



THE GEORGE
WASHINGTON
UNIVERSITY

WASHINGTON, DC

CS 4364/6364

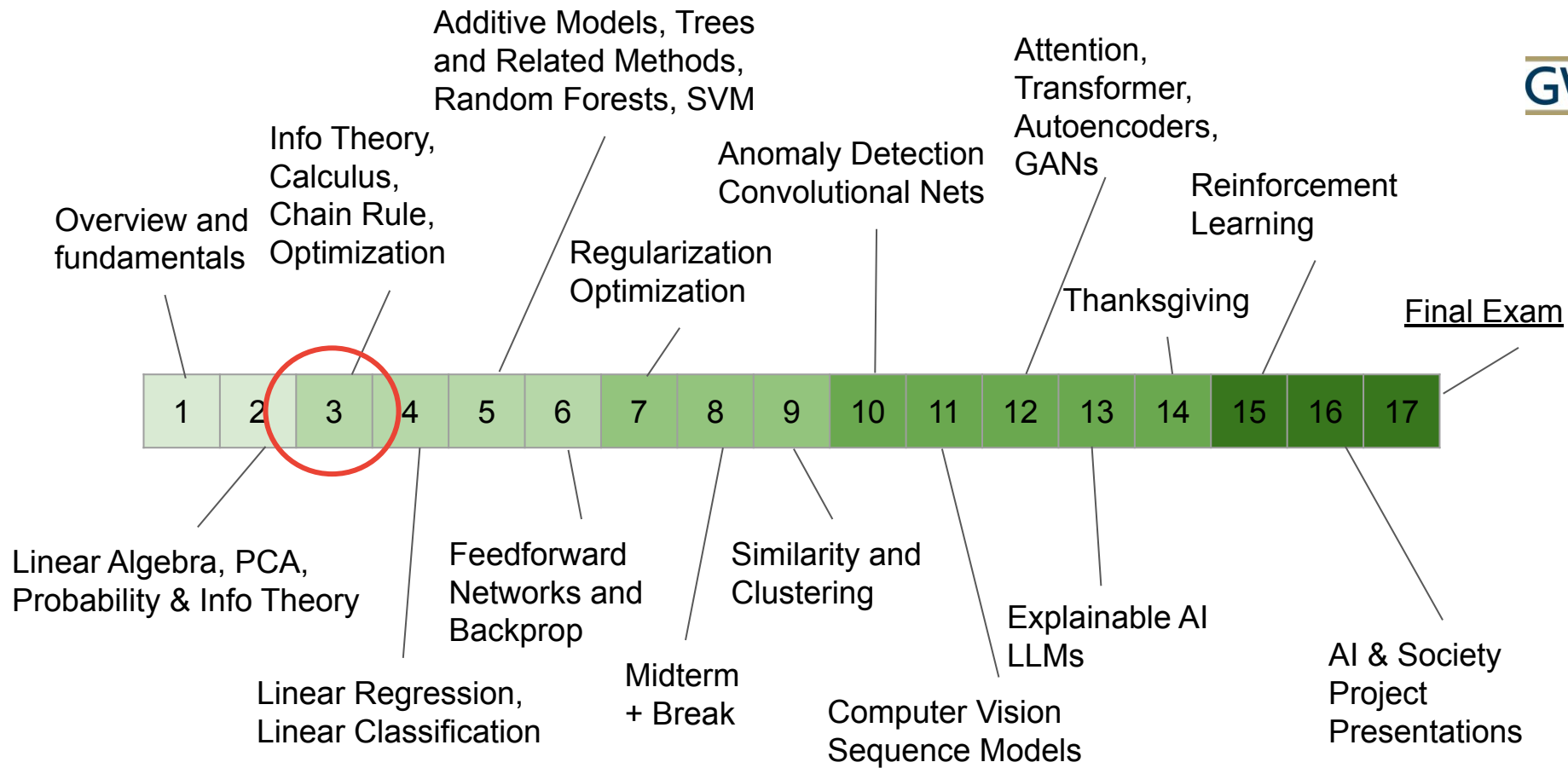
Machine Learning

Fall Semester 9/7/2023

Lecture 5.

