Gradient Descent & its variations $\in$ Convex Optimization

Set function

$$\min_{w} f(w)$$

$$\text{s.t.} \quad w \in K$$

Convex optimization if $\begin{cases} f: \text{convex} \\ K: \text{convex} \end{cases}$
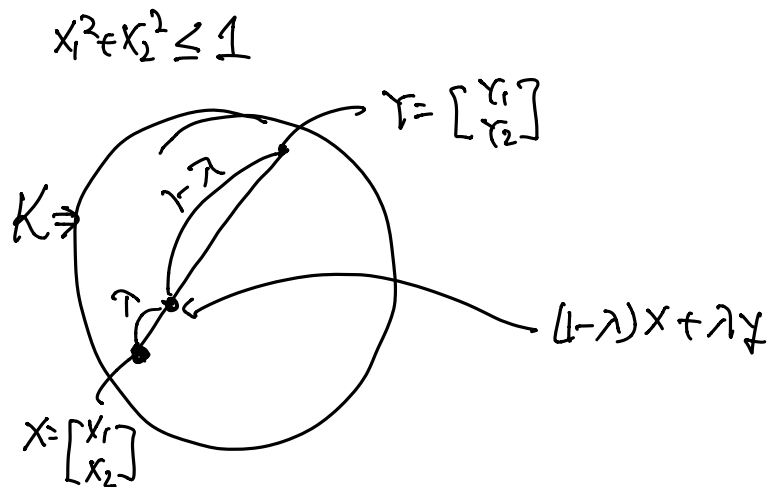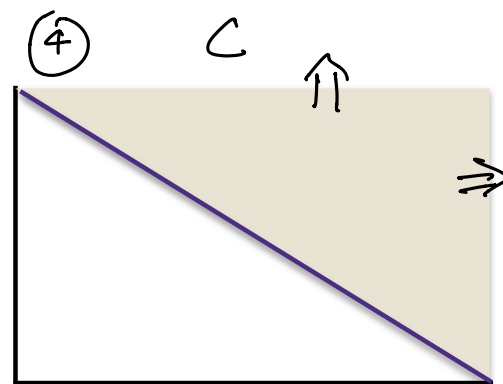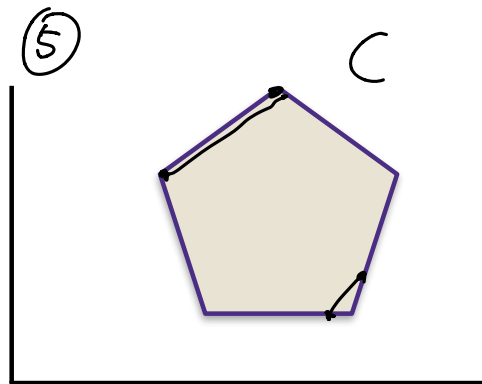
First-order methods

# Convexity

# What is a convex set?

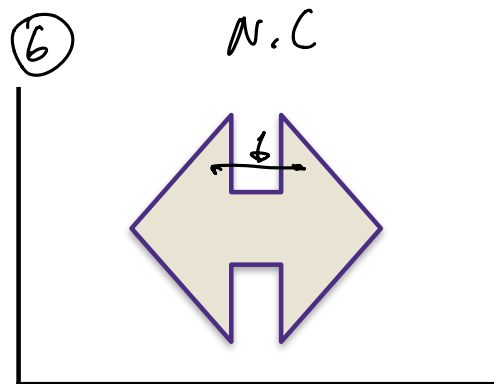A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$
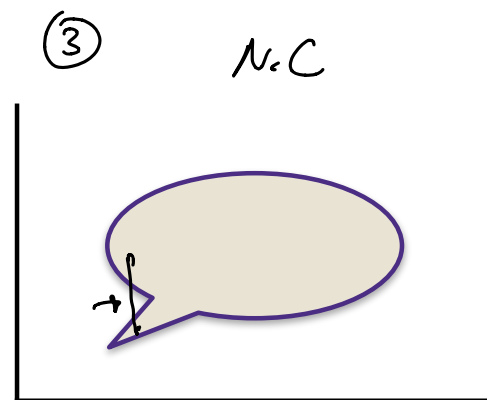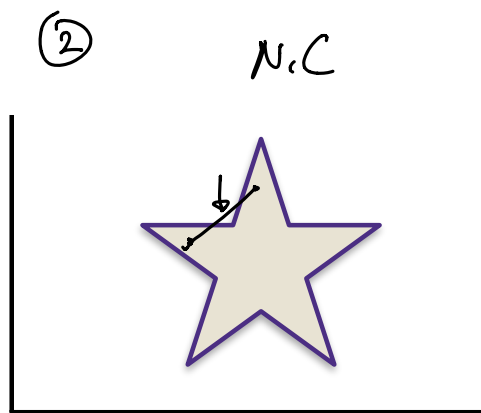
# What is a convex set?

