



Nonlinear Programming:

Zero-order, First-order, and Second-order Optimization Methods

Assoc. Prof. Karl Ezra Pilario, Ph.D.

Process Systems Engineering Laboratory
Department of Chemical Engineering
University of the Philippines Diliman

Outline

- Introduction to NLP
- Necessary and Sufficient Conditions for Optimality
- Convex Programming
- Methods for Solving NLP
 - One-Dimensional, Unconstrained NLP
 - Multivariable, Unconstrained NLP
 - Zero-order, First-order, Second-order Methods
 - Constrained NLP

Nonlinear Programming

- The general form of a nonlinear program (NLP) is:

Minimize: $f(\mathbf{x})$

Subject to: $h_i(\mathbf{x}) = 0 \quad i = 1, 2, \dots, l$

$g_j(\mathbf{x}) \leq 0 \quad j = 1, 2, \dots, m$

and $x_k^L \leq x_k \leq x_k^U \quad k = 1, 2, \dots, n$

- At least one of $f(\mathbf{x})$, $g_j(\mathbf{x})$, or $h_i(\mathbf{x})$ is *nonlinear*.
- x_k^L and x_k^U are lower and upper bounds on x_k .

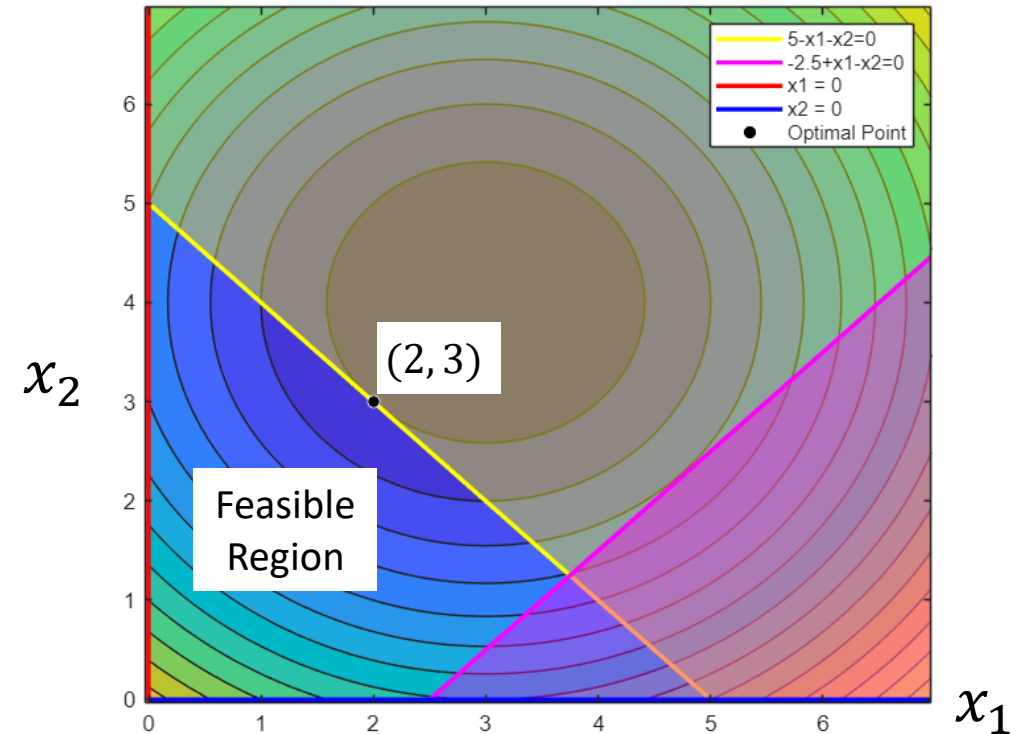
Example:

Minimize: $f(x_1, x_2) = (x_1 - 3)^2 + (x_2 - 4)^2$

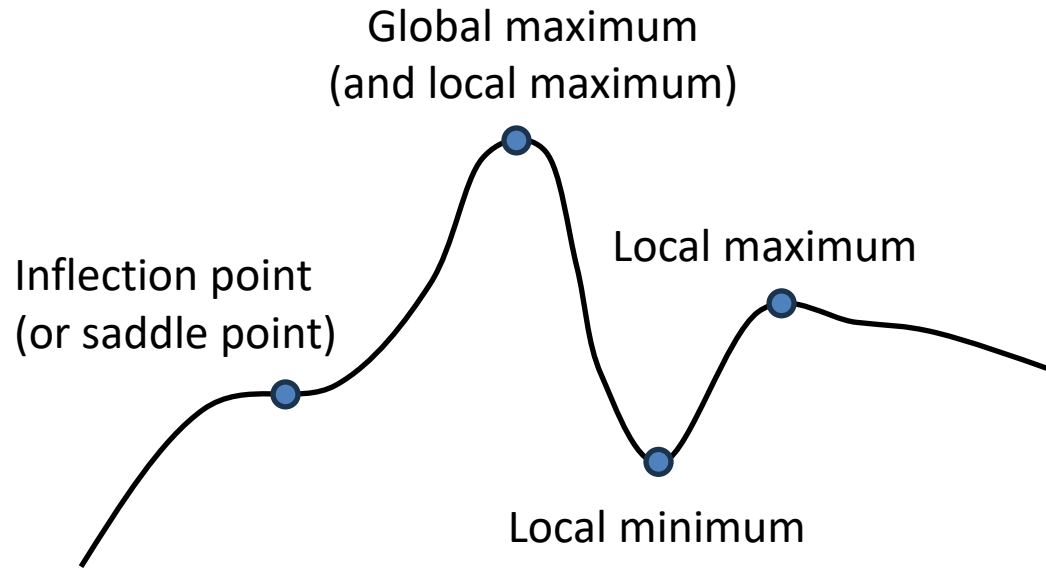
Subject to: $x_1, x_2 \geq 0$

$5 - x_1 - x_2 \geq 0$

$-2.5 + x_1 - x_2 \geq 0$



Nonlinear Programming



Local minimum

A point \mathbf{x}^* such that no other point in the *vicinity* of \mathbf{x}^* yields a value of $f(\mathbf{x})$ less than $f(\mathbf{x}^*)$:

$$f(\mathbf{x}) \geq f(\mathbf{x}^*)$$

Global minimum

A point \mathbf{x}^* such that $f(\mathbf{x}) \geq f(\mathbf{x}^*)$ holds for any \mathbf{x} in the n -dimensional space of \mathbf{x} .