Numerical Computation and Gradient Optimization

John Sipple
jsipple@gwu.edu



Overflow & Underflow

Overflow and Underflow

Optimization often runs up to the finite precision limits of the computer

Underflow:

- Very small number rounds down to zero
- Can cause problems with division or logarithm

Overflow:

- Large magnitudes approximated as +/- infinity
- Next operation will usually fail with Not A Number...

Gradient-Based Optimization

Optimization

Task of minimizing or maximizing some function $f(\mathbf{x})$ by altering x

$f(\mathbf{x})$ is called a **cost function**, **loss function**, or **error function**

Optimum: $\mathbf{x}^* = \arg \min f(\mathbf{x})$ (the point where $f(x)$ is least)

The Derivative – a tool for optimization

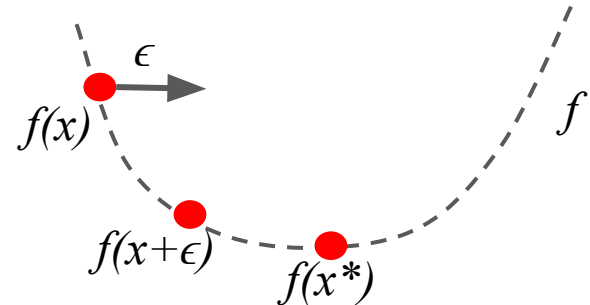
For a real-valued function $y = f(x)$, the derivative is

$$f'(x) \text{ or } \frac{dy}{dx}$$

Approximates the value of f some small displacement ϵ from x :

$$f(x + \epsilon) \approx f(x) + \epsilon f'(x)$$

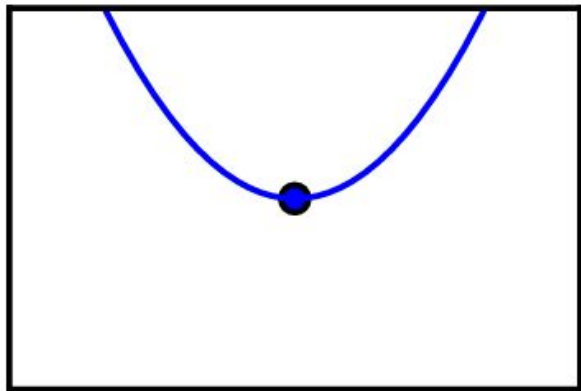
Gradient descent - follow the derivative to x^*



