

# Convex optimization

CS 446 / ECE 449

2022-02-01 04:33:47 -0600 (6fe12f5)

# Why are we learning about convexity?

## We've already seen it twice:

- ▶ If  $\hat{\mathcal{R}}$  is convex, then any  $\mathbf{w}$  with  $\nabla \hat{\mathcal{R}}(\mathbf{w})$  is **globally optimal**:  $\hat{\mathcal{R}}(\mathbf{w}) = \inf_{\mathbf{v}} \hat{\mathcal{R}}(\mathbf{v})$ .  
Matches the normal equations!
- ▶ Gradient descent on convex  $\hat{\mathcal{R}}$  is guaranteed to work.

Convexity is pervasive throughout mathematics.

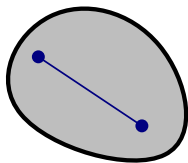
**Warning:**  $\hat{\mathcal{R}}$  is **not** convex for deep networks.

# Plan for today

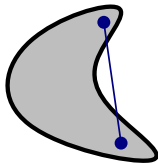
- ▶ Convex sets and functions.
- ▶ Minimizing convex functions.

# Convex sets

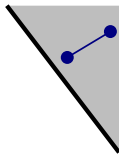
A set  $S$  is convex if, for every pair of points  $\{x, x'\}$  in  $S$ , the line segment between  $x$  and  $x'$  is also contained in  $S$ .



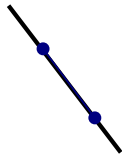
convex



not convex



convex



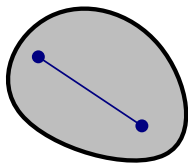
convex

In symbols:

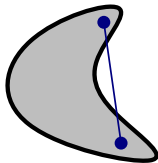
$$\{x, x'\} \in S \quad \implies \quad [x, x'] = \{\alpha x + (1 - \alpha)x' : \alpha \in [0, 1]\} \subseteq S.$$

# Convex sets

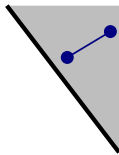
A set  $S$  is convex if, for every pair of points  $\{\mathbf{x}, \mathbf{x}'\}$  in  $S$ , the line segment between  $\mathbf{x}$  and  $\mathbf{x}'$  is also contained in  $S$ .



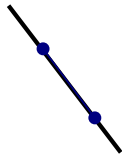
convex



not convex



convex



convex

In symbols:

$$\{\mathbf{x}, \mathbf{x}'\} \in S \quad \Longrightarrow \quad [\mathbf{x}, \mathbf{x}'] = \{\alpha \mathbf{x} + (1 - \alpha) \mathbf{x}' : \alpha \in [0, 1]\} \subseteq S.$$

## Examples:

- ▶ All of  $\mathbb{R}^d$ .
- ▶ Empty set.
- ▶ Half-spaces:  $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{a}^\top \mathbf{x} \leq b\}$ .
- ▶ Intersections of convex sets.
- ▶ Polyhedra:  $\{\mathbf{x} \in \mathbb{R}^d : \mathbf{A}\mathbf{x} \leq \mathbf{b}\} = \bigcap_{i=1}^m \{\mathbf{x} \in \mathbb{R}^d : \mathbf{a}_i^\top \mathbf{x} \leq b_i\}$ .
- ▶ Convex hulls:  $\text{conv}(S) := \{\sum_{i=1}^k \alpha_i \mathbf{x}_i : k \in \mathbb{N}, \mathbf{x}_i \in S, \alpha_i \geq 0, \sum_{i=1}^k \alpha_i = 1\}$ .

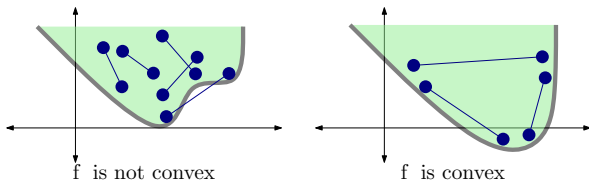
Let's verify that any halfspace  $H := \{x \in \mathbb{R}^d : a^\top x \leq b\}$  is convex.