Conditions for Optimality in Unconstrained NLPs

A point \mathbf{x}^* is an **optimal solution** to an unconstrained NLP:

- ONLY IF $f(\mathbf{x})$ is twice differentiable at \mathbf{x}^* . **NECESSARY**
- ONLY IF $\nabla f(\mathbf{x}^*) = 0$ or \mathbf{x}^* is a stationary point. **NECESSARY**
- IF $\mathbf{H}(\mathbf{x}^*)$ is *positive-definite* for a minimum to exist, and *negative-definite* for a maximum to exist. **SUFFICIENT**

$\mathbf{H}(\mathbf{x})$ is the Hessian matrix.

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

If all the eigenvalues of $\mathbf{H}(\mathbf{x})$ are...

$\mathbf{H}(\mathbf{x})$ is...

Positive-definite	> 0
Positive-semidefinite	≥ 0
Negative-definite	< 0
Negative-semidefinite	≤ 0

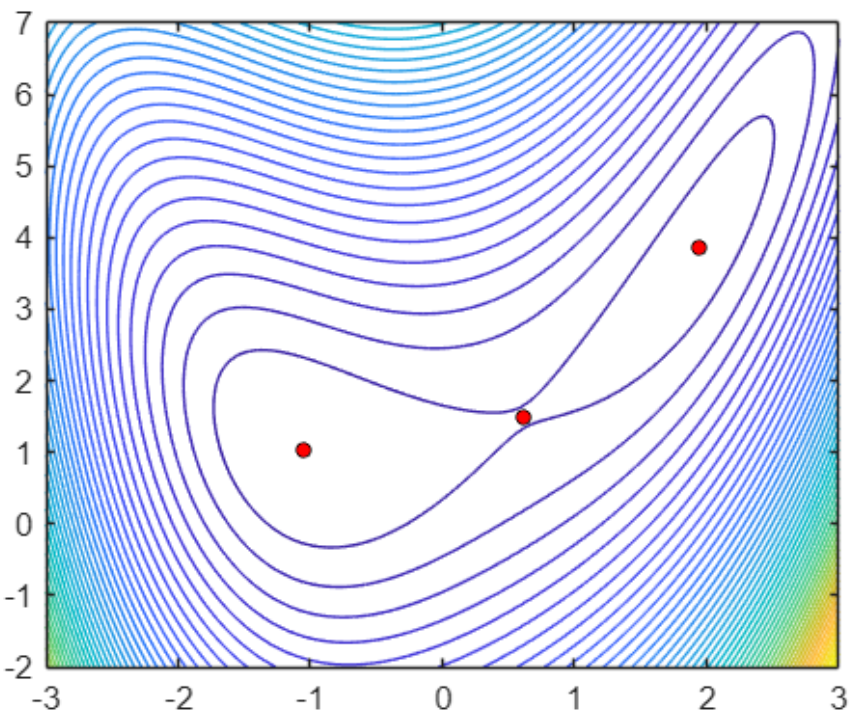
Nonlinear Programming

Example:

$$f(x_1, x_2) = 4 + 4.5x_1 - 4x_2 + x_1^2 + 2x_2^2 - 2x_1x_2 + x_1^4 - 2x_1^2x_2$$

This function has three stationary points (1.941, 3.854), (−1.053, 1.028), and (0.6117, 1.4929). Evaluate the optimality at these points.

$$H(x_1, x_2) = \begin{bmatrix} 2 + 12x_1^2 - 4x_2 & -2 - 4x_1 \\ -2 - 4x_1 & 4 \end{bmatrix}$$



Stationary point	$f(x)$	$f'(x)$	$H(x)$	Eigenvalues		Evaluation
(1.941, 3.854)	0.9856	0.00	$\begin{bmatrix} 31.79 & -9.76 \\ -9.76 & 4.00 \end{bmatrix}$	0.9128	34.8810	Local minimum
(−1.053, 1.028)	−0.5134	0.00	$\begin{bmatrix} 11.19 & 2.21 \\ 2.21 & 4.00 \end{bmatrix}$	3.3743	11.8194	Local minimum
(0.6117, 1.4929)	2.8091	0.00	$\begin{bmatrix} 0.52 & -4.45 \\ -4.45 & 4.00 \end{bmatrix}$	-2.5161	7.0346	Saddle point

Nonlinear Programming

Unlike an LP, solving an NLP is hard because a minimum is **not guaranteed** to be the global minimum.

But if the NLP is **convex**, we have the following result.

Theorem

For an NLP whose objective function $f(x)$ is a **convex function** and each inequality constraint $g(x)$ is a convex function so that they form a **convex set**:

The local minimum of $f(x)$ is also a global minimum.

Analogously, if $f(x)$ is a **concave function** and the constraints $g(x)$ still form a **convex set**, then:

The local maximum of $f(x)$ is also a global maximum.

