# Convex Optimization
## A Gentle Introduction

S.V. N. (vishy) Vishwanathan
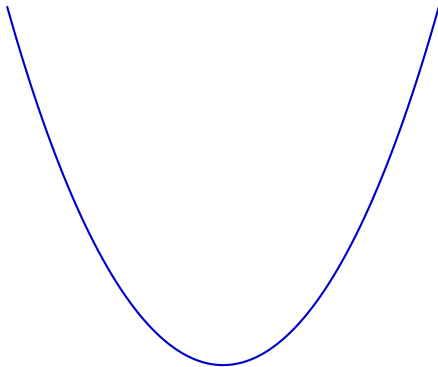
Purdue University
vishy@purdue.edu

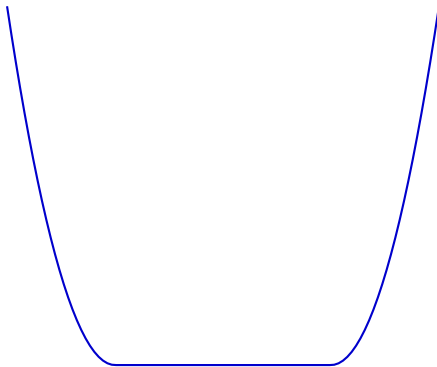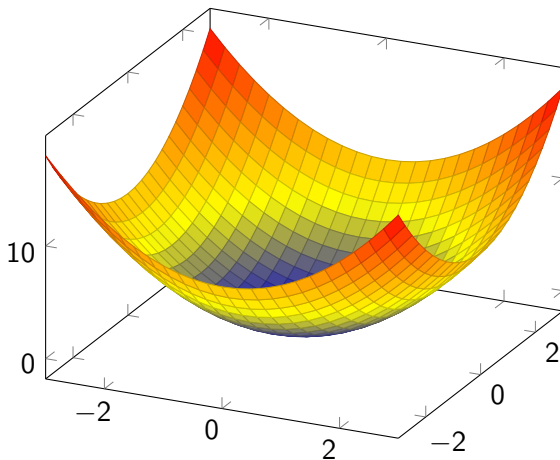April 15, 2013

## Outline

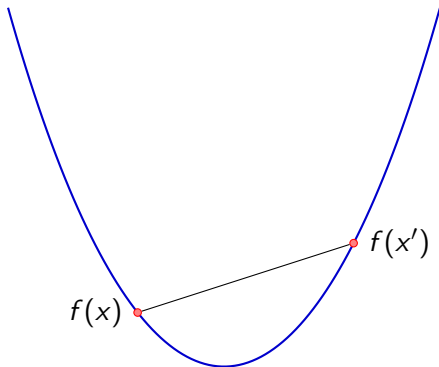## Convex Functions

## Convex Functions

# Convex Functions

**Disclaimer**

- My focus is on intuition
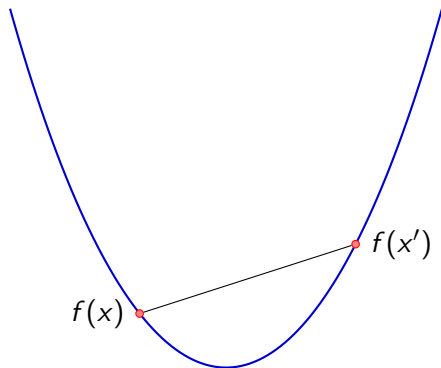- Not mathematically rigorous

**Convex Function**



A function $f$ is convex if, and only if, for all $x, x'$ and $\lambda \in (0, 1)$

$$f(\lambda x + (1 - \lambda)x') \leq \lambda f(x) + (1 - \lambda)f(x')$$

**Convex Function**



A function $f$ is <span style="color:magenta">strictly</span> convex if, and only if, for all $x, x'$ and $\lambda \in (0, 1)$
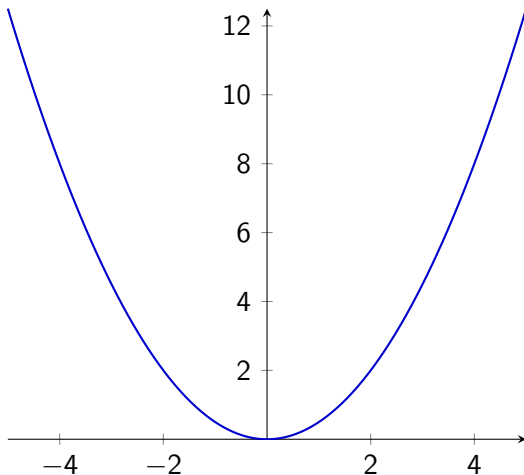
$$f(\lambda x + (1 - \lambda)x') < \lambda f(x) + (1 - \lambda)f(x')$$

