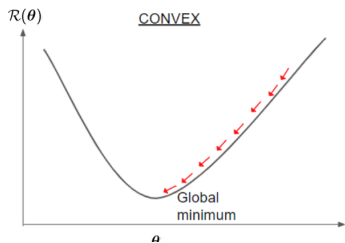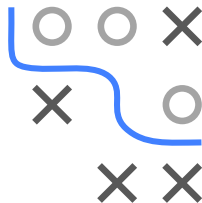# Optimization in Machine Learning

## Deep dive
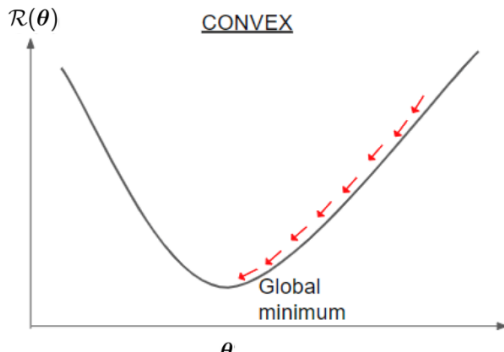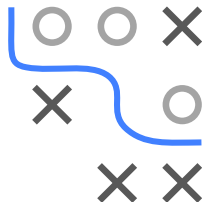## Gradient descent and optimality



**Learning goals**

- Convergence of GD
- Proof strategy and tools
- Descent lemma

# SETTING

- GD is **greedy**: **locally optimal** moves in each iteration

- If $f$ is **convex**, **differentiable** and has a **Lipschitz gradient**, GD converges to global minimum for sufficiently small step sizes.
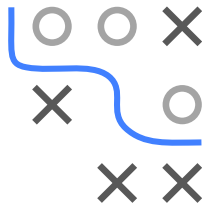
# SETTING

- **Assumptions:**
    - *f* convex and differentiable
    - Global minimum $\boldsymbol{x}^*$ exists
    - *f* has Lipschitz gradient ($\nabla f$ does not change too fast)

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\tilde{\boldsymbol{x}})\| \leq L\|\boldsymbol{x} - \tilde{\boldsymbol{x}}\| \quad \text{for all } \boldsymbol{x}, \tilde{\boldsymbol{x}}$$

**Theorem** (Convergence of GD)**.** GD with step size $\alpha \leq 1/L$ yields

$$f(\boldsymbol{x}^{[k]}) - f(\boldsymbol{x}^*) \leq \frac{\|\boldsymbol{x}^{[0]} - \boldsymbol{x}^*\|^2}{2\alpha k}.$$

In other words: GD converges with rate $\mathcal{O}(1/k)$.

## PROOF STRATEGY

**1** Show that $f(\mathbf{x}^{[t]})$ **strictly decreases** with each iteration $t$
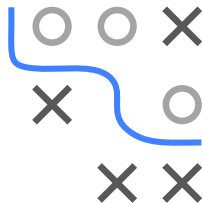
**Descent lemma:**
$$f(\mathbf{x}^{[t+1]}) \leq f(\mathbf{x}^{[t]}) - \frac{\alpha}{2}\|\nabla f(\mathbf{x}^{[t]})\|^2$$

**2** Bound **error of one step**

$$f(\mathbf{x}^{[t+1]}) - f(\mathbf{x}^*) \leq \frac{1}{2\alpha}\left(\|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t+1]} - \mathbf{x}^*\|^2\right)$$

**3** Finalize by **telescoping** argument
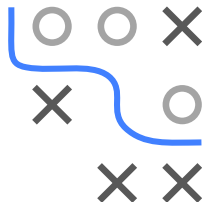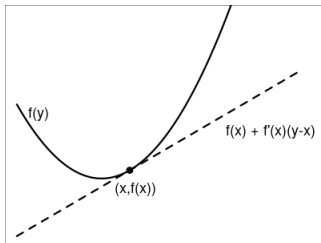
# MAIN TOOL

- **Recall:** First order condition of convexity

> Every tangent line of $f$ is always below $f$.
>
> $$f(\boldsymbol{y}) \geq f(\boldsymbol{x}) + \nabla f(\boldsymbol{x})^\top (\boldsymbol{y} - \boldsymbol{x})$$
>
>

# DESCENT LEMMA

- **Recall:** $\nabla f$ Lipschitz $\implies \nabla^2 f(\boldsymbol{x}) \preccurlyeq L \cdot \boldsymbol{I}$ for all $\boldsymbol{x}$
- This gives convexity of $g(\boldsymbol{x}) := \frac{L}{2}\|\boldsymbol{x}\|^2 - f(\boldsymbol{x})$ since

$$\nabla^2 g(\boldsymbol{x}) = L \cdot \boldsymbol{I} - \nabla^2 f(\boldsymbol{x}) \succcurlyeq 0.$$

- First order condition of convexity of $g$ yields

$$g(\boldsymbol{x}) \geq g(\boldsymbol{x}^{[t]}) + \nabla g(\boldsymbol{x}^{[t]})^\top (\boldsymbol{x} - \boldsymbol{x}^{[t]})$$

$$\Leftrightarrow \quad \frac{L}{2}\|\boldsymbol{x}\|^2 - f(\boldsymbol{x}) \geq \frac{L}{2}\|\boldsymbol{x}^{[t]}\|^2 - f(\boldsymbol{x}^{[t]}) + (L\boldsymbol{x}^{[t]} - \nabla f(\boldsymbol{x}^{[t]}))^\top (\boldsymbol{x} - \boldsymbol{x}^{[t]})$$

$$\Leftrightarrow \qquad \vdots$$

$$\Leftrightarrow \quad f(\boldsymbol{x}) \leq f(\boldsymbol{x}^{[t]}) + \nabla f(\boldsymbol{x}^{[t]})^\top (\boldsymbol{x} - \boldsymbol{x}^{[t]}) + \frac{L}{2}\|\boldsymbol{x} - \boldsymbol{x}^{[t]}\|^2$$