# Convex functions from convex sets

The [epigraph](#) of a function  $f$  is the region above the curve:

$$\text{epi}(f) := \left\{ (\mathbf{x}, y) \in \mathbb{R}^{d+1} : y \geq f(\mathbf{x}) \right\}.$$

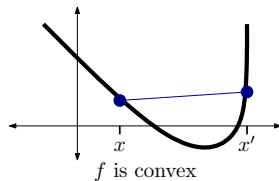
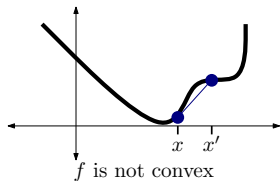
A function is convex if its epigraph is a convex set.



# Convex functions (standard definitions)

(Secants lie above.)  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ ,

$$f((1 - \alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1 - \alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$

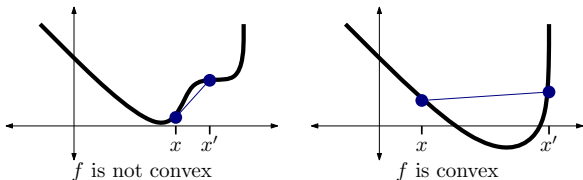




# Convex functions (standard definitions)

(Secants lie above.)  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ ,

$$f((1 - \alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1 - \alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$



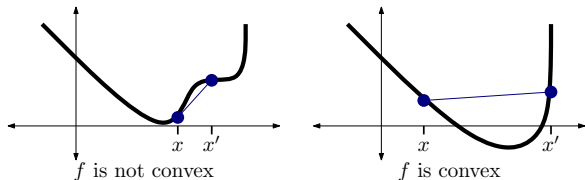
## Other characterizations.

- ▶ (Tangents lie below.) Given  $\mathbf{x}, \mathbf{z}$ , then  $f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x})$ .
- ▶ (Increasing slopes.) Given  $\mathbf{x}, \mathbf{z}$ , then  $(\nabla f(\mathbf{x}) - \nabla f(\mathbf{z}))^\top (\mathbf{x} - \mathbf{z}) \geq 0$ .
- ▶ (Curves upwards.) Given  $\mathbf{x}$ , then  $\nabla^2 f(\mathbf{x}) \succeq 0$ .

# Convex functions (standard definitions)

(Secants lie above.)  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if for any  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and  $\alpha \in [0, 1]$ ,

$$f((1 - \alpha)\mathbf{x} + \alpha\mathbf{x}') \leq (1 - \alpha) \cdot f(\mathbf{x}) + \alpha \cdot f(\mathbf{x}').$$



## Other characterizations.

- ▶ (Tangents lie below.) Given  $\mathbf{x}, \mathbf{z}$ , then  $f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{z} - \mathbf{x})$ .
- ▶ (Increasing slopes.) Given  $\mathbf{x}, \mathbf{z}$ , then  $(\nabla f(\mathbf{x}) - \nabla f(\mathbf{z}))^\top (\mathbf{x} - \mathbf{z}) \geq 0$ .
- ▶ (Curves upwards.) Given  $\mathbf{x}$ , then  $\nabla^2 f(\mathbf{x}) \succeq 0$ .

## Examples.

- ▶  $f(\mathbf{x}) = \mathbf{b}^\top \mathbf{x}$  for any  $\mathbf{b} \in \mathbb{R}^d$ .
- ▶  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$  for symmetric positive semidefinite  $\mathbf{A}$ .
- ▶  $f(\mathbf{x}) = \ln \left( \sum_{i=1}^d \exp(x_i) \right)$ , which approximates  $\max_i x_i$ .

Is  $f(\mathbf{x}) := \frac{1}{2}\mathbf{x}^\top \mathbf{M}\mathbf{x}$  convex? Does the answer depend on dimension?

# Operations preserving convexity

**Summations:** if  $(f_1, \dots, f_k)$  convex and  $(\alpha_1, \dots, \alpha_k)$  nonnegative,

$$\mathbf{x} \mapsto \alpha_1 f_1(\mathbf{x}) + \dots + \alpha_k f_k(\mathbf{x}) \quad \text{is convex.}$$

**Affine composition:** if  $f$  is convex, then for any  $\mathbf{A} \in \mathbb{R}^{m \times d}$  and  $\mathbf{b} \in \mathbb{R}^m$ ,

$$\mathbf{x} \mapsto f(\mathbf{Ax} + \mathbf{b}) \quad \text{is convex.}$$

**Maxima:** if  $(f_1, \dots, f_k)$  are convex,

$$\mathbf{x} \mapsto \max_i f_i(\mathbf{x}) \quad \text{is convex.}$$

## Examples: losses and empirical risks.

Our standard losses are convex.

