



Lecture 9

Math Foundations Team



BITS Pilani

Pilani | Dubai | Goa | Hyderabad

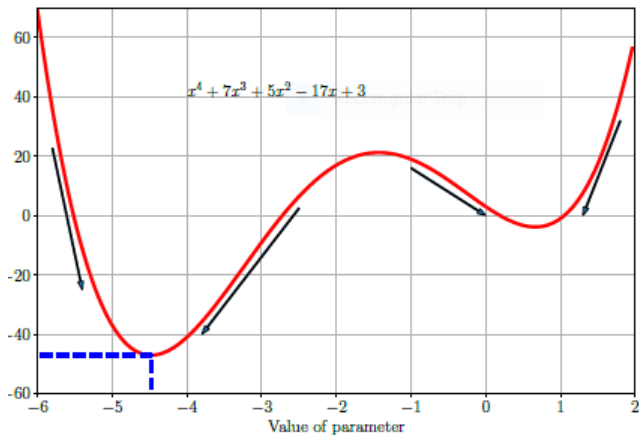


- ▶ We will look at continuous optimization concepts in this lecture.
- ▶ There are two main branches of continuous optimization - constrained and unconstrained optimization.
- ▶ We seek the minimum of an objective function which we assume is differentiable.
- ▶ This is like finding the valleys of the objective function, and since the objective function is differentiable, the gradient tells us the direction to move to get the maximum increase in the objective function



- ▶ We move in the direction of the negative gradient to decrease the objective function.
- ▶ We move until we encounter a point at which the gradient is zero.
- ▶ This constitutes a valley of the objective function, known as a local minimum
- ▶ For unconstrained optimization, this is enough for the big picture.

Example



- ▶ Let $l(x) = x^4 + 7x^3 + 5x^2 - 17x + 3$.
- ▶ The gradient is $\frac{dl(x)}{dx} = 4x^3 + 21x^2 + 10x - 17$.
- ▶ Setting the gradient to zero identifies points corresponding to a local minimum or local maximum - there are three such points since this is a cubic equation.
- ▶ How do we check if we are dealing with a local minimum or local maximum? Look at the second derivative!
- ▶ If the second derivative at the point where the gradient vanishes is negative we are talking about a local maximum, otherwise it is a local minimum.
- ▶ The second derivative is $12x^2 + 42x + 10$



- ▶ For low-order polynomials we can solve the equations analytically and find points at which the gradient is zero. Then we can do the second derivative test and identify whether these points are local minima/local maxima.
- ▶ In general we need to follow the negative gradient.
- ▶ Consider the problem of solving for the minimum of a real-valued function $\min_{\mathbf{x}} f(\mathbf{x})$ where $f : R^d \rightarrow R$ is an objective loss function that captures the cost of using the current set of parameters.
- ▶ We assume our function f is differentiable but that the minimum cannot be found analytically in closed form.



- ▶ The main idea of gradient-descent, a first-order optimization algorithm, is to take a step from the current point of magnitude proportional to the negative gradient of the function at the current point.
- ▶ In mathematical terms if $\mathbf{x}_1 = \mathbf{x}_0 - \gamma((\nabla f)(\mathbf{x}_0))^T$ for a small step-size $\gamma > 0$ then $f(\mathbf{x}_1) \leq f(\mathbf{x}_0)$.
- ▶ The above suggests that in order to find the optimum of f , say at $f(\mathbf{x}_*)$, we can start at some initial point \mathbf{x}_0 and then iterate according to $\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i((\nabla f)(\mathbf{x}_i))^T$
- ▶ For a suitable step-size γ_i , the sequence of points $f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq \dots$ converges to some local minimum.

Example of gradient descent



- ▶ Let $f = \frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x}$ where $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$, $\mathbf{A} = \begin{bmatrix} 2 & 1 \\ 1 & 20 \end{bmatrix}$ and $\mathbf{b} = \begin{bmatrix} -5 \\ -3 \end{bmatrix}$.
- ▶ From a previous lecture on Vector Calculus, we know that $\nabla f = \mathbf{x}^T \mathbf{A} + \mathbf{b}^T$
- ▶ Starting at $\mathbf{x}_0 = \begin{bmatrix} -3 \\ -1 \end{bmatrix}$ and iterating according to the equation given on a previous slide, we will eventually encounter the minimum value.