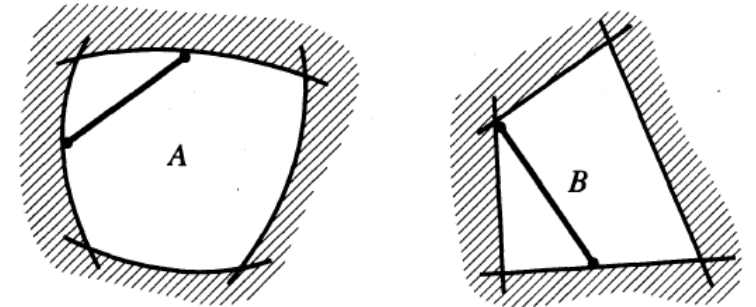
Definition: Convex Set

For every pair of points x_1 and x_2 in a convex set, the point x given by a linear combination of the two points:

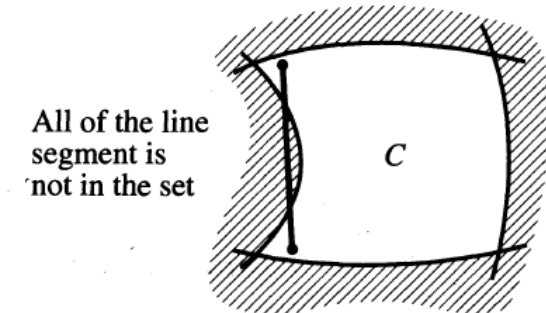
$$x = \gamma x_1 + (1 - \gamma)x_2, \quad 0 \leq \gamma \leq 1$$

must also be in the set. Also, the intersection of any number of convex sets is a convex set.

Examples of a
Convex Set



Example of a
Nonconvex Set



Nonlinear Programming

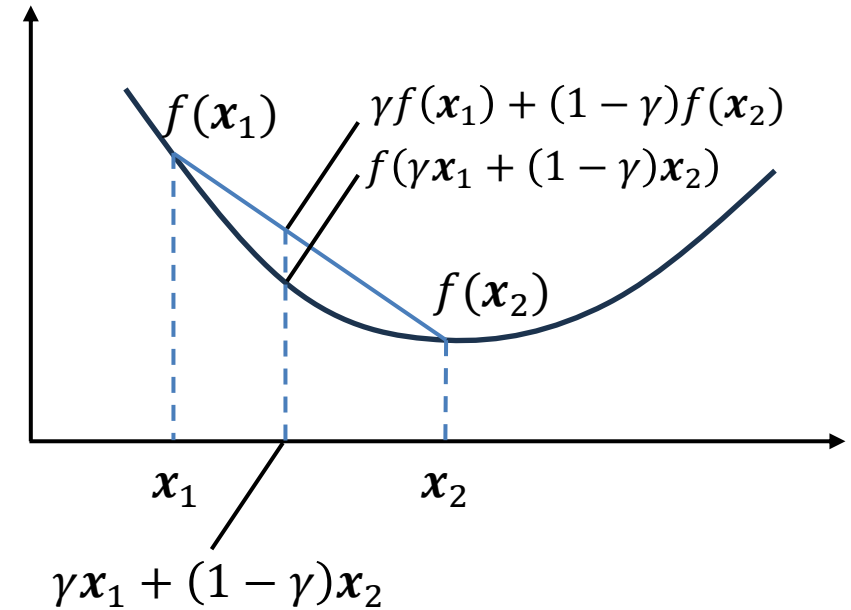
Definition: Convex Function

A function $f(x)$ defined on a convex set F is said to be a **convex function** if the following relation holds:

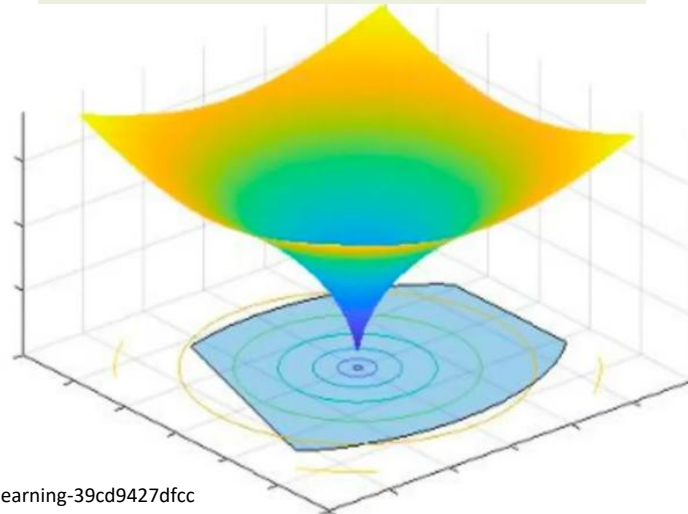
$$f(\gamma x_1 + (1 - \gamma)x_2) \leq \gamma f(x_1) + (1 - \gamma)f(x_2)$$

Where $0 \leq \gamma \leq 1$.

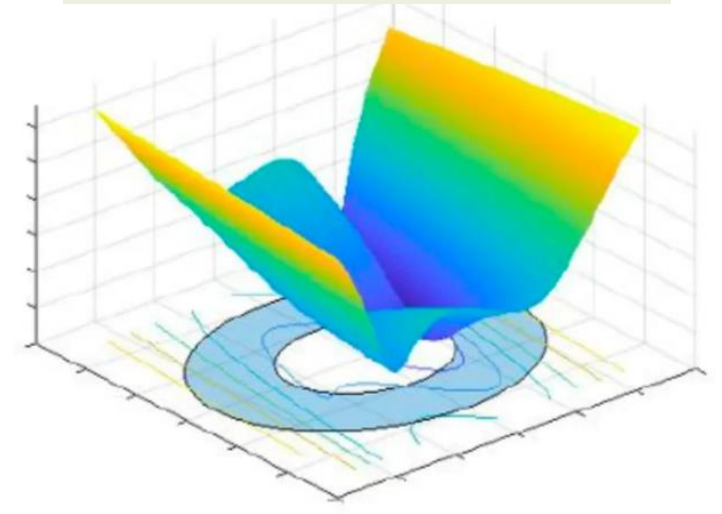
- If $f(x)$ is convex, then $-f(x)$ is concave.
- If only the inequality sign holds ($<$), the function is said to be *strictly convex*.
- If $f(x)$ is strictly convex, then $-f(x)$ is strictly concave.



A **convex** objective with
convex constraints.