**Exercise: Jensen's Inequality**

- Extend the definition of convexity to show that if $f$ is convex, then for all $\lambda_i \geq 0$ such that $\sum_i \lambda_i = 1$ we have
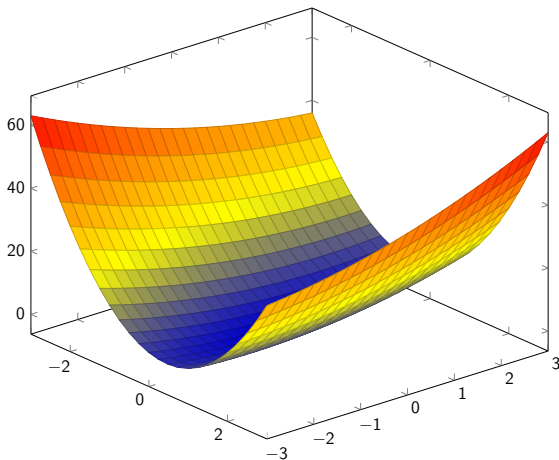
$$f\left(\sum_i \lambda_i x_i\right) \leq \sum_i \lambda_i f(x_i)$$
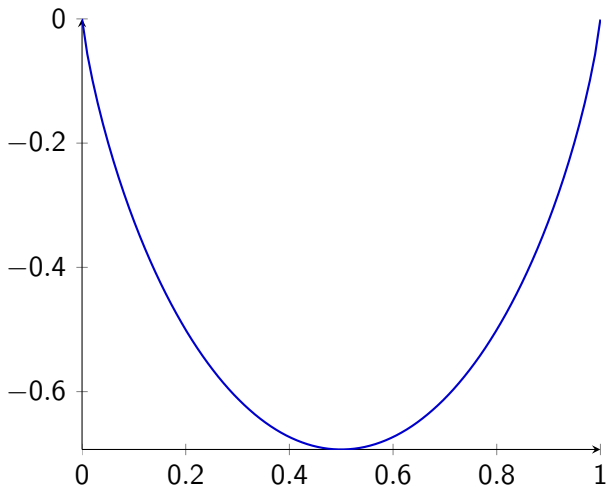
## Some Familiar Examples



$$f(x) = \frac{1}{2}x^2 \text{ (Square norm)}$$
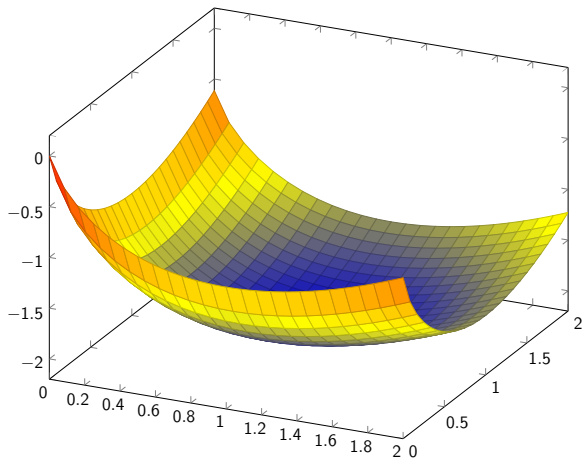
# Some Familiar Examples



$$f(x,y) = \frac{1}{2} \begin{bmatrix} x, y \end{bmatrix} \begin{bmatrix} 10, 1 \\ 2, 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

## Some Familiar Examples
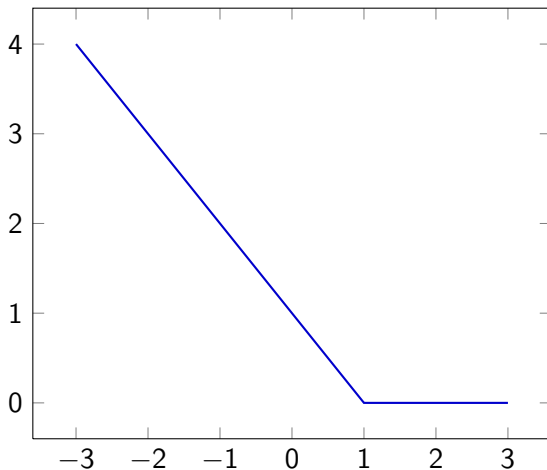


$f(x) = x \log x + (1 - x) \log(1 - x)$ (Negative entropy)