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

# What is a convex function?

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$

# What is a convex function?

> A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$

$f(x)$

① C

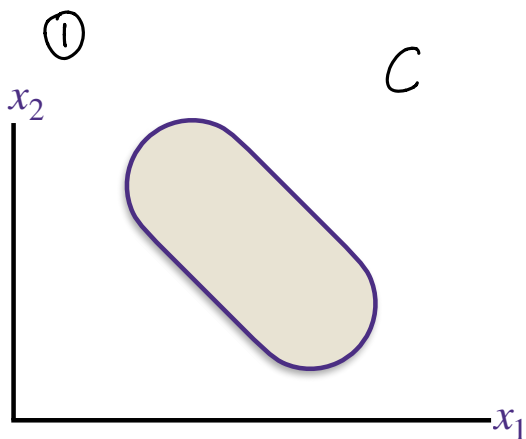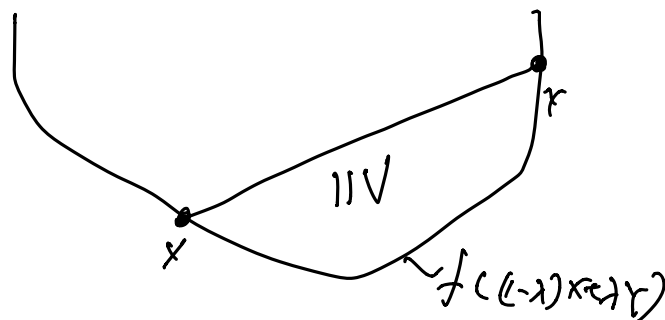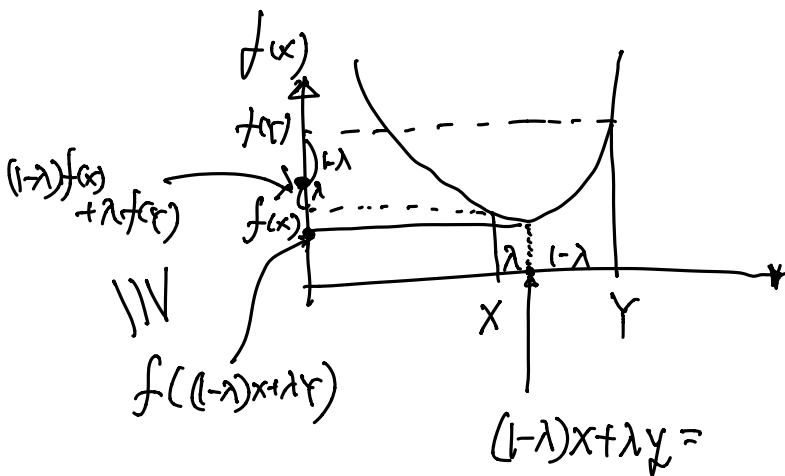② N.C