- ▶ **(Logistic loss.)**  $\ell_{\log}(z) := \ln(1 + \exp(-z))$  is convex.
- ▶ **(Squared loss.)**  $\ell_1(z) := \frac{1}{2}(1 - z)^2$  is convex in  $z$ ,  
 $\ell_2(y, \hat{y}) := \frac{1}{2}(y - \hat{y})^2$  is convex in  $\hat{y}$ .

## Examples: losses and empirical risks.

### Our standard losses are convex.

- ▶ **(Logistic loss.)**  $\ell_{\log}(z) := \ln(1 + \exp(-z))$  is convex.
- ▶ **(Squared loss.)**  $\ell_1(z) := \frac{1}{2}(1 - z)^2$  is convex in  $z$ ,  
 $\ell_2(y, \hat{y}) := \frac{1}{2}(y - \hat{y})^2$  is convex in  $\hat{y}$ .

### Empirical risks.

- ▶ If  $\ell$  is convex over  $\mathbb{R}$ , then  $\widehat{\mathcal{R}}(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \ell(y_i \mathbf{x}_i^T \mathbf{w})$  is convex over  $\mathbb{R}^d$ .
- ▶ If  $\ell$  is convex, and we use a nonlinear predictor  $f_{\mathbf{w}} : \mathbb{R}^d \rightarrow \mathbb{R}$ ,  
then  $\widehat{\mathcal{R}}(f_{\mathbf{w}}) := \frac{1}{n} \sum_{i=1}^n \ell(y_i f_{\mathbf{w}}(\mathbf{x}_i))$  may be non-convex.

# Convexity and optimization

Gradient descent behaves nicely for convex functions.

Many principles of its behavior seem to carry over to deep network optimization.

**Fact.**

If  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex and differentiable, then  $x$  is globally optimal iff  $\nabla f(x) = 0$ .

**Implication:**  $w$  satisfies the normal equations iff  $w$  is globally optimal for  $\hat{\mathcal{R}}_{\text{sq}}$ .



# Differentiability

Many convex functions are not differentiable

(e.g.,  $\mathbf{x} \mapsto \max_i x_i$ , or  $z \mapsto |z|$ , or  $\text{ReLU}(z) := \max\{0, z\}$ .)

**(Tangents lie below.)** Given  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  and  $\mathbf{x}$ , subdifferential set  $\partial f(\mathbf{x})$  is

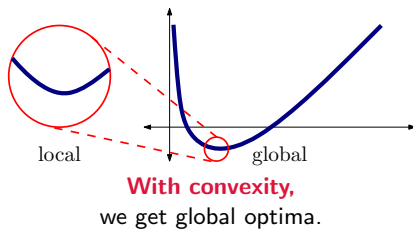
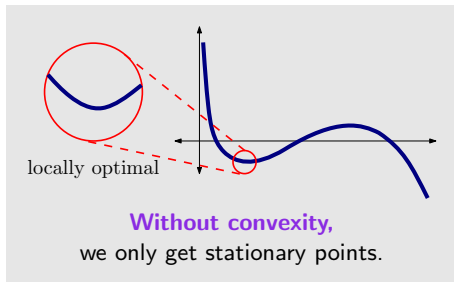
$$\partial f(\mathbf{x}) := \left\{ \mathbf{s} \in \mathbb{R}^d : \forall \mathbf{z} \in \mathbb{R}^d, f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{s}^\top (\mathbf{z} - \mathbf{x}) \right\}.$$

## Properties.

- ▶ If  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, then  $\partial f$  is nonempty everywhere.
- ▶  $\mathbf{x}$  globally optimizes  $f$  iff  $0 \in \partial f(\mathbf{x})$ .
- ▶ Gives an easy proof of **Jensen's inequality**:  $f(\mathbb{E}X) \leq \mathbb{E}f(X)$  for convex  $f$ .

# Behavior of gradient descent

Recall that gradient descent rolls down hills:  $x' := x - \eta \nabla f(x)$ .