## Some Familiar Examples



$f(x, y) = x \log x + y \log y - x - y$ (Un-normalized negative entropy)
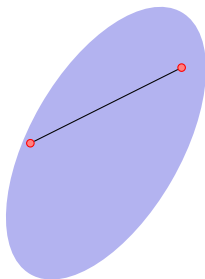
## Some Familiar Examples



$$f(x) = \max(0, 1 - x) \text{ (Hinge Loss)}$$

**Some Other Important Examples**

- Linear functions: $f(x) = ax + b$
- Softmax: $f(x) = \log \sum_i \exp(x_i)$
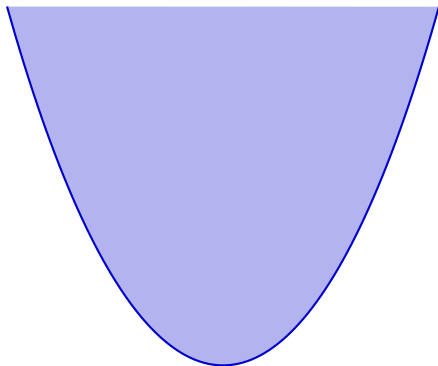- Norms: For example the 2-norm $f(x) = \sqrt{\sum_i x_i^2}$

## Convex Sets



A set $C$ is convex if, and only if, for all $x, x' \in C$ and $\lambda \in (0, 1)$ we have

$$\lambda x + (1 - \lambda)x' \in C$$
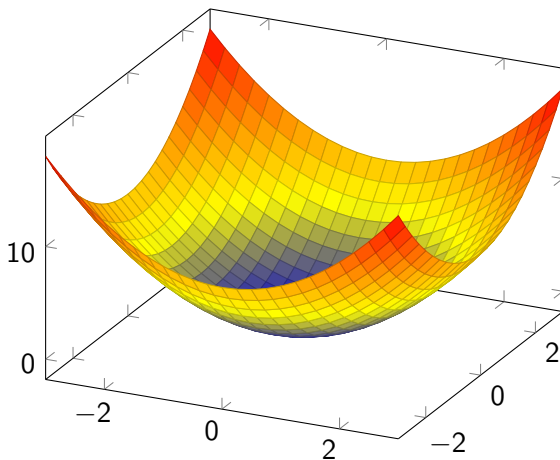
**Convex Sets and Convex Functions**



A function $f$ is convex if, and only if, its epigraph is a convex set

**Convex Sets and Convex Functions**

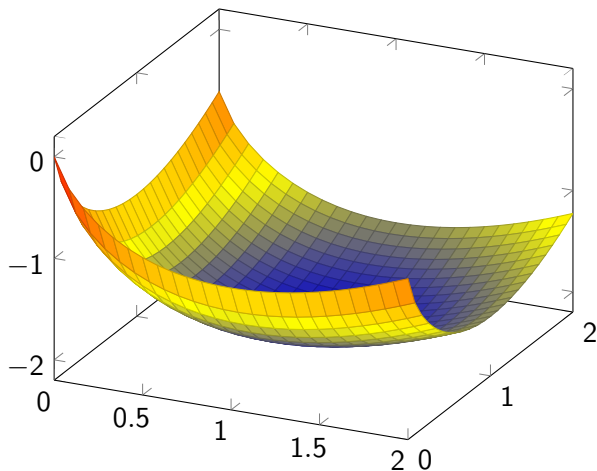- Indicator functions of convex sets are convex

$$I_C(x) = \begin{cases} 0 & \text{if } x \in C \\ \infty & \text{otherwise.} \end{cases}$$

## Below sets of Convex Functions



$$f(x, y) = x^2 + y^2$$
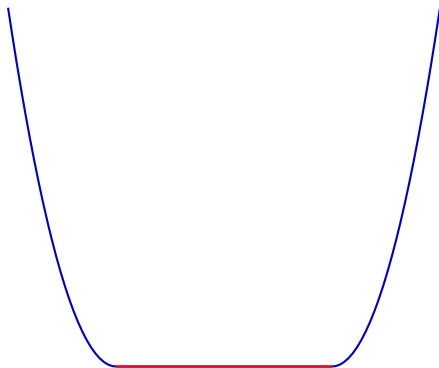
## Below sets of Convex Functions



$$f(x, y) = x \log x + y \log y - x - y$$

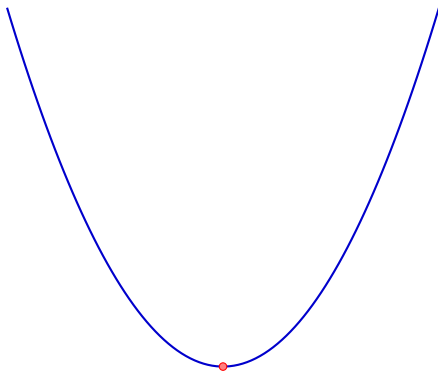**Below sets of Convex Functions**

- If $f$ is convex, then all its level sets are convex
- Is the converse true? (Exercise: construct a counter-example)

**Minima on Convex Sets**



- Set of minima of a convex function is a convex set
- Proof: Consider the set $\{x : f(x) \leq f^*\}$

**Minima on Convex Sets**



- Set of minima of a strictly convex function is a singleton
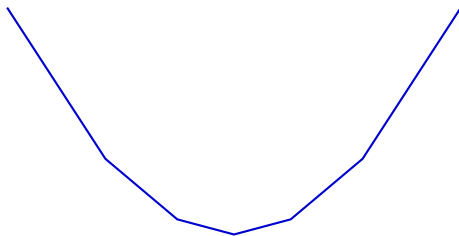- Proof: try this at home!