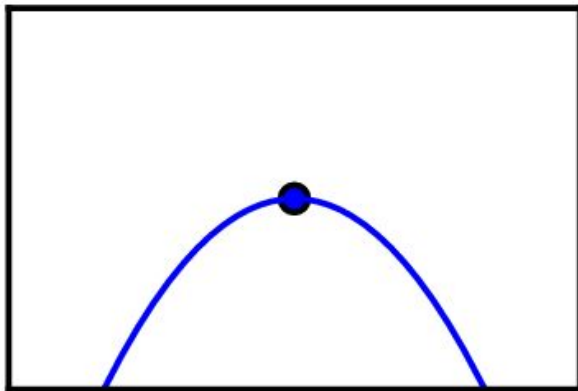
Critical Points

Points where $f'(x) = 0$

Minimum



Maximum



Saddle point

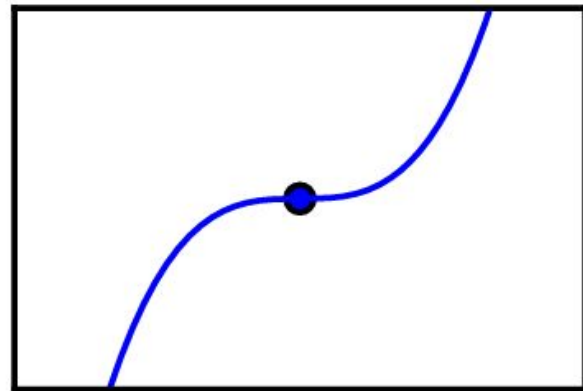


Figure 4.2

Approximate Optimization

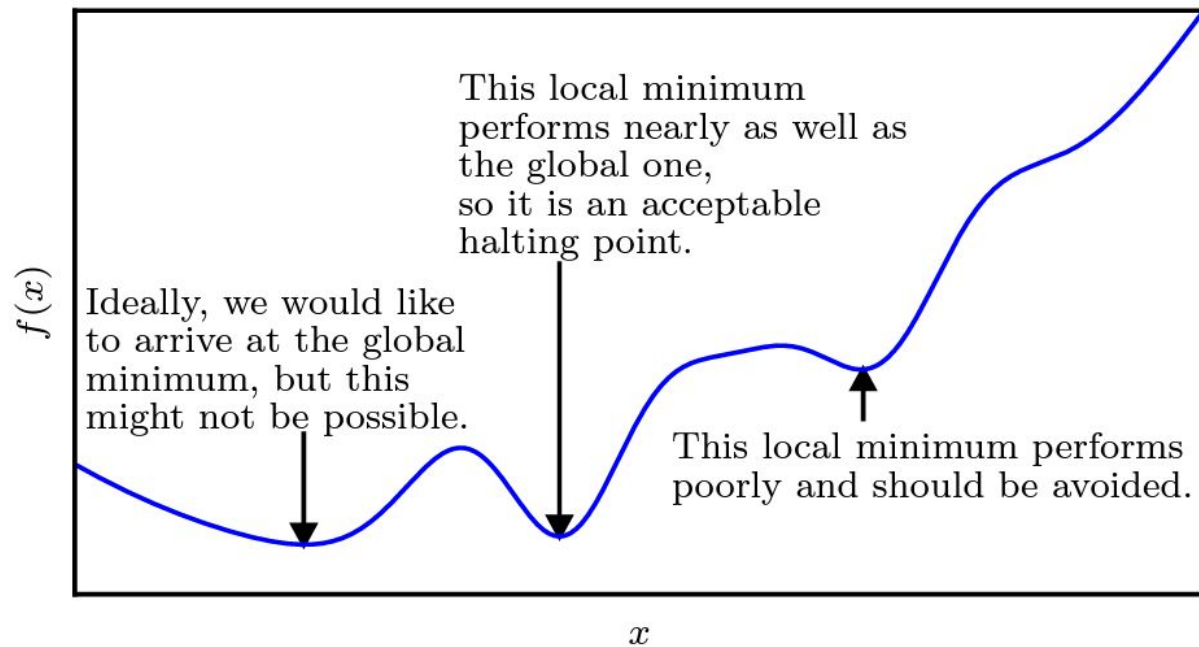


Figure 4.3

The gradient

ML Loss functions have multiple inputs:

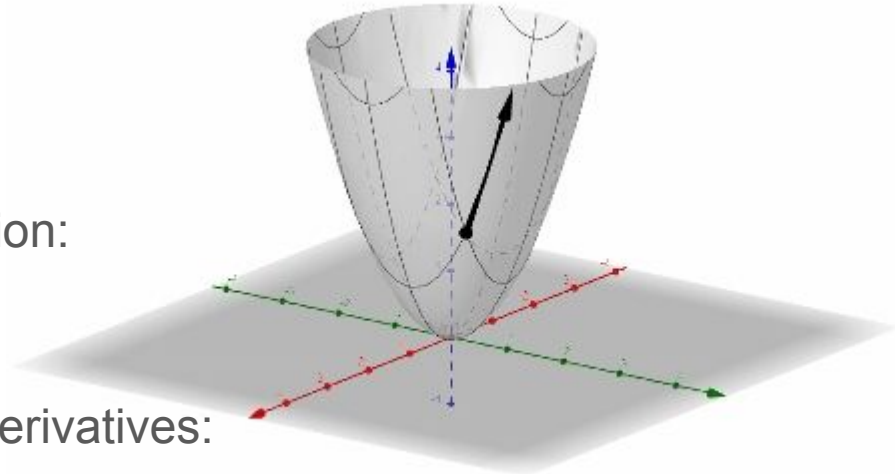
$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

