \sum_{t=1}^T f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \\ & \leq \sum_{t=1}^T \left( \|\mathbf{w}_t - \mathbf{u}\|^2 \left( \frac{1}{2\eta_t} - \frac{\alpha}{2} \right) - \frac{\|\mathbf{w}_{t+1} - \mathbf{u}\|^2}{2\eta_t} + \frac{\eta_t \|\nabla f_t(\mathbf{w}_t)\|^2}{2} \right) \\ & = \|\mathbf{w}_1 - \mathbf{u}\|^2 \left( \frac{1}{2\eta_1} - \frac{\alpha}{2} \right) + \sum_{t=2}^T \|\mathbf{w}_t - \mathbf{u}\|^2 \left( \frac{1}{2\eta_t} - \frac{\alpha}{2} - \frac{1}{2\eta_{t-1}} \right) \\ & \quad - \frac{\|\mathbf{w}_{T+1} - \mathbf{u}\|^2}{2\eta_T} + \sum_{t=1}^T \frac{\eta_t \|\nabla f_t(\mathbf{w}_t)\|^2}{2} \end{aligned}$$

Key idea for telescoping is to cancel coefficient on  $\|\mathbf{w}_t - \mathbf{u}\|^2$  in the sum:

$$\frac{1}{\eta_{t+1}} - \alpha = \frac{1}{\eta_t}$$

## Proof III

So

$$\eta_{t+1} = \frac{1}{\frac{1}{\eta_t} + \alpha}$$

A good starting point (cancelling the first term) is  $\eta_1 = \frac{1}{\alpha}$ . This leads to  $\eta_t = \frac{1}{\alpha t}$ . We then find

$$\sum_{t=1}^T f_t(\mathbf{w}_t) - f_t(\mathbf{u}) \leq \sum_{t=1}^T \frac{\|\nabla f_t(\mathbf{w}_t)\|^2}{2\alpha t} \leq \frac{G^2}{2\alpha} (1 + \ln T)$$

# Conclusion

Tools for learning in convex settings.

- ▶ Guaranteed robustness against adversarial losses
- ▶ Efficient
- ▶ Building block for
  - ▶ Learning in non-convex settings (AdaGrad for DNN)
  - ▶ Learning in games
  - ▶ Non-convex games (GANs)
  - ▶ ...