## Outline

**Set Operations**

- Intersection of convex sets is convex
- Image of a convex set under a linear transformation is convex
- Inverse image of a convex set under a linear transformation is convex

**Function Operations**

- Linear Combination with non-negative weights: $f(x) = \sum_i w_i f_i(x)$ s.t. $w_i \geq 0$
- Pointwise maximum: $f(x) = \max_i f_i(x)$
- Composition with affine function: $f(x) = g(Ax + b)$
- Projection along a direction: $f(\eta) = g(x_0 + \eta d)$
- Restricting the domain on a convex set: $f(x)$s.t. $x \in \mathcal{C}$
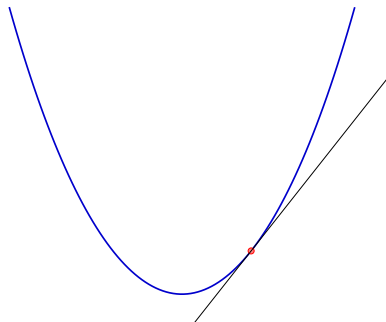
**One Quick Example**



The piecewise linear function $f(x) := \max_i \langle u_i, x \rangle$ is convex

## Outline

1. **Convex Functions and Sets**

2. **Operations Which Preserve Convexity**

3. **First Order Properties**

4. **Minimizing a 1-d Convex Function**

5. **Coordinate Descent**

6. **Gradient Descent**

## First Order Taylor Expansion

The First Order Taylor approximation globally lower bounds the function



For any $x$ and $x'$ we have

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

**Identifying the Minimum**

- Let $f : X \to \mathbb{R}$ be a differentiable convex function. Then $x$ is a minimizer of $f$, if, and only if,
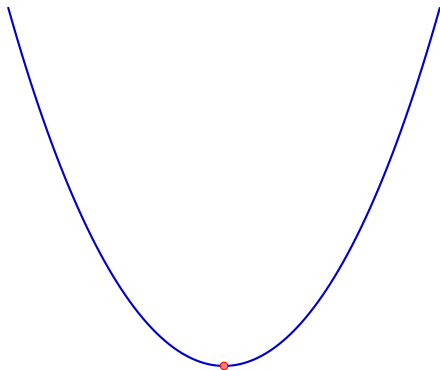
$$\langle x' - x, \nabla f(x) \rangle \geq 0 \text{ for all } x'.$$

- One way to ensure this is to set $\nabla f(x) = 0$

- Minimizing a smooth convex function is the same as finding an $x$ such that $\nabla f(x) = 0$

## Outline

1. Convex Functions and Sets

2. Operations Which Preserve Convexity

3. First Order Properties

4. **Minimizing a 1-d Convex Function**

5. Coordinate Descent

6. Gradient Descent

**Problem Statement**



- Given a black-box which can compute $J : \mathbb{R} \to \mathbb{R}$ and $J' : \mathbb{R} \to \mathbb{R}$ find the minimum value of $J$

**Increasing Gradients**

- From the first order conditions

$$J(w) \geq J(w') + (w - w') \cdot J'(w')$$

and
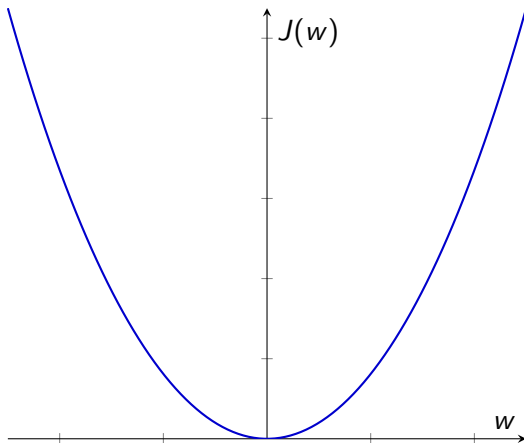
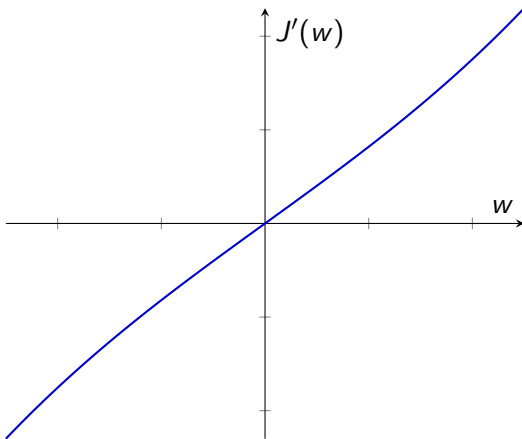$$J(w') \geq J(w) + (w' - w) \cdot J'(w)$$