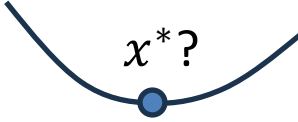


A **nonconvex** objective with
nonconvex constraints.

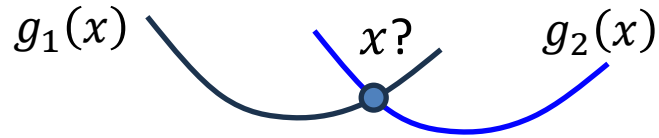


One-Dimensional, Unconstrained NLP

Minimize: $f(x)$



Recall: Newton's Method for Root-Finding



Given: $f(x) = g_1(x) - g_2(x) = 0$

Initialization: One initial guess (x_0)
and a Tolerance, Tol

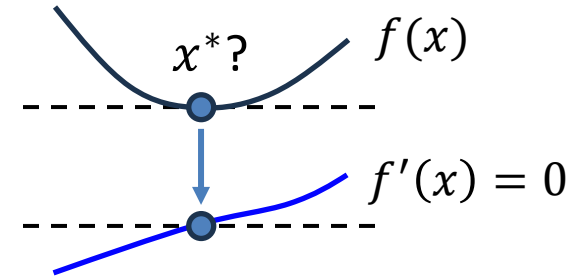
- 1: Set $x_{old} \leftarrow \text{Inf}$.
- 2: Set $x_{new} \leftarrow x_0$.
- 3: **WHILE** $|x_{new} - x_{old}| > Tol$:
- 4: Assign: $x_{old} \leftarrow x_{new}$.
- 5: Assign: $x_{new} \leftarrow x_{old} - \frac{f(x_{old})}{f'(x_{old})}$.
- 6: **END WHILE**
- 7: **OUTPUT** x_{new}

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

Newton's Method for Optimization

Minimize: $f(x)$

Find the roots:
 $f'(x) = 0$



Newton's
Method

$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}$$

Quasi-Newton Method

$$x_{k+1} = x_k - \frac{(x_{k-1} - x_k)f'(x_k)}{f'(x_{k-1}) - f'(x_k)}$$

$$f''(x_k) \approx \frac{f'(x_{k-1}) - f'(x_k)}{x_{k-1} - x_k} \quad \text{2nd Order finite difference}$$

One-Dimensional, Unconstrained NLP

Example:

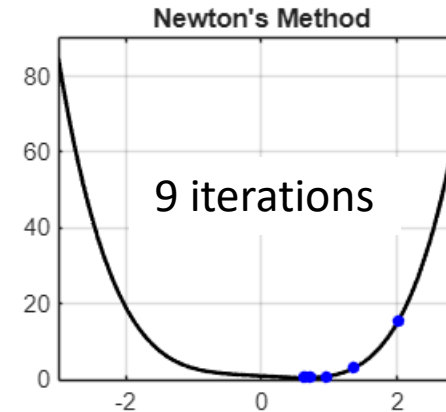
Minimize: $f(x) = x^4 - x + 1$

- Use **Newton's Method** starting with $x_0 = 3$.
- Use **Quasi-newton method** starting with $x_0 = 3$ and $x_1 = 2$.
- For both methods, iterate until the change in x is less than 10^{-7} .

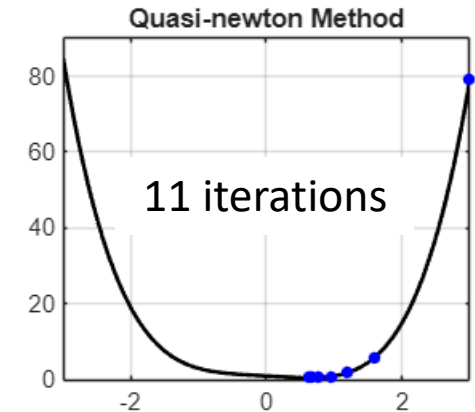
```
f = @(x) x.^4-x+1;  
df = @(x) 4*x.^3-1;  
ddf = @(x) 12*x.^2;
```

```
% Use Newton's method  
x0 = 3; x1 = Inf; ctr = 1;  
while abs(x0-x1) > 1e-7  
    x1 = x0;  
    x0 = x1 - df(x1)/ddf(x1);  
    ctr = ctr + 1;  
end  
disp(x0);  
fprintf('%d\n', ctr);
```

```
% Use quasi-newton method  
x0 = 3; x1 = 2; ctr = 1;  
while abs(x0-x1) > 1e-7  
    xn = x1 - (x0-x1)*df(x1)/(df(x0)-df(x1));  
    x0 = x1; x1 = xn;  
    ctr = ctr + 1;  
end  
disp(x0);  
fprintf('%d\n', ctr);
```



$$x^* = 0.63$$



$$x^* = 0.63$$

Outline

- Introduction to NLP
- Necessary and Sufficient Conditions for Optimality
- Convex Programming
- Methods for Solving NLP
 - One-Dimensional, Unconstrained NLP
 - **Multivariable, Unconstrained NLP**
 - **Zero-order, First-order, Second-order Methods**
 - Constrained NLP

A Taxonomy of NLP Solvers

Zero-order Methods

Derivative-free, black-box

- Grid Search / Exhaustive Search
- Random Search
- Nelder-Mead Simplex
- Metaheuristic Search
 - Genetic Algorithms
 - Particle Swarm
 - Simulated Annealing
 - Differential Evolution
 - CMAES
- Bayesian Optimization / Surrogate-based

First-order Methods

Uses 1st derivative information

- Steepest Descent or Gradient Descent
- Stochastic Gradient Descent (SGD)
 - RMSProp, Adam, etc.
- Conjugate Gradient methods (e.g. Fletcher-Reeves)
- Coordinate Descent

Second-order Methods

Uses 1st and 2nd derivative information

- *Newton's method*
- *Quasi-newton method (BFGS)*
- Levenberg-Marquardt
- Trust Region Newton Methods
- Active Set and Interior-Point Methods
- Sequential Quadratic Programming (SQP)
- IPOPT