③ C

$x$

⑥ N.C

⑤ C

④ non-convex concave

X     Y

# Convex functions and convex sets?

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$
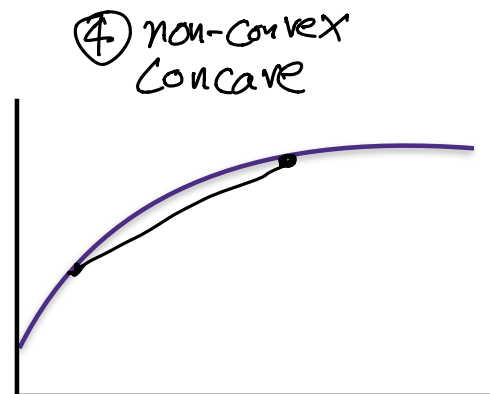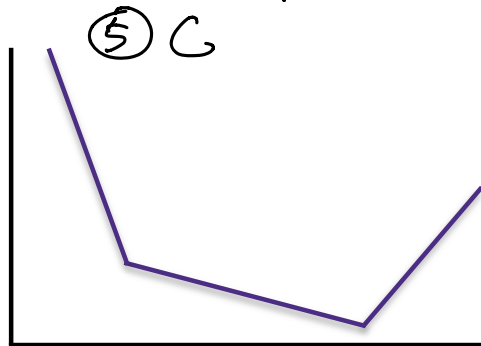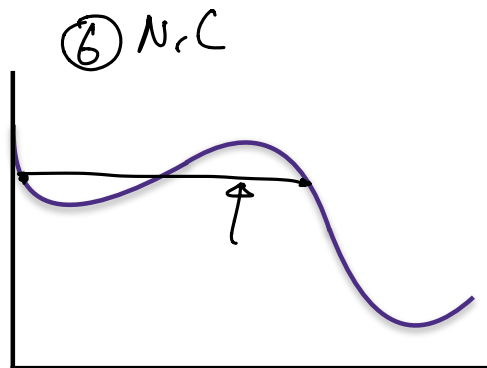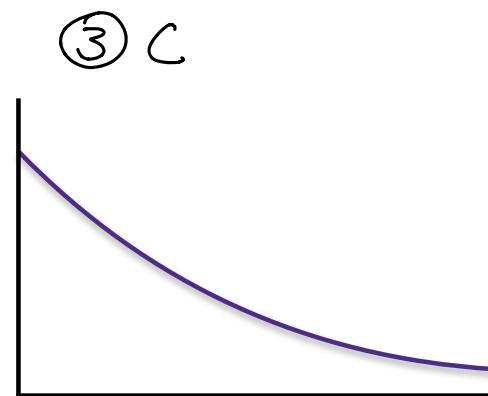
A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex **Set**

Definition: Epigraph $f = \{ (x,t) \mid f(x) \leq t \}$

$\mathbb{R}^d \quad \mathbb{R}$

$f(x)$

graph $\{ (x,t) \mid f(x) = t \}$

$\mathbb{R}^d \quad \mathbb{R}$

$d = 1$

$x$

# Convex functions and convex sets?
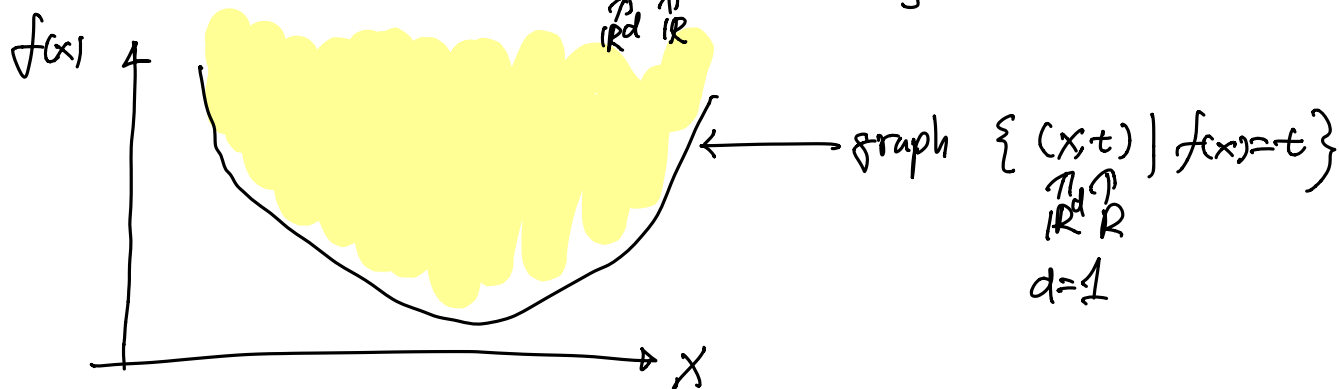
A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