- Add the two

$$(w - w') \cdot (J'(w) - J'(w')) \geq 0$$
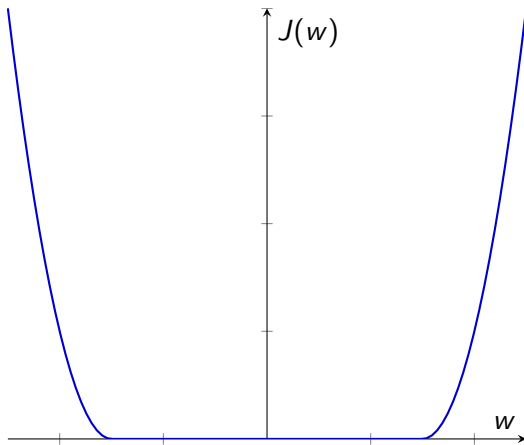
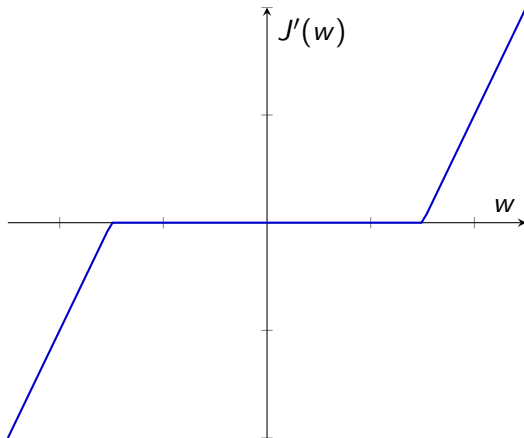$w \geq w'$ implies that $J'(w) \geq J'(w')$

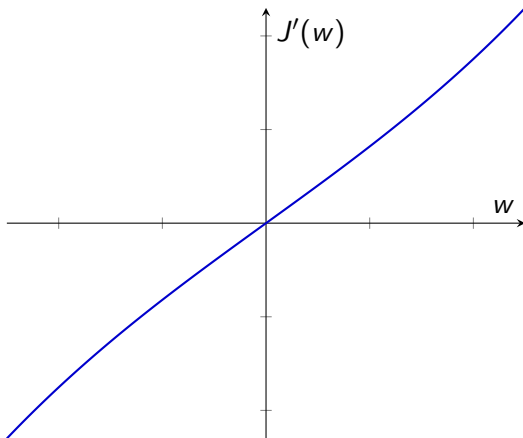## Increasing Gradients

## Increasing Gradients

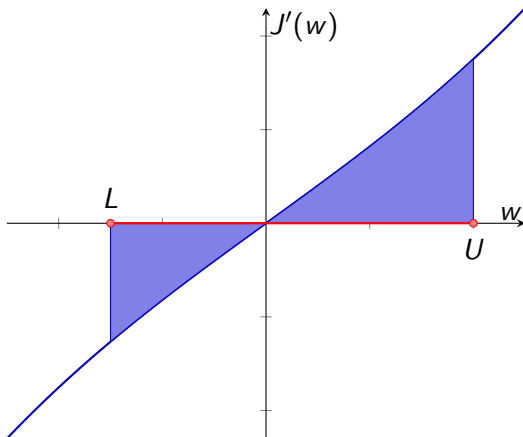## Increasing Gradients

## Increasing Gradients
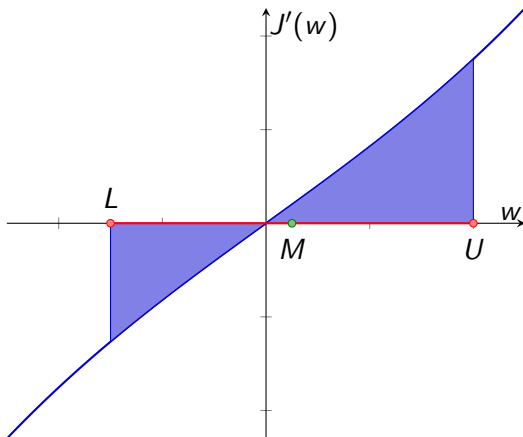
## Problem Restatement



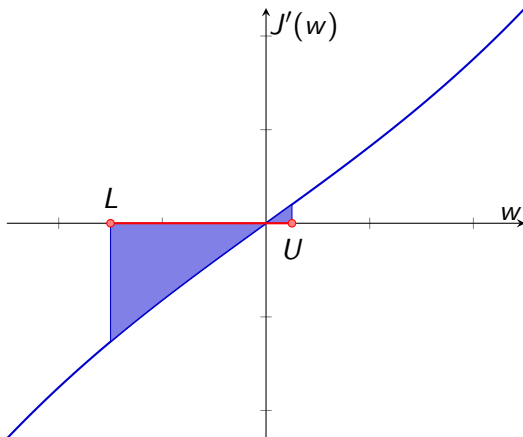- Identify the point where the increasing function $J'$ crosses zero