Multivariable, Unconstrained NLP

- **Steepest Descent or Gradient Descent**
- Stochastic Gradient Descent (SGD)
 - RMSProp, Adam, etc.
- Conjugate Gradient methods (e.g. Fletcher-Reeves)
- Coordinate Descent

Steepest Descent / Gradient Descent

At each iteration, descend at the direction of the greatest rate of decrease in $f(\mathbf{x})$.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent

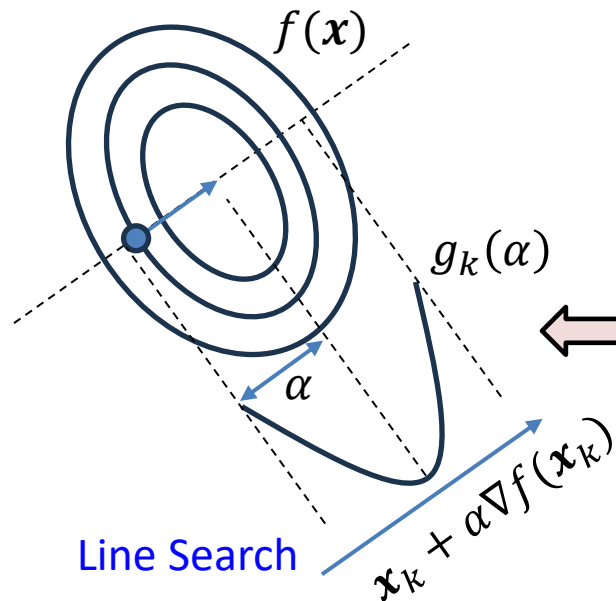
Steps:

- Make an initial guess, \mathbf{x}_0
- Calculate the search direction $\nabla f(\mathbf{x}_k)$
- Calculate step size α by minimizing:

$$g_k(\alpha) = f(\mathbf{x}_k + \alpha \nabla f(\mathbf{x}_k))$$

(e.g. using quasi-newton method)

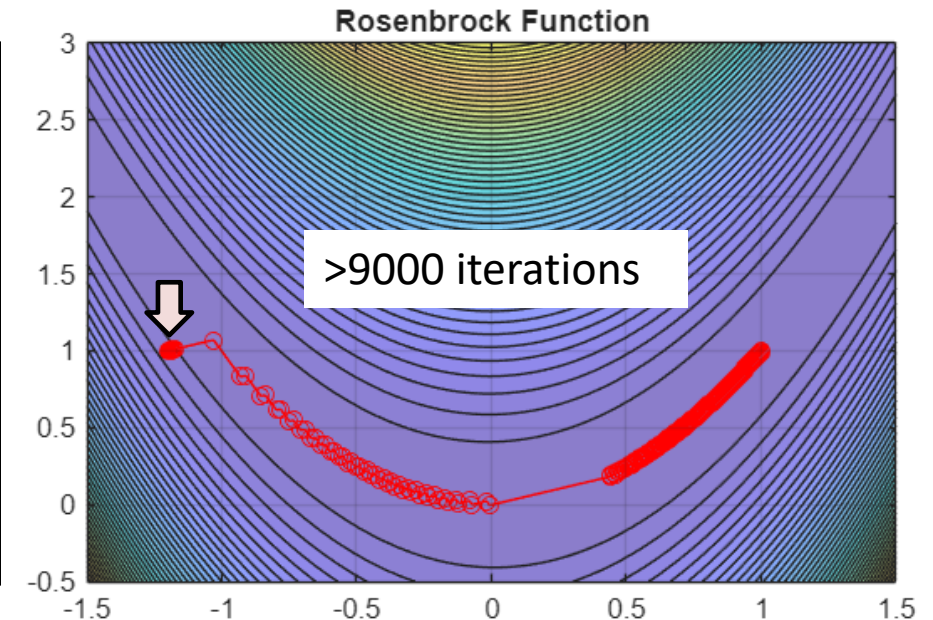
- Update \mathbf{x}_k



Example:

Minimize $f(\mathbf{x})$ from $\mathbf{x}_0 = (-1.2, 1.0)$:

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$$



Multivariable, Unconstrained NLP

- Steepest Descent or Gradient Descent
- **Stochastic Gradient Descent (SGD)**
 - **RMSProp, Adam, etc.**
- Conjugate Gradient methods (e.g. Fletcher-Reeves)
- Coordinate Descent

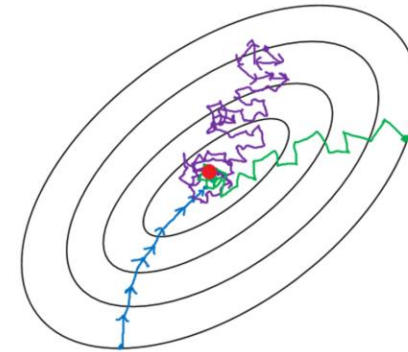
Batch vs. Stochastic Gradient Descent

Gradient descent (GD) is used to train neural nets.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent

- **Batch GD** = all N data pts are used in $C(\mathbf{W}, \mathbf{b})$
- **Mini-batch GD** = a batch of $M < N$ data pts are used.
- **Stochastic GD** = only one random data pt is used.



Also, α_k is now called a learning rate. **RMSProp** and **Adam** are ways to decay α_k at every k .