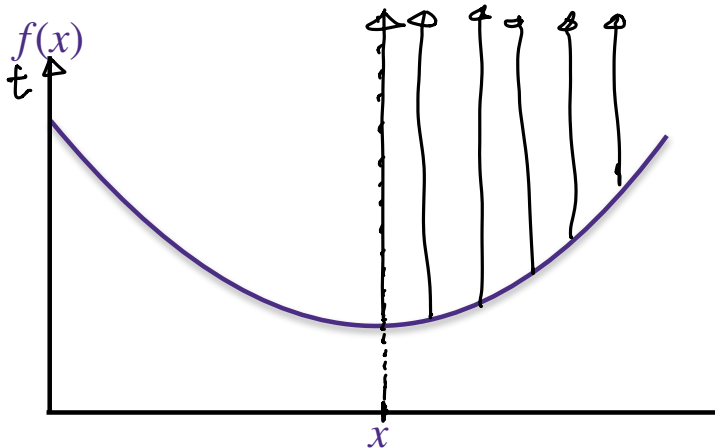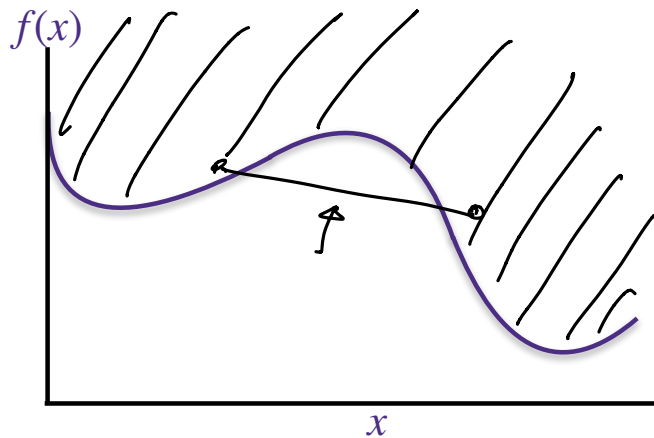$C$ $\{t \mid t \geq f(x)\}$ $N_{\epsilon} C$

$f(x)$
$t$

$f(x)$

$x$

$x$

# More definitions of convexity

A set $K \subset \mathbb{R}^d$ is convex if $(1 - \lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0, 1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x, t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex

*

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is <u>differentiable everywhere</u> is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in dom(f)$

LHS          RHS



$f(y)$=LHS
IIV

⅃RHS

N.C.

$f(y)$

$y$          $x$          $y$

# More definitions of convexity

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in dom(f)$

Positive semidefiniteness.