## Bisection Algorithm
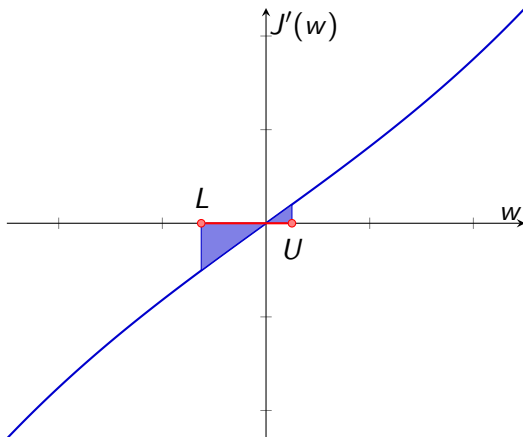
## Bisection Algorithm

# Bisection Algorithm

# Bisection Algorithm

# Bisection Algorithm

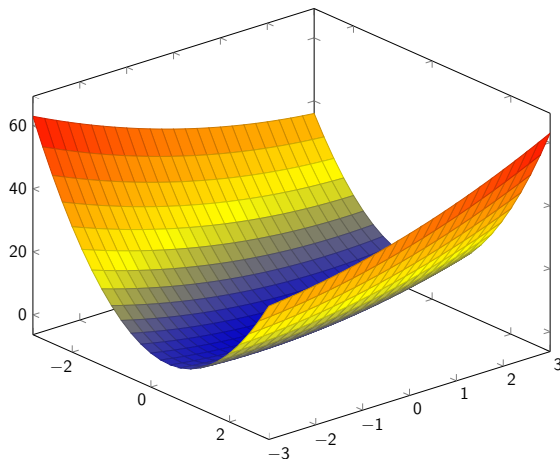**Interval Bisection**

**Require:** $L, U, \epsilon$
1: $maxgrad \leftarrow J'(U)$
2: **while** $(U - L) \cdot maxgrad > \epsilon$ **do**
3:     $M \leftarrow \frac{U+L}{2}$
4:     **if** $J'(M) > 0$ **then**
5:        $U \leftarrow M$
6:     **else**
7:        $L \leftarrow M$
8:     **end if**
9: **end while**
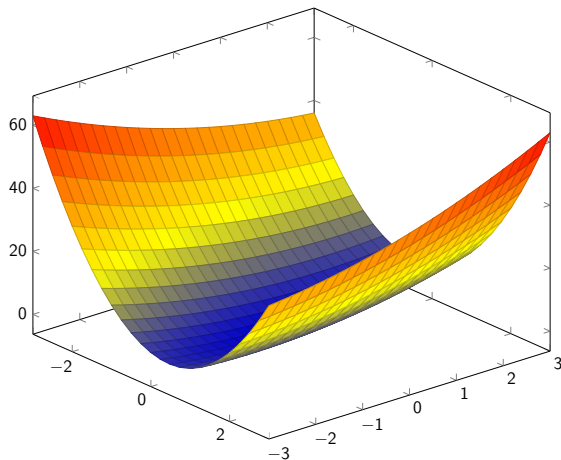10: **return** $\frac{U+L}{2}$

## Outline

1. **Convex Functions and Sets**

2. **Operations Which Preserve Convexity**

3. **First Order Properties**

4. **Minimizing a 1-d Convex Function**

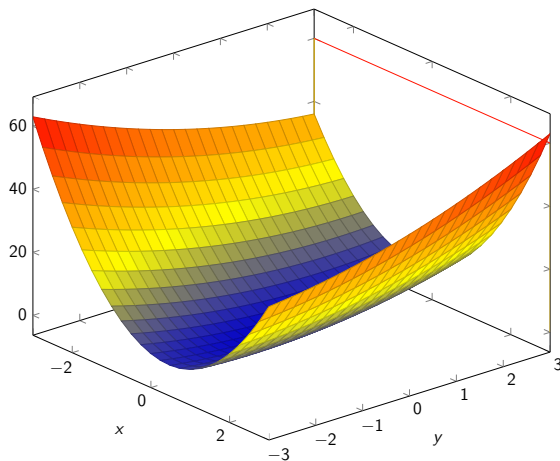5. **Coordinate Descent**

6. **Gradient Descent**

## Problem Statement



- Given a black-box which can compute $J : \mathbb{R}^n \to \mathbb{R}$ and $\nabla J : \mathbb{R}^n \to \mathbb{R}^n$ find the minimum value of $J$

## Concrete Example



$$f(x, y) = \frac{1}{2} \begin{bmatrix} x, y \end{bmatrix} \begin{bmatrix} 10, 1 \\ 2, 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