Partial derivative with respect to i^{th} dimension:

$$\frac{\partial}{\partial x_i} f(\mathbf{x})$$

Gradient: n -dimensional vector of partial derivatives:

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x}) \\ \frac{\partial}{\partial x_2} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial x_n} f(\mathbf{x}) \end{bmatrix}$$



Source:
<https://study.com/academy/lesson/directional-derivatives-gradient-of-f-and-the-min-max.html>

Tracking to the minimum with the Gradient

Define a unit vector \mathbf{u} in the direction of the gradient:

$$\frac{\partial}{\partial \epsilon} f(\mathbf{x} + \epsilon \mathbf{u})$$

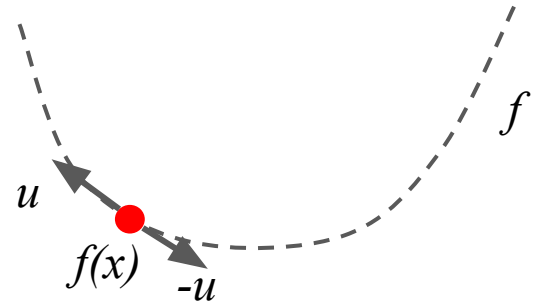
When $\epsilon = 0$, evaluates to $\mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x})$

Since we want to minimize f :

$$\min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \mathbf{u}^\top \nabla_{\mathbf{x}} f(\mathbf{x}) = \min_{\mathbf{u}, \mathbf{u}^\top \mathbf{u} = 1} \|\mathbf{u}\|_2 \|\nabla_{\mathbf{x}} f(\mathbf{x})\|_2 \cos \theta$$

Min achieved when following the opposite direction:

$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x})$$



The Jacobian

If the function has m inputs and n outputs:

$$\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

The derivative m -by- n matrix is called the Jacobian $\mathbf{J} \in \mathbb{R}^{n \times m}$ of \mathbf{f} :

$$\mathbf{J} = \begin{bmatrix} \frac{\partial}{\partial x_1} f(\mathbf{x})_1 & \frac{\partial}{\partial x_2} f(\mathbf{x})_1 & \cdots & \frac{\partial}{\partial x_m} f(\mathbf{x})_1 \\ \frac{\partial}{\partial x_1} f(\mathbf{x})_2 & \frac{\partial}{\partial x_2} f(\mathbf{x})_2 & \cdots & \frac{\partial}{\partial x_m} f(\mathbf{x})_2 \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial}{\partial x_1} f(\mathbf{x})_n & \frac{\partial}{\partial x_2} f(\mathbf{x})_n & \cdots & \frac{\partial}{\partial x_m} f(\mathbf{x})_n \end{bmatrix}$$

The Hessian

Derivative of a derivative, **Curvature**, is the rate of change of slope.

For a loss function:

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

We have an n -by- n matrix:

$$\mathbf{H}(f)(\mathbf{x})_{i,j} = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} f(\mathbf{x}) & \frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_m} f(\mathbf{x}) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_2^2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f(\mathbf{x}) \\ \cdots & & & \\ \frac{\partial^2}{\partial x_n \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_n \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_n^2} f(\mathbf{x}) \end{bmatrix}$$