**Local-to-global principle:** both cases pass near stationary points, but for convex functions this implies optimality.

# Convergence rates for gradient descent

**Additional assumption:** suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is “ $\beta$ -smooth” : for any  $w, w'$ ,

$$f(w') \leq f(w) + \nabla f(w)^\top (w' - w) + \frac{\beta}{2} \|w' - w\|^2.$$

(In words:  $f$  never curves upwards faster than a quadratic.)

# Convergence rates for gradient descent

**Additional assumption:** suppose  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is “ $\beta$ -smooth” : for any  $\mathbf{w}, \mathbf{w}'$ ,

$$f(\mathbf{w}') \leq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{\beta}{2} \|\mathbf{w}' - \mathbf{w}\|^2.$$

(In words:  $f$  never curves upwards faster than a quadratic.)

**Gradient descent convergence rates for  $\beta$ -smooth  $f$  with  $\eta := \frac{1}{\beta}$ .**

1. Without convexity,

$$\min_{i \leq t} \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{2\beta}{t} \left( f(\mathbf{w}_0) - \inf_{\mathbf{w}} f(\mathbf{w}) \right).$$

2. With convexity, for any  $\mathbf{u} \in \mathbb{R}^d$ ,

$$f(\mathbf{w}_t) - f(\mathbf{u}) \leq \frac{\beta}{2t} \left( \|\mathbf{w}_0 - \mathbf{u}\|^2 - \|\mathbf{w}_t - \mathbf{u}\|^2 \right).$$

Proofs are in appendix.

# Stochastic gradients

Gradient requires time  $\mathcal{O}(n)$ :

$$\nabla \widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell'(y_i f(\mathbf{x}_i; \mathbf{w})) y_i \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}).$$

Instead use a **stochastic minibatch**  $S$  of size  $|S| = B$ :

$$\frac{1}{B} \sum_{i \in S} \ell'(y_i f(\mathbf{x}_i; \mathbf{w})) y_i \nabla_{\mathbf{w}} f(\mathbf{x}_i; \mathbf{w}).$$

- ▶ In the convex case, can still prove convergence rates.
- ▶ Increasing  $B$  improves gradient accuracy;  
in practice, hardware issues like memory size dictate  $B$ .  
Common choices in deep learning are [256, 2048].

## Other practical considerations

There are many other gradient-based (“first-order”) methods.  
`torch.optim` contains many variants (e.g., Adam).

## Other practical considerations

There are many other gradient-based (“first-order”) methods.  
`torch.optim` contains many variants (e.g., Adam).

Deep networks are not differentiable.

They are not convex, so we can't even use subgradients.

Pytorch does not correctly compute gradients, even for convex functions.

No one cares, since `torch.optim` seems to work.

A **convex program** is minimization of a convex function  $f$  over a convex set  $S$ :

$$\min_{\mathbf{x} \in S} f(\mathbf{x}).$$

- ▶ Gradient methods can be adjusted to work well here.
- ▶ This used to be a big research area, but deep learning doesn't make any serious use of this formulation or its tools.



## Summary for today

- ▶ Convex sets and functions.
- ▶ Minimizing convex functions.

(Appendix.)

(Didn't mention:  $f$  is concave means  $-f$  is convex.)

Convex functions can lack curvature:  $x \mapsto x$  and  $x \mapsto |x|$  are convex.

There are two standard notions of curved convex functions:

- ▶ **strict convexity**: no flat regions;
- ▶ **strong convexity**: more curved than a quadratic, everywhere.

## Strict convexity

Function values:  $\forall \mathbf{x}, \mathbf{y}, \forall \alpha \in [0, 1]$ :

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

Derivatives:  $\forall \mathbf{x}, \mathbf{y}$ ,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

Hessians:  $\forall \mathbf{x}$ ,

$$\nabla^2 f(\mathbf{x}) \succeq 0.$$

# Strict convexity

Function values:  $\forall \mathbf{x} \neq \mathbf{y}, \forall \alpha \in (0, 1)$ :

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) < \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}).$$