## Concrete Example

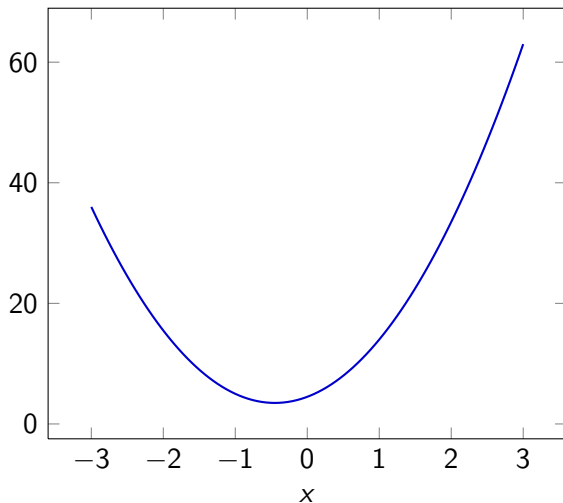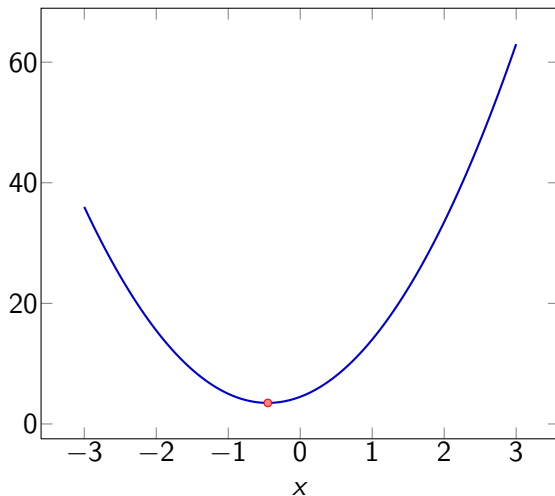

$$f(x, 3) = \frac{1}{2} \begin{bmatrix} x, 3 \end{bmatrix} \begin{bmatrix} 10, 1 \\ 2, 1 \end{bmatrix} \begin{bmatrix} x \\ 3 \end{bmatrix}$$
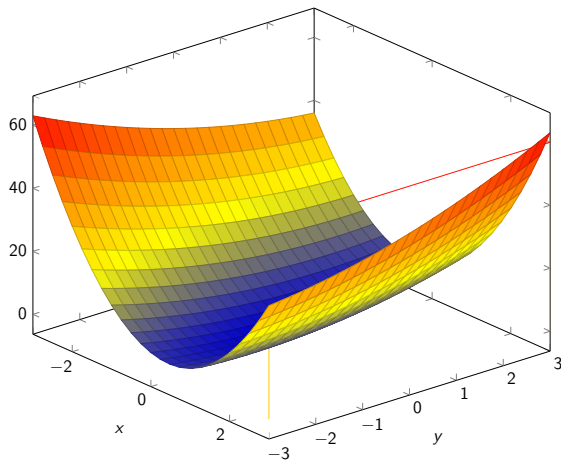
## Concrete Example



$$f(x, 3) = 5x^2 + \frac{9}{2}x + \frac{9}{2}$$
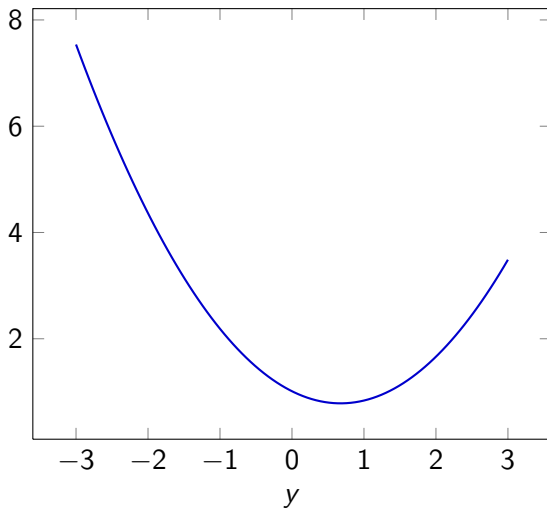
## Concrete Example



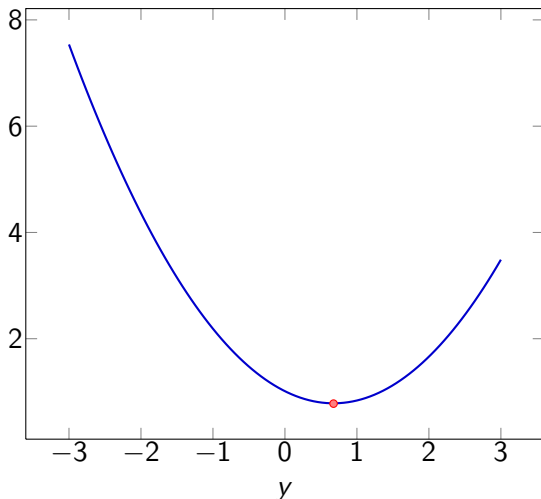$$f(x, 3) = 5x^2 + \frac{9}{2}x + \frac{9}{2} \quad \text{Minima: } x = -\frac{9}{20}$$