To train **neural nets** (from *machine learning*), we need to minimize $C(\mathbf{W}, \mathbf{b})$:

$$\min_{\mathbf{W}, \mathbf{b}} C(\mathbf{W}, \mathbf{b}) = \frac{1}{N} \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2$$

$C(\mathbf{W}, \mathbf{b})$ = cost (mean squared error)

$[\mathbf{W}; \mathbf{b}]$ = weights and biases

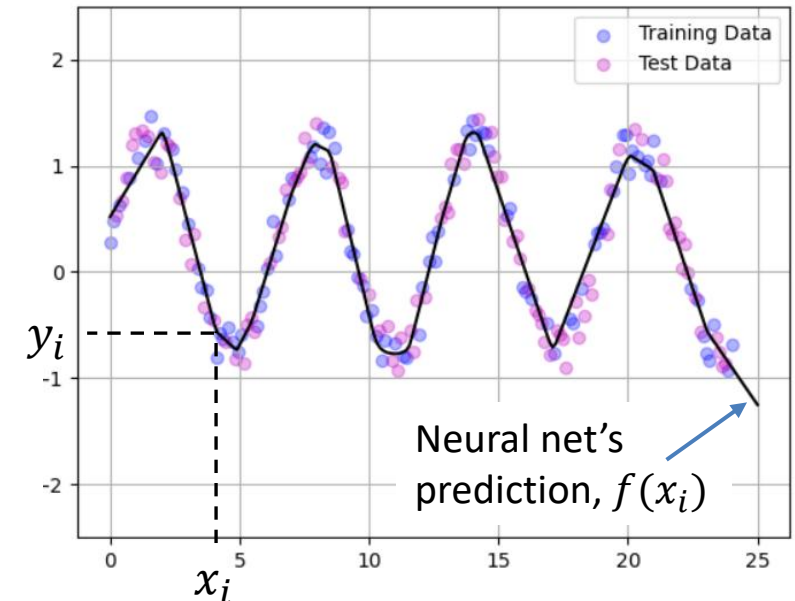
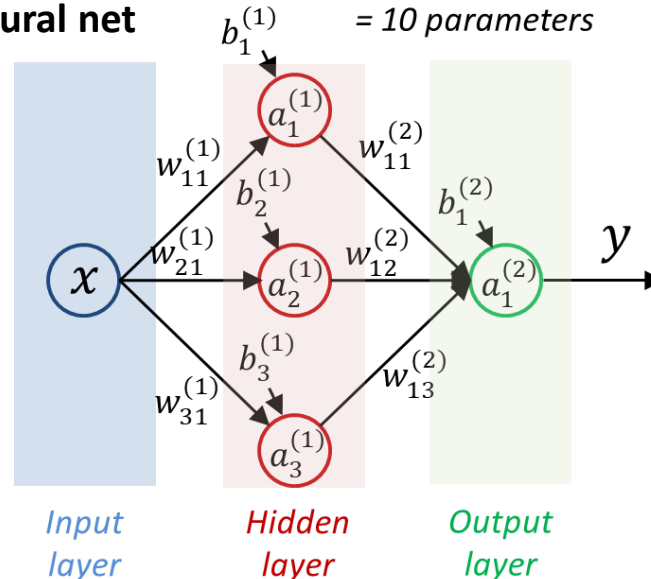
y_i = true values to be predicted

$f(\mathbf{x}_i)$ = neural net's prediction

N = no. of data points

Example: 1 hidden-layer neural net

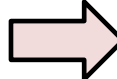
6 weights, 4 biases
= 10 parameters



Multivariable, Unconstrained NLP

- Steepest Descent or Gradient Descent
- Stochastic Gradient Descent (SGD)
 - RMSProp, Adam, etc.
- **Conjugate Gradient methods (e.g. Fletcher-Reeves)**
- Coordinate Descent

Conjugate Gradient Descent

This is the **steepest descent** update:  This is the **conjugate gradient** update:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent

$$\mathbf{s}_0 = -\nabla f(\mathbf{x}_0) \quad (\text{search direction})$$

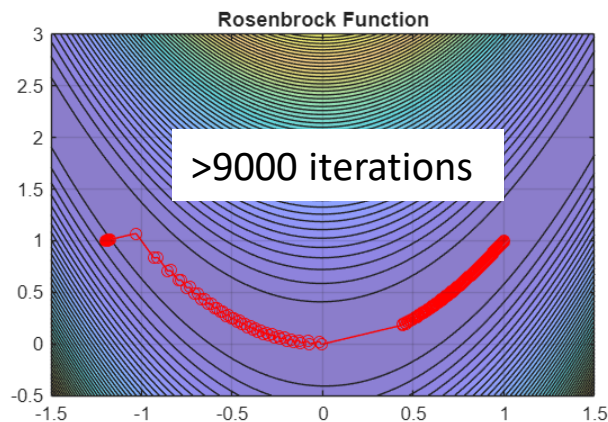
$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{s}_k$$

$$\mathbf{s}_{k+1} = -\nabla f(\mathbf{x}_{k+1}) + \mathbf{s}_k \frac{\nabla f(\mathbf{x}_{k+1})^T \nabla f(\mathbf{x}_{k+1})}{\nabla f(\mathbf{x}_k)^T \nabla f(\mathbf{x}_k)}$$

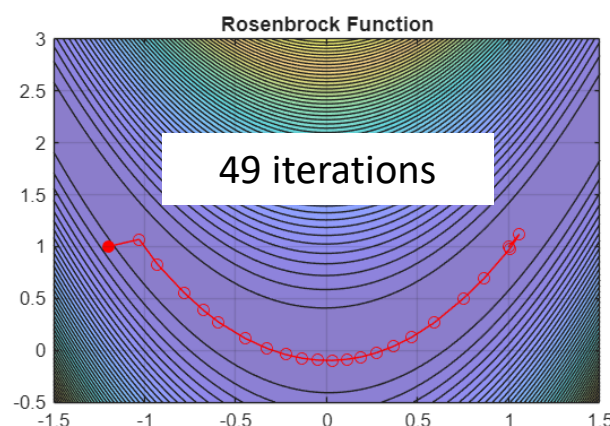
Example:

Minimize $f(\mathbf{x})$ from $\mathbf{x}_0 = (-1.2, 1.0)$: $f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$

Steepest Descent



Conjugate Gradient Descent

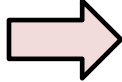


- New search directions are *conjugate* to previous ones.
- Prevents “zigzagging”.
- Prevents undoing the progress made in previous steps.
- Allows convergence in at most N steps for an N -dimensional quadratic function.