Derivatives:  $\forall \mathbf{x} \neq \mathbf{y}$ ,

$$f(\mathbf{y}) > f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

Hessians:  $\forall \mathbf{x}$ ,

$$\nabla^2 f(\mathbf{x}) \succ 0.$$

## $\lambda$ -Strong-Convexity.

Function values:  $\forall \mathbf{x}, \mathbf{y}, \forall \alpha \in [0, 1]$

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

Derivatives:  $\forall \mathbf{x}, \mathbf{y}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}).$$

Hessians:  $\forall \mathbf{x},$

$$\nabla^2 f(\mathbf{x}) \succeq 0.$$

## $\lambda$ -Strong-Convexity.

Function values:  $\forall \mathbf{x}, \mathbf{y}, \forall \alpha \in [0, 1]$

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) - \frac{\lambda \alpha (1 - \alpha)}{2} \|\mathbf{x} - \mathbf{y}\|^2.$$

Derivatives:  $\forall \mathbf{x}, \mathbf{y}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{\lambda}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

Hessians:  $\forall \mathbf{x}$ ,

$$\nabla^2 f(\mathbf{x}) \succeq \lambda \mathbf{I}.$$



## Convexity of key losses.

**Logistic loss**  $z \mapsto \ln(1 + \exp(-z))$  is **strictly convex**.  
(e.g., verify that second derivative is positive.)

## Convexity of key losses.

**Logistic loss**  $z \mapsto \ln(1 + \exp(-z))$  is **strictly convex**.  
(e.g., verify that second derivative is positive.)

**Squared loss**  $z \mapsto \frac{1}{2}(1 - z)^2$  is **1-strongly-convex**.  
(e.g., second derivative is 1.)

## Convexity of key losses.

**Logistic loss**  $z \mapsto \ln(1 + \exp(-z))$  is **strictly convex**.  
(e.g., verify that second derivative is positive.)

**Squared loss**  $z \mapsto \frac{1}{2}(1 - z)^2$  is **1-strongly-convex**.  
(e.g., second derivative is 1.)

Combined with our earlier linear prediction calculation,  
logistic regression and least squares are convex!

# GD for smooth, non-convex functions.

**Theorem** (smoothness implies stationary points).

## GD for smooth, non-convex functions.

**Theorem (smoothness implies stationary points).**

Let  $\mathbf{w}_0$  be given, and  $\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla f(\mathbf{w}_i)$ .

If  $f$  is  $\beta$ -smooth and  $\eta = 1/\beta$ ,

$$\min_{i \leq t} \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{1}{t} \sum_{i=1}^t \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{2\beta}{t} \left( f(\mathbf{w}_0) - \min_{\mathbf{w}} f(\mathbf{w}) \right).$$

## GD for smooth, non-convex functions.

**Theorem (smoothness implies stationary points).**

Let  $\mathbf{w}_0$  be given, and  $\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla f(\mathbf{w}_i)$ .

If  $f$  is  $\beta$ -smooth and  $\eta = 1/\beta$ ,

$$\min_{i \leq t} \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{1}{t} \sum_{i=1}^t \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{2\beta}{t} \left( f(\mathbf{w}_0) - \min_{\mathbf{w}} f(\mathbf{w}) \right).$$

**Proof.** Combining the definitions with choice of iterates gives (for each  $i \leq t$ )

$$\begin{aligned} f(\mathbf{w}_i) &\leq f(\mathbf{w}_{i-1}) + \nabla f(\mathbf{w}_{i-1})^\top (\mathbf{w}_i - \mathbf{w}_{i-1}) + \frac{\beta}{2} \|\mathbf{w}_i - \mathbf{w}_{i-1}\|^2 \\ &= f(\mathbf{w}_{i-1}) - \frac{1}{2\beta} \|\nabla f(\mathbf{w}_{i-1})\|^2. \end{aligned}$$

Averaging these inequalities (over  $i \leq t$ ) gives

$$\frac{1}{t} \sum_{i=1}^t \|\nabla f(\mathbf{w}_{i-1})\|^2 \leq \frac{2\beta}{t} (f(\mathbf{w}_0) - f(\mathbf{w}_t)).$$

□

## GD for smooth, convex functions.

**Theorem** (smoothness and convexity imply global optimality).

## GD for smooth, convex functions.

**Theorem** (smoothness and convexity imply global optimality).

Let  $\mathbf{w}_0$  be given, and  $\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla f(\mathbf{w}_i)$ .