# Concrete Example

## Concrete Example
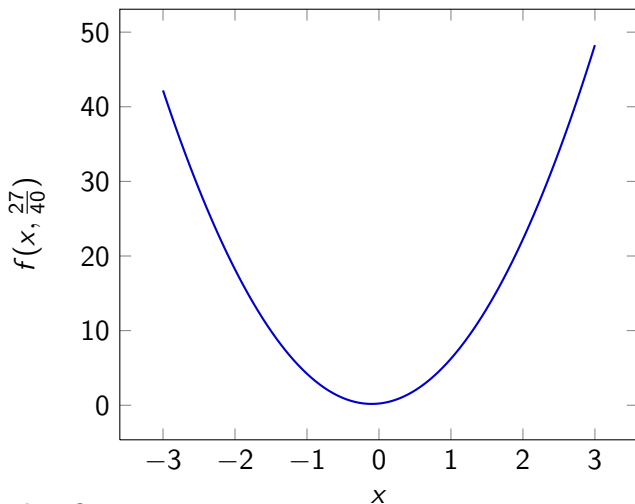


$$f(-\frac{9}{20}, y) = \frac{1}{2}y^2 - \frac{27}{40}y + \frac{81}{80}$$

## Concrete Example
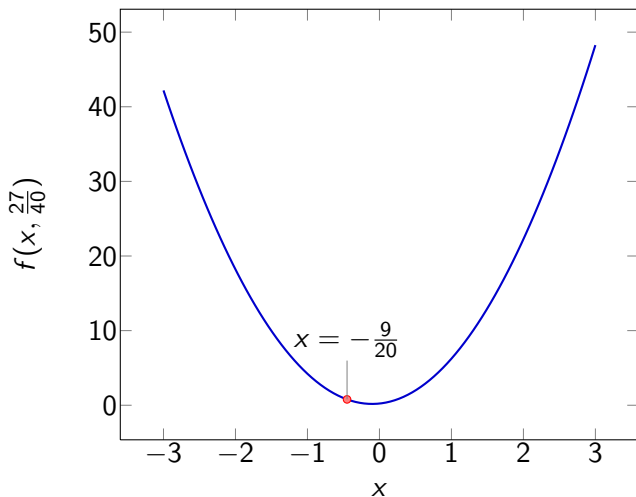


$$f(-\frac{9}{20}, y) = \frac{1}{2}y^2 - \frac{27}{40}y + \frac{81}{80} \quad \text{Minima: } y = \frac{27}{40}$$

## Concrete Example



- Are we done?

## Concrete Example



• Are we done?

## Outline

## Problem Statement



- Given a black-box which can compute $J : \mathbb{R}^n \to \mathbb{R}$ and $\nabla J : \mathbb{R}^n \to \mathbb{R}^n$ find the minimum value of $J$

**Basic Idea**

- Given a location $w_t$ at iteration $t$ update

$$w_{t+1} = w_t - \eta_t \nabla J(w_t),$$

- $\eta_t$ is a scalar stepsize

## Gradient Descent Algorithm

1: **Input:** Initial point $w_0$, gradient norm tolerance $\epsilon$
2: Set $t = 0$
3: **while** $\|\nabla J(w_t)\| \geq \epsilon$ **do**
4: $\quad w_{t+1} = w_t - \eta_t \nabla J(w_t)$
5: $\quad t = t + 1$
6: **end while**
7: **Return:** $w_t$

**Line Search Strategies - I**

- **Exact:** $J(w_t - \eta \nabla J(w_t))$ is a one dimensional convex function in $\eta$.
- **Inexact:** Armijio-Goldstein (Wolfe) conditions

$$J(w_{t+1}) \ \leq \ J(w_t) + c_1 \eta_t \langle \nabla J(w_t), w_{t+1} - w_t \rangle \ \text{(sufficient decrease)}$$
$$\langle \nabla J(w_{t+1}), w_{t+1} - w_t \rangle \ \geq \ c_2 \langle \nabla J(w_t), w_{t+1} - w_t \rangle \ \text{(curvature)}$$

with $0 < c_1 < c_2 < 1$.

**Line Search Strategies - II**

- **Decaying Stepsize:** Use a stepsize which decays according to a fixed schedule, for example, $\eta_t = 1/\sqrt{t}$
- **Fixed Stepsize:** Suppose $J$ has a Lipschitz continuous gradient with modulus $L$. Set $\eta_t = \frac{1}{L}$.