- **Now:** One GD step with step size $\alpha \leq 1/L$:

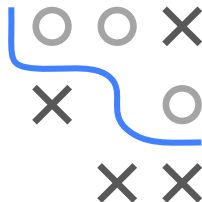$$\boldsymbol{x} \leftarrow \boldsymbol{x}^{[t+1]} = \boldsymbol{x}^{[t]} - \alpha \nabla f\left(\boldsymbol{x}^{[t]}\right)$$

## DESCENT LEMMA

$$f(\mathbf{x}^{[t+1]}) \leq f(\mathbf{x}^{[t]}) + \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t+1]} - \mathbf{x}^{[t]}) + \frac{L}{2}\|\mathbf{x}^{[t+1]} - \mathbf{x}^{[t]}\|^2$$

$$= f(\mathbf{x}^{[t]}) + \nabla f(\mathbf{x}^{[t]})^\top (\mathbf{x}^{[t]} - \alpha \nabla f(\mathbf{x}^{[t]}) - \mathbf{x}^{[t]})$$
$$+ \frac{L}{2}\|\mathbf{x}^{[t]} - \alpha \nabla f(\mathbf{x}^{[t]}) - \mathbf{x}^{[t]}\|^2$$

$$= f(\mathbf{x}^{[t]}) - \nabla f(\mathbf{x}^{[t]})^\top \alpha \nabla f(\mathbf{x}^{[t]}) + \frac{L}{2}\|\alpha \nabla f(\mathbf{x}^{[t]})\|^2$$

$$= f(\mathbf{x}^{[t]}) - \alpha \|\nabla f(\mathbf{x}^{[t]})\|^2 + \frac{L\alpha^2}{2}\|\nabla f(\mathbf{x}^{[t]})\|^2$$

$$\leq f(\mathbf{x}^{[t]}) - \frac{\alpha}{2}\|\nabla f(\mathbf{x}^{[t]})\|^2$$

- **Note:** $\alpha \leq 1/L$ yields $L\alpha^2 \leq \alpha$
  - $\|\nabla f(\mathbf{x}^{[t]})\|^2 > 0$ unless $\nabla f(\mathbf{x}) = \mathbf{0}$
  - $f$ **strictly decreases** with each GD iteration until optimum reached
  - Descent lemma yields bound on **guaranteed progress** if $\alpha \leq 1/L$
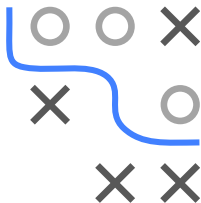    (explains why GD may diverge if step sizes too large)

# ONE STEP ERROR BOUND

- Again, first order condition of convexity gives

$$f(\boldsymbol{x}^{[t]}) - f(\boldsymbol{x}^*) \leq \nabla f(\boldsymbol{x}^{[t]})^\top (\boldsymbol{x}^{[t]} - \boldsymbol{x}^*).$$

- This and the descent lemma yields

$$
\begin{aligned}
f(\boldsymbol{x}^{[t+1]}) - f(\boldsymbol{x}^*) &\leq f(\boldsymbol{x}^{[t]}) - \frac{\alpha}{2}\|\nabla f(\boldsymbol{x}^{[t]})\|^2 - f(\boldsymbol{x}^*) \\
&= f(\boldsymbol{x}^{[t]}) - f(\boldsymbol{x}^*) - \frac{\alpha}{2}\|\nabla f(\boldsymbol{x}^{[t]})\|^2 \\
&\leq \nabla f(\boldsymbol{x}^{[t]})^\top (\boldsymbol{x}^{[t]} - \boldsymbol{x}^*) - \frac{\alpha}{2}\|\nabla f(\boldsymbol{x}^{[t]})\|^2 \\
&= \frac{1}{2\alpha}\left(\|\boldsymbol{x}^{[t]} - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}^{[t]} - \boldsymbol{x}^* - \alpha\nabla f(\boldsymbol{x}^{[t]})\|^2\right) \\
&= \frac{1}{2\alpha}\left(\|\boldsymbol{x}^{[t]} - \boldsymbol{x}^*\|^2 - \|\boldsymbol{x}^{[t+1]} - \boldsymbol{x}^*\|^2\right)
\end{aligned}
$$

- **Note:** Line $3 \to 4$ is hard to see (just expand line 4).

## FINALIZATION

- Summing over iterations yields

$$k(f(\mathbf{x}^{[k]}) - f(\mathbf{x}^*)) \leq \sum_{t=1}^{k} [f(\mathbf{x}^{[t]}) - f(\mathbf{x}^*)]$$

$$\leq \sum_{t=1}^{k} \frac{1}{2\alpha} \left[ \|\mathbf{x}^{[t-1]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[t]} - \mathbf{x}^*\|^2 \right]$$

$$= \frac{1}{2\alpha} \left( \|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2 - \|\mathbf{x}^{[k]} - \mathbf{x}^*\|^2 \right)$$

$$\leq \frac{1}{2\alpha} \left( \|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2 \right).$$

- **Arguments:** Descent lemma (line 1).
  Telescoping sum (line $2 \rightarrow 3$).

$$f(\mathbf{x}^{[t+1]}) - f(\mathbf{x}^*) \leq \frac{\|\mathbf{x}^{[0]} - \mathbf{x}^*\|^2}{2\alpha k}$$