Multivariable, Unconstrained NLP

- Steepest Descent or Gradient Descent
- Stochastic Gradient Descent (SGD)
 - RMSProp, Adam, etc.
- Conjugate Gradient methods (e.g. Fletcher-Reeves)
- **Coordinate Descent**

Coordinate Descent

This is the **steepest descent** update:  Like steepest descent, but the search direction is aligned to one coordinate at a time.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \frac{\partial f(\mathbf{x}_k)}{\partial x_i}$$

What's the difference?

$$k = 1, \quad 2, \quad 3, \dots$$
$$\nabla f(\mathbf{x}_k) \quad \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix}, \begin{bmatrix} -4 \\ 2 \\ 0 \end{bmatrix}, \begin{bmatrix} -2 \\ -5 \\ 1 \end{bmatrix} \dots$$
$$\frac{\partial f(\mathbf{x}_k)}{\partial x_i} \quad \begin{bmatrix} -3 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ -2 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix} \dots$$

Steps:

- Make an initial guess, \mathbf{x}_0
- Calculate the search direction $\nabla f(\mathbf{x}_k)$
- Calculate step size α by minimizing:

$$g_k(\alpha) = f(\mathbf{x}_k + \alpha \nabla f(\mathbf{x}_k))$$

(e.g. using quasi-newton method)

- Update \mathbf{x}_k

Steps:

- Make an initial guess, \mathbf{x}_0
- Choose any coordinate among x_i .
- Calculate step size α by minimizing:

$$g_k(\alpha) = f\left(\mathbf{x}_k + \alpha \frac{\partial f(\mathbf{x}_k)}{\partial x_i}\right)$$

(e.g. using quasi-newton method)

- Update \mathbf{x}_k

Multivariable, Unconstrained NLP

- Steepest Descent or Gradient Descent
- Stochastic Gradient Descent (SGD)
 - RMSProp, Adam, etc.
- Conjugate Gradient methods (e.g. Fletcher-Reeves)
- **Coordinate Descent**

Coordinate Descent

This is the **steepest descent** update:  Like steepest descent, but the search direction is aligned to one coordinate at a time.

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

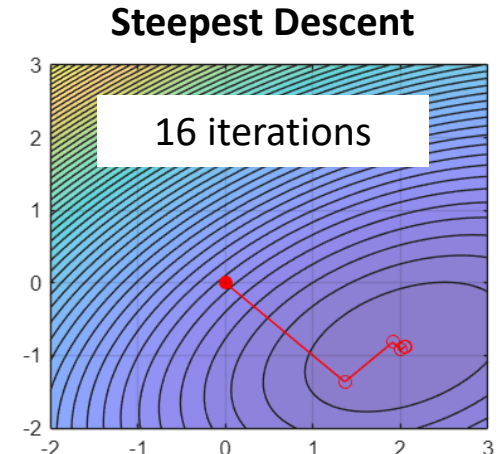
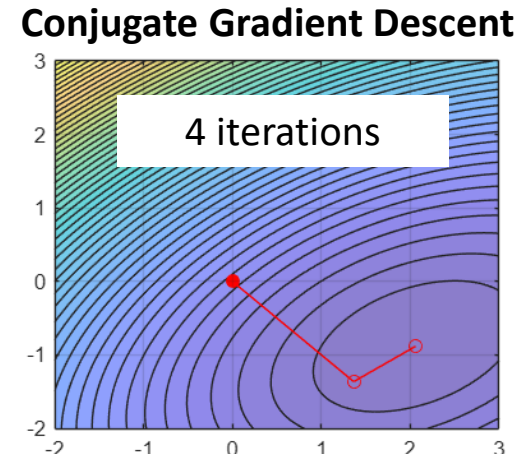
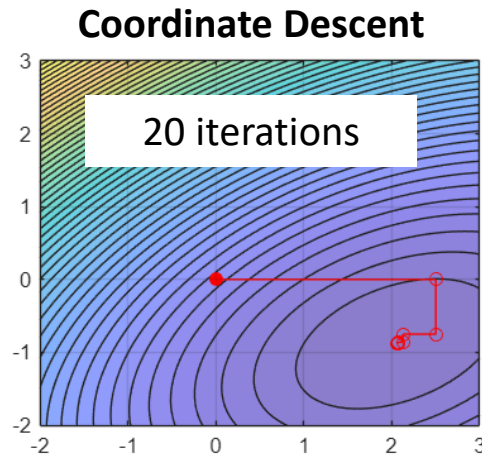
α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \frac{\partial f(\mathbf{x}_k)}{\partial x_i}$$

Example:

Minimize $f(\mathbf{x})$ from $\mathbf{x}_0 = (0.0, 0.0)$:

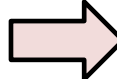
$$f(\mathbf{x}) = 0.06(x_1^2 - x_1x_2) + 0.1x_2^2 - 0.3(x_1 - x_2)$$



Multivariable, Unconstrained NLP

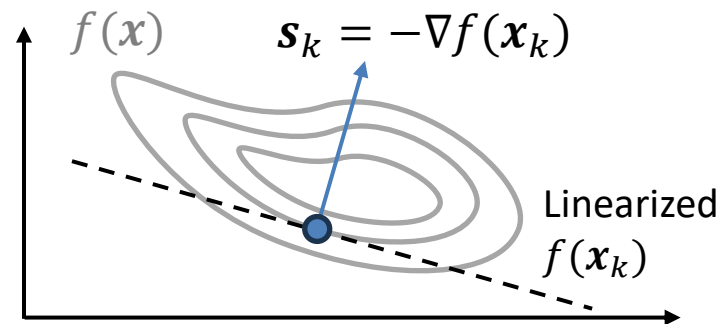
- Newton's method
- Quasi-newton Method (BFGS)
- Levenberg-Marquardt
- Trust Region Newton Methods
- Active Set and Interior-Point Methods
- Sequential Quadratic Programming (SQP)
- IPOPT