The Hessian is the *Jacobian of the Gradient*

Curvature

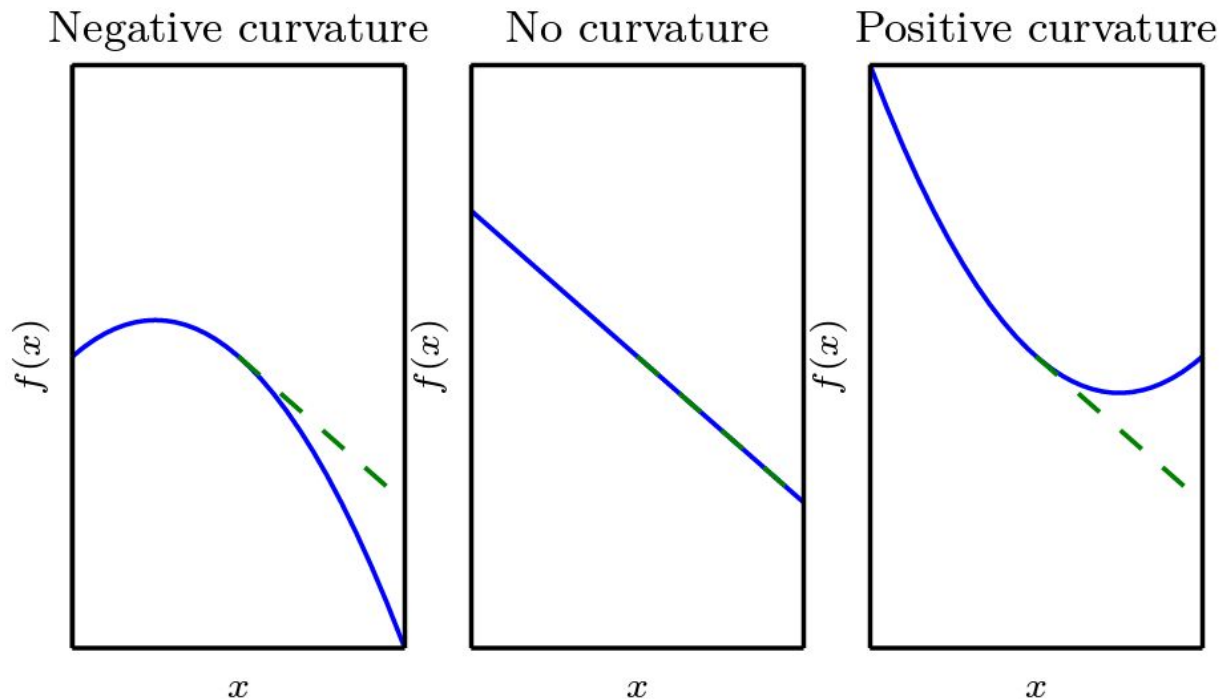


Figure 4.4

Properties of the Hessian

Since differential operators are commutative:

$$\frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}) = \frac{\partial^2}{\partial x_j \partial x_i} f(\mathbf{x})$$

The Hessian is real & symmetric

$$H_{i,j} = H_{j,i}$$

We can get real eigenvalues and real eigenvectors

2nd Order Taylor Series Approximation

In the neighborhood of $\mathbf{x}^{(0)}$, $f(\mathbf{x})$ can be approximated as:

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{g} + \frac{1}{2}(\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}(\mathbf{x} - \mathbf{x}^{(0)})$$

With a learning rate of ϵ , the new point is at $\mathbf{x}^{(0)} - \epsilon \mathbf{g}$, and substituting:

$$f(\mathbf{x}^{(0)} - \epsilon \mathbf{g}) \approx f(\mathbf{x}^{(0)}) - \epsilon \mathbf{g}^\top \mathbf{g} + \frac{1}{2} \epsilon^2 \mathbf{g}^\top \mathbf{H} \mathbf{g}$$



Original point

Expected
Improvement

Curvature
Correction

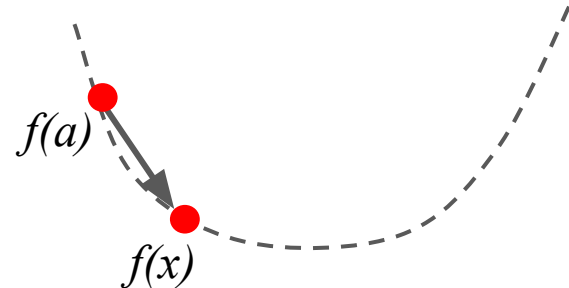
Taylor Series Approximation

Taylor series of a function is an infinite sum of terms that are expressed in terms of the function's derivatives at a single point.