$f(x)$

$x$

$\dfrac{df(x)}{dx}$

$x$

$\dfrac{d^2 f(x)}{dx^2}$  $> 0$

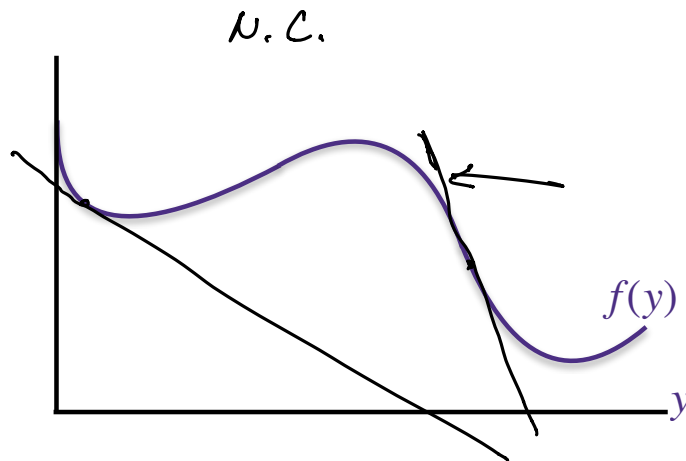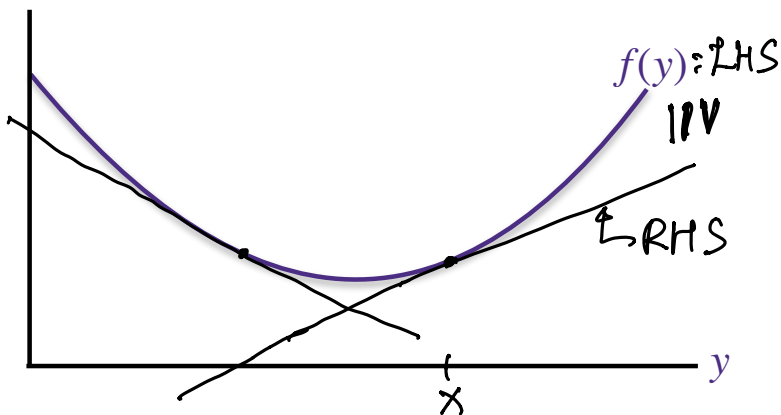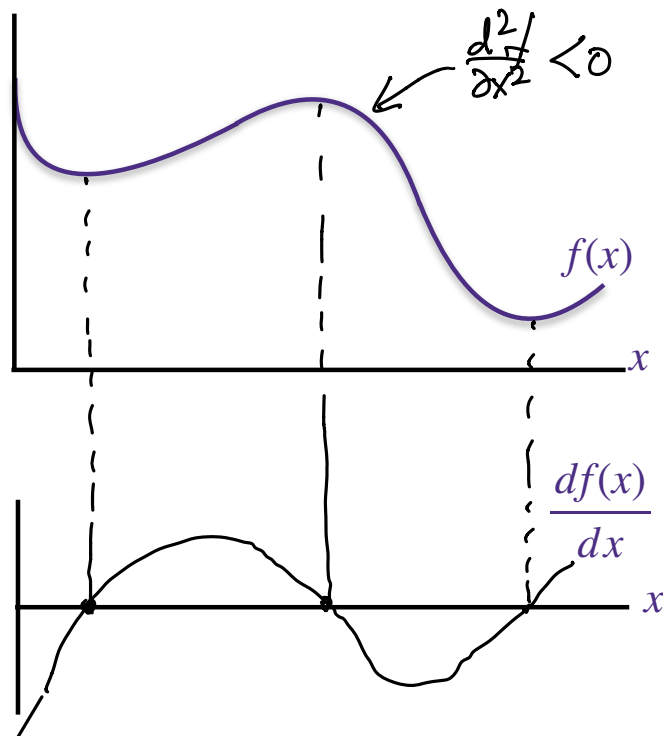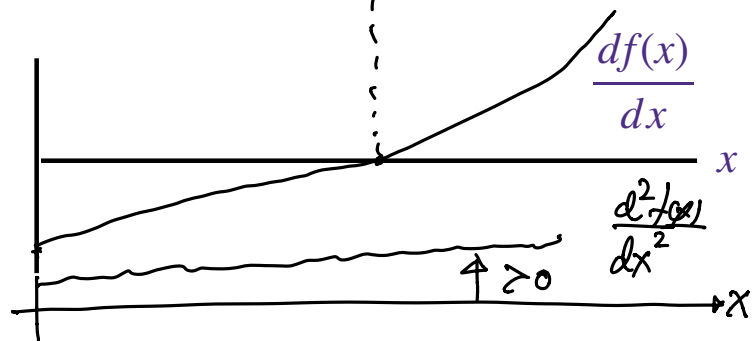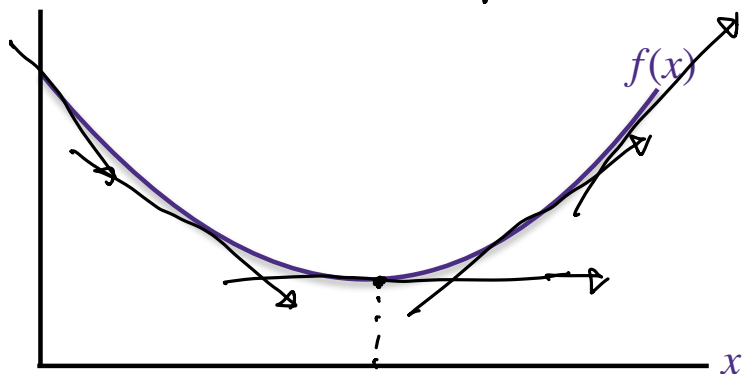$\dfrac{d^2 f}{dx^2} < 0$

$f(x)$

$x$

$\dfrac{df(x)}{dx}$

$x$

# More definitions of convexity

A set $K \subset \mathbb{R}^d$ is convex if $(1-\lambda)x + \lambda y \in K$ for all $x, y \in K$ and $\lambda \in [0,1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if $f((1-\lambda)x + \lambda y) \leq (1-\lambda)f(x) + \lambda f(y)$ for all $x, y \in \mathbb{R}^d$ and $\lambda \in [0,1]$