If  $f$  is convex and  $\beta$ -smooth, and  $\eta = 1/\beta$ , then  $\forall \mathbf{u} \in \mathbb{R}^d$ ,

$$f(\mathbf{w}_t) - f(\mathbf{u}) \leq \frac{1}{t} \sum_{i=1}^t (f(\mathbf{w}_i) - f(\mathbf{u})) \leq \frac{\beta}{2t} (\|\mathbf{w}_0 - \mathbf{u}\|^2 - \|\mathbf{w}_t - \mathbf{u}\|^2).$$



## GD for smooth, convex functions.

**Theorem** (smoothness and convexity imply global optimality).

Let  $\mathbf{w}_0$  be given, and  $\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla f(\mathbf{w}_i)$ .

If  $f$  is convex and  $\beta$ -smooth, and  $\eta = 1/\beta$ , then  $\forall \mathbf{u} \in \mathbb{R}^d$ ,

$$f(\mathbf{w}_t) - f(\mathbf{u}) \leq \frac{1}{t} \sum_{i=1}^t (f(\mathbf{w}_i) - f(\mathbf{u})) \leq \frac{\beta}{2t} (\|\mathbf{w}_0 - \mathbf{u}\|^2 - \|\mathbf{w}_t - \mathbf{u}\|^2).$$

**Proof.** For each  $i \leq t$ , using the previous proof,

$$\begin{aligned} \|\mathbf{w}_i - \mathbf{u}\|^2 &= \|\mathbf{w}_{i-1} - \mathbf{u}\|^2 - 2\eta \nabla f(\mathbf{w}_{i-1})^\top (\mathbf{w}_{i-1} - \mathbf{u}) + \eta^2 \|\nabla f(\mathbf{w}_{i-1})\|^2 \\ &\leq \|\mathbf{w}_{i-1} - \mathbf{u}\|^2 + 2\eta (f(\mathbf{u}) - f(\mathbf{w}_{i-1})) + 2\eta^2 \beta (f(\mathbf{w}_{i-1}) - f(\mathbf{w}_i)) \\ &= \|\mathbf{w}_{i-1} - \mathbf{u}\|^2 + \frac{2}{\beta} (f(\mathbf{u}) - f(\mathbf{w}_i)). \end{aligned}$$

Rearranging and then averaging these inequalities over  $i \leq t$  gives the bound.  $\square$

## Example: OLS

As a consequence, for least squares, starting from  $\mathbf{w}_0 = 0$ ,

$$\widehat{\mathcal{R}}(\mathbf{w}_t) - \widehat{\mathcal{R}}(\hat{\mathbf{w}}_{\text{ols}}) \leq \frac{\beta \|\hat{\mathbf{w}}_{\text{ols}}\|^2}{2t},$$

where  $\beta$  and  $\|\hat{\mathbf{w}}_{\text{ols}}\|$  depend on  $\mathbf{X}$  and  $\mathbf{y}$ .

# Jensen's inequality

**Jensen's inequality:** if  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is convex, then  $\mathbb{E}f(\mathbf{X}) \geq f(\mathbb{E}\mathbf{X})$ .

**Proof.** Set  $\mathbf{y} := \mathbb{E}\mathbf{X}$ , and pick any  $\mathbf{s} \in \partial f(\mathbb{E}\mathbf{X})$ . Then

$$\mathbb{E}f(\mathbf{X}) \geq \mathbb{E}\left(f(\mathbf{y}) + \mathbf{s}^\top(\mathbf{X} - \mathbf{y})\right) = f(\mathbf{y}) + \mathbf{s}^\top\mathbb{E}(\mathbf{X} - \mathbf{y}) = f(\mathbf{y}).$$

- ▶ Shalev-Shwartz/Ben-David: chapter 12.
- ▶ Further advanced reading on convexity and convex optimization:
  - ▶ Convex optimization: course lecture slides by Lieven Vandenberghe part 1 part 2 part 3, “introductory lectures on convex optimization” by Nesterov.
  - ▶ Convexity: “fundamentals of convex analysis” by Hiriart-Urruty and Lemaréchal, “convex analysis and nonlinear optimization” by Borwein and Lewis.