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$$

In practice, rarely use more than 2nd derivative.

$$f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2$$



2nd Derivative Test

Helps to identify **minimum**, **maximum** or **saddle points**.

Minimum: if the Hessian is positive definite (all positive eigenvalues)

Maximum: if the Hessian is negative definite (all negative eigenvalues)

Saddle point: if some positive and negative eigenvalues in Hessian

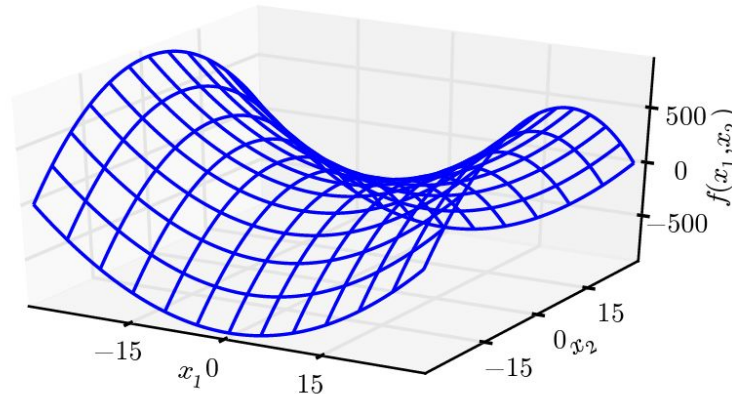


Figure 4.5

Condition Number of a Matrix

For n -by- n matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ the **condition number** is the ratio of its largest to smallest eigenvalue:

$$\max_{i,j} \left| \frac{\lambda_i}{\lambda_j} \right|$$

Large condition numbers lead to numeric instability.

If the condition number of the Hessian is large, gradient descent performs poorly.

Newton's Method

Use the curvature to adjust the step size:

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)}) + \left(\frac{1}{2} \mathbf{x} - \mathbf{x}^{(0)}\right)^\top \mathbf{H}(f)(\mathbf{x}^{(0)}) (\mathbf{x} - \mathbf{x}^{(0)})$$

If f is positive definite and quadratic, Newton's Method jumps to the minimum:

$$\mathbf{x}^* = \mathbf{x}^{(0)} - \mathbf{H}(f)(\mathbf{x}^{(0)})^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}^{(0)})$$

Otherwise, multiple iterations are required to converge of \mathbf{x}^*

Constrained Optimization

Constrained Optimization

Constrained Optimization: rather than optimize over all possible points, apply constraints to some subset of the solution space $\mathcal{S}^n \subset \mathbb{R}^n$

feasible point: $\mathbf{x} \in \mathcal{S}$

- Guide gradient descent to stay inside \mathcal{S}
- Transform into a different unconstrained optimization
- Karush-Kuhn-Tucker (KKT), generalized Lagrangian

Constrained Optimization with KKT

“Generalized Lagrangian”

Constraint S in terms of **equations** $g^{(i)}(\mathbf{x})$ and **inequalities** $h^{(j)}(\mathbf{x})$

$$S = \{\mathbf{x} | \forall i, g^{(i)}(\mathbf{x}) = 0 \text{ and } \forall j, h^{(j)}(\mathbf{x}) \leq 0\}$$

Add in new variables, λ_i , and α_j , for each constraint:

$$L(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\alpha}) = f(\mathbf{x}) + \underbrace{\sum_i \lambda_i g^{(i)}(\mathbf{x}) + \sum_j \alpha_j h^{(j)}(\mathbf{x})}_{\text{Constraints}}$$