A function $f : \mathbb{R}^d \to \mathbb{R}$ is convex if the set $\{(x,t) \in \mathbb{R}^{d+1} : f(x) \leq t\}$ is convex
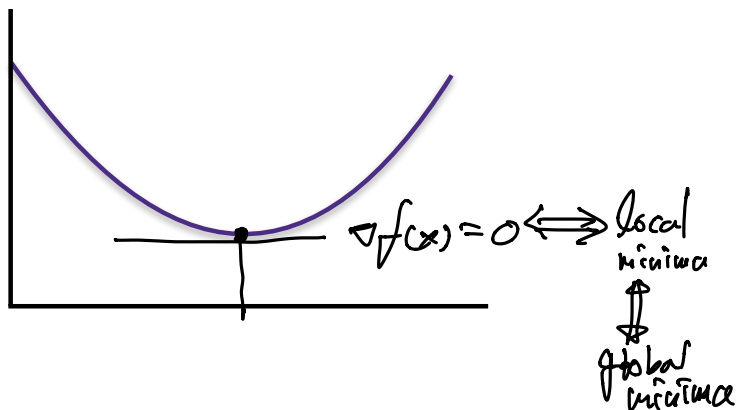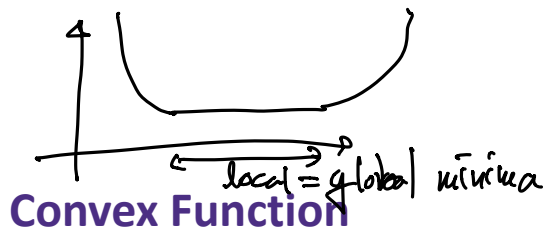
A function $f : \mathbb{R}^d \to \mathbb{R}$ that is differentiable everywhere is convex if $f(y) \geq f(x) + \nabla f(x)^\top (y - x)$ for all $x, y \in dom(f)$

A function $f : \mathbb{R}^d \to \mathbb{R}$ that is twice-differentiable everywhere is convex if $\nabla^2 f(x) \succeq 0$ for all $x \in dom(f)$
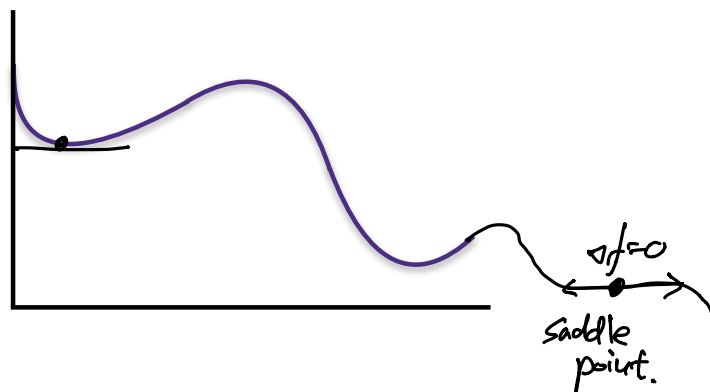
# Why do we care about convexity?

Convex functions

- All local minima are global minima

- Efficient to optimize (e.g., gradient descent)



**Convex Function**

$local = global\ minima$

$\nabla f(x) = 0 \Longleftrightarrow local\ minima$

$\Updownarrow$

$global\ minima$

**Non-convex Function**

$\nabla f = 0$

$saddle\ point.$

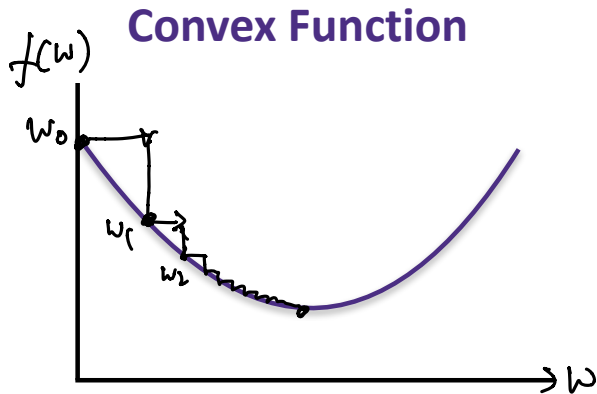# Gradient Descent on $\min_w f(w)$

Initialize: $w_0 = 0$

for $t = 1, 2, \ldots$

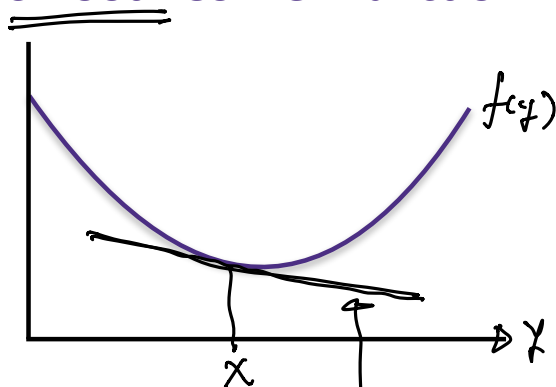$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