Multivariable Quasi-Newton Methods

This is the **steepest descent** update:  This is the **Newton** update:

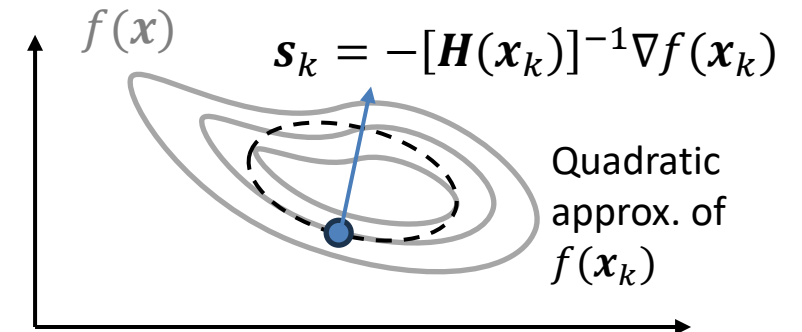
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent



$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k [\mathbf{H}(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent
 $\mathbf{H}(\mathbf{x}_k)$ = Hessian matrix

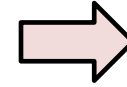


Multivariable, Unconstrained NLP

- Newton's method
- Quasi-newton Method (BFGS)
- Levenberg-Marquardt
- Trust Region Newton Methods
- Active Set and Interior-Point Methods
- Sequential Quadratic Programming (SQP)

Multivariable Quasi-Newton Methods

This is the **Newton** update:



$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k [\mathbf{H}(\mathbf{x}_k)]^{-1} \nabla f(\mathbf{x}_k)$$

α_k = step size at the k th iteration
 $\nabla f(\mathbf{x}_k)$ = direction of steepest descent
 $\mathbf{H}(\mathbf{x}_k)$ = Hessian matrix

Quasi-Newton: BFGS Algorithm (Broyden-Fletcher-Goldfarb-Shanno)

Steps:

- Make an initial guess, \mathbf{x}_0 , and initial guess of inverse Hessian $\mathbf{\Gamma}_0 = \mathbf{I}$. [*steepest descent*]
- Calculate *search direction*, $\mathbf{s}_k = -\mathbf{\Gamma}_k \nabla f(\mathbf{x}_k)$
- Calculate step size, α , using line search:

$$g_k(\alpha) = f(\mathbf{x}_k + \alpha \mathbf{s}_k)$$

- Update $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{s}_k$.
- Update $\mathbf{\Gamma}_k$ to $\mathbf{\Gamma}_{k+1}$.

$\mathbf{\Gamma}_k$ = k th approximation of the inverse Hessian

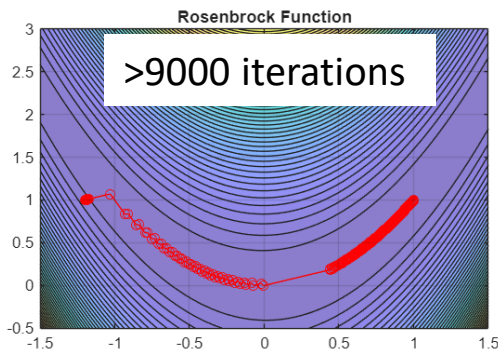
$$\mathbf{\Gamma}_{k+1} = \left(\mathbf{I} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) \mathbf{\Gamma}_k \left(\mathbf{I} - \frac{\mathbf{y}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{y}_k^T \mathbf{s}_k}$$

where $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$

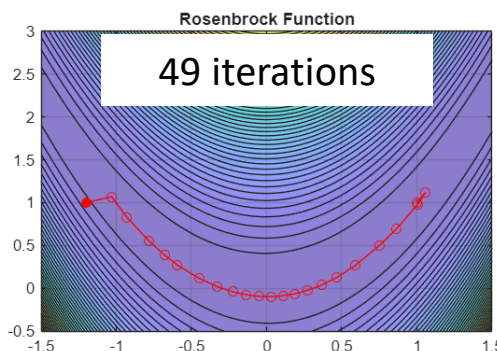
Example:

Minimize $f(\mathbf{x})$ from $\mathbf{x}_0 = (-1.2, 1.0)$: $f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$

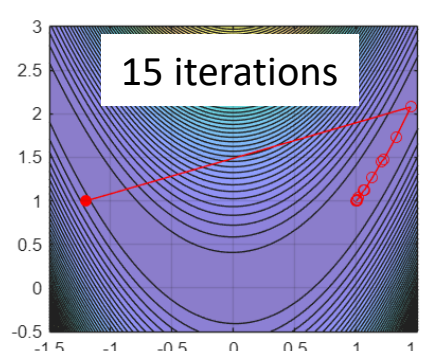
Steepest Descent



Conjugate Gradient Descent



BFGS



Outline

- Introduction to NLP
- Necessary and Sufficient Conditions for Optimality
- Convex Programming
- Methods for Solving NLP
 - One-Dimensional, Unconstrained NLP
 - Multivariable, Unconstrained NLP
 - Zero-order, First-order, Second-order Methods
 - Constrained NLP

Class Exercises

1. Minimize the **Beale function** from $\mathbf{x}_0 = (-2, -2)$:

$$f(\mathbf{x}) = (1.5 - x_1 + x_1x_2)^2 + (2.25 - x_1 + x_1x_2^2)^2 + (2.625 - x_1 + x_1x_2^3)^2$$

2. Minimize the **Booth function** from $\mathbf{x}_0 = (0, 0)$:

$$f(\mathbf{x}) = (x_1 + 2x_2 - 7)^2 + (2x_1 + x_2 - 5)^2$$

Answers:

1. $f(3, 0.5) = 0$
2. $f(1, 3) = 0$