Original minimization function

Constraints

Constrained Optimization with KKT

If $f(x)$ has at least one feasible point and $f(x)$ cannot be ∞ we can write the optimization problem:

$$\min_x \max_{\lambda} \max_{\alpha, \alpha \geq 0} L(x, \lambda, \alpha)$$

- Has the same objective function value and set of optimal point x as

$$\min_{x \in \mathbb{S}} f(x)$$

- When constraints are satisfied:

$$\max_{\lambda} \max_{\alpha, \alpha \geq 0} L(x) = f(x)$$

- When constraints are not satisfied:

$$\max_{\lambda} \max_{\alpha, \alpha \geq 0} L(x) = \infty$$

Guarantee: no
infeasible point
can be optimal

Constrained Optimization with KKT

Can reformulate the problem as a maximization:

$$\min_x \max_{\lambda} \max_{\alpha, \alpha \geq 0} -f(x) + \sum_i \lambda_i g^{(i)}(x) + \sum_j \alpha_j h^{(j)}(x)$$

or

$$\max_x \min_{\lambda} \min_{\alpha, \alpha \geq 0} f(x) + \sum_i \lambda_i g^{(i)}(x) - \sum_j \alpha_j h^{(j)}(x)$$

Constrained Optimization with KKT

KKT conditions for a point to be optimal:

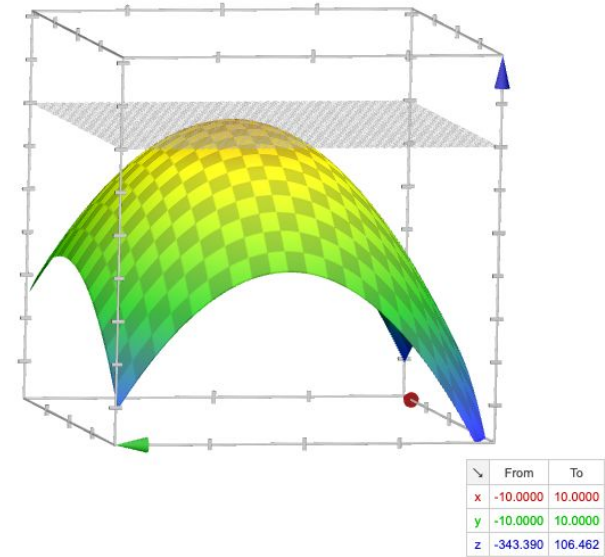
- Gradient is zero at x^*
- All constraints on both x and the KKT multipliers (λ , α) are satisfied
- Inequality constraints exhibit “complementary slackness”:

$$\alpha \odot h(x) = \mathbf{0}$$

Constrained Optimization Example

The Problem

$$\begin{aligned} &\text{maximize} && -(x - 2)^2 - 2(y - 1)^2 \\ &\text{subject to:} && \\ &&& x + 4y \leq 3 \\ &&& -x + y \leq 0 \end{aligned}$$



Constrained Optimization Example

The Lagrangian is:

$$L(x, y, \alpha_1, \alpha_2) = -(x - 2)^2 - 2(y - 1)^2 + \alpha_1(3 - x - 4y) + \alpha_2(x - y)$$

This gives the following KKT Conditions:

$$\frac{\partial L}{\partial x} = -2(x - 2) - \alpha_1 + \alpha_2 = 0$$

$$\frac{\partial L}{\partial y} = -4(y - 1) - 4\alpha_1 - \alpha_2 = 0$$

$$\alpha_1(3 - x - 4y) = 0$$

$$\alpha_2(x - y) = 0$$

$$\alpha_1, \alpha_2 \geq 0$$

Constrained Optimization Example

We have two complementary conditions, i.e., Complementary Slackness ($\alpha \odot h(x) = \mathbf{0}$) with four possible cases:

1. $\alpha_1 = \alpha_2 = 0 \rightarrow x = 2, y = 1$

2. $\alpha_1 = 0, x - y = 0 \rightarrow x = \frac{4}{3}, \alpha_2 = -\frac{4}{3}$

3. $3 - x - 4y = 0, \alpha_2 = 0 \rightarrow x = \frac{5}{3}, y = \frac{1}{3}, \alpha_1 = \frac{2}{3}$

4. $3 - x - 4y = 0, x - y = 0 \rightarrow x = \frac{3}{5}, y = \frac{3}{5}, \alpha_1 = \frac{22}{25}, \alpha_2 = -\frac{48}{25}$

Optimal Solution: $x^* = \frac{5}{3}, y^* = \frac{1}{3}, f(x^*, y^*) = -\frac{4}{9}$