Step size / learning rate

**Convex Function**

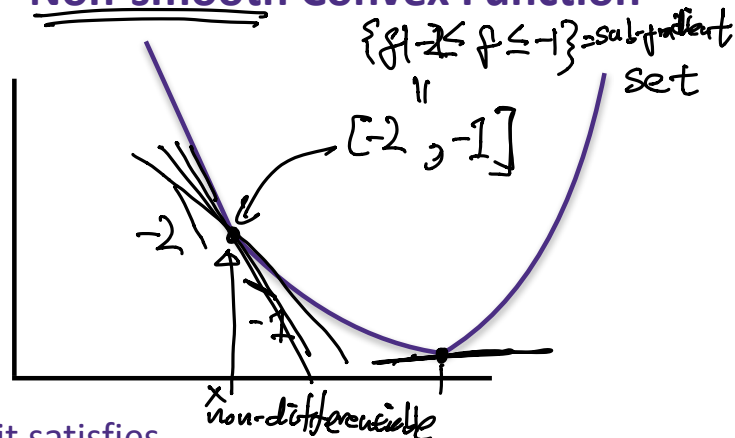**Non-convex Function**

$\nabla f \approx 0$

# Sub-Gradient

Definition: a function is **non-smooth** if it is not differentiable everywhere

**Smooth Convex Function**

**Non-smooth Convex Function**

$\{g: -2 \leq g \leq -1\}$ =subgradient set

$\Uparrow$

$[-2, -1]$

$-2$

$-1$

$f(y)$

$x$

$y$

$x$ non-differentiable

Definition: a vector $g \in \mathbb{R}^d$ is a **sub-gradient** at $x$ if it satisfies

$$\left\{ g: f(y) \geq \boxed{f(x) + g^T(y-x)} \text{ for all } y \in \mathbb{R}^d \right\} = \text{subgradient set}$$

$\nabla f(x)$

for smooth convex functions, the minimum is achieved at points where gradient is zero

for non-smooth convex functions, the minimum is achieved at points where sub-gradient set includes the zero vector

# Sub-Gradient Descent

Initialize: $w_0 = 0$
for $t = 1, 2, \ldots$

    Find any $g_t$ such that $f(y) \geq f(w_t) + g_t^\top (y - w_t)$

$$w_{t+1} = w_t - \eta g_t$$

**Convex Function**



**Non-convex Function**

# Coordinate descent

Initialize: $w_0 = 0$

for $t = 1, 2, \ldots$

    Let $i_t = t \ \% \ d$

$$w_{t+1}^{(i_t)} = w_t^{(i_t)} - \eta_t \frac{\partial f(w)}{\partial w^{(i_t)}} \bigg|_{w=w_t}$$

# Machine Learning Problems

- **Given data:**

$$\{(x_i, y_i)\}_{i=1}^n \qquad x_i \in \mathbb{R}^d \qquad y_i \in \mathbb{R}$$

- **Learning a model's parameters:** $\displaystyle\sum_{i=1}^n \ell_i(w)$

$$\text{Logistic Loss: } \ell_i(w) = \log(1 + \exp(-y_i \, x_i^T w))$$

$$\text{Squared error Loss: } \ell_i(w) = (y_i - x_i^T w)^2$$

Gradient Descent:

$$w_{t+1} = w_t - \eta \nabla_w \left( \frac{1}{n} \sum_{i=1}^n \ell_i(w) \right) \Big|_{w=w_t}$$

# Optimization summary

- You can always run gradient descent whether f is convex or not. But you only have guarantees if f is convex

- Many bells and whistles can be added onto gradient descent such as momentum and dimension-specific step-sizes (Nesterov, Adagrad, ADAM, etc.)