- ▶ If the step-size is too small, gradient descent will become too slow.
- ▶ If the step-size is too large, gradient descent might overshoot the target, fail to converge, or even diverge.
- ▶ One way to handle the issue of choosing step-sizes is to choose a different step-size at each step. If a given step-size leads to an increase in the value of the cost function, we have been too aggressive, so a better step-size would be half the current step-size. On the other hand, if the value of the function has decreased, the step-size could be larger. In each case we undo the current step and change the step-size.
- ▶ This heuristic guarantees monotonic convergence.



- ▶ Gradient descent may be slow if the curvature of the function is such that the gradient descent steps hop between the walls of the valley of contours and approaches the optimum slowly.
- ▶ If we endow the optimization procedure with memory, we can improve convergence.
- ▶ We use an additional term in the step-update to remember what happened in the previous iteration, so that we can dampen oscillations and speed up convergence. This is a momentum term - the name momentum comes from a comparison to a rolling ball whose direction becomes more and more difficult to change as its velocity increases.
- ▶ In terms of equations we have
$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i((\nabla f)(\mathbf{x}_i))^T + \alpha \Delta \mathbf{x}_i.$$



- ▶ The momentum term helps us deal with a noisy estimate of the current gradient as the memory term remembers the past direction and prevents the current gradient estimate from leading the optimization procedure in a completely wrong direction.
- ▶ Why do we need to work with noisy or approximate gradients? This is because computing the true gradient may be computationally very demanding.
- ▶ We can deal with an approximate gradient and then use the momentum term to smooth out the noise.
- ▶ Stochastic gradient descent is a stochastic way of approximating the gradient of an objective function which is expressed as a sum of differentiable functions.



- ▶ Let us consider a machine learning problem consisting of loss functions incurred at N data points. Let the loss function at the n th data point be $L_n(\theta)$, and let the total loss be $L(\theta) = \sum_{n=1}^{n=N} L_n(\theta)$. Here θ is the parameter vector of interest and the goal of machine learning is to find the parameter vector θ that minimizes the loss function.
- ▶ The standard gradient descent procedure is a batch optimization method in that the update step considers the gradient of the entire loss function $L(\theta)$., i.e
$$\theta_{i+1} = \theta_i - \gamma_i \nabla L(\theta_i)^T = \theta_i - \gamma_i \sum_{n=1}^{n=N} \nabla L_n(\theta_i)^T.$$
- ▶ This step is expensive since we may have many data points in the batch.



- ▶ We can perform a mini-batch gradient descent by taking only a small subset of the data points.
- ▶ In the extreme case we can choose only one L_n to estimate the gradient.
- ▶ Why does taking only a subset of data points work? It works because we get an unbiased estimate of the true gradient.
- ▶ We can show that when the learning rate decreases at a suitable rate and some mild assumptions can be made, stochastic gradient descent almost surely converges to a local minimum.
- ▶ Estimating the gradient out of large mini-batches will lead to lower variance in the estimate, in contrast to estimating the gradient from small batch sizes.



- ▶ Consider the following problem: $\min_{\mathbf{x}} f(\mathbf{x}), f: \mathbb{R}^D \rightarrow \mathbb{R}$, subject to additional constraints - so we are looking at a minimization problem except that the set of all \mathbf{x} over which minimization is performed is not all of \mathbb{R}^D .
- ▶ The constrained problem becomes $\min_{\mathbf{x}} f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0 \forall i=1, 2, \dots, m$.
- ▶ Since we have a method of finding a solution to the unconstrained optimization problem, one way to proceed now is to convert the given constrained optimization problem into an unconstrained one.
- ▶ We construct $J(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \mathbf{1}(g_i(\mathbf{x}))$, where $\mathbf{1}(z) = 0$ for $z \leq 0$ and $\mathbf{1}(z) = \infty$ for $z > 0$.



- ▶ The formulation of $J(\mathbf{x})$ in the previous slide ensures that its value is infinity if any one of the constraints $g_i(\mathbf{x})$ is not satisfied. This ensures that the optimal solution to the unconstrained problem is the same as the constrained problem.
- ▶ The step-function is also difficult to optimize, and our solution is to replace the step-function by a linear function using Lagrange multipliers.
- ▶ We create the Lagrangian of the given constrained optimization problem as follows:
$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) = f(\mathbf{x}) + \boldsymbol{\lambda}^T \mathbf{g}(\mathbf{x}), \text{ where } \lambda_i \geq 0 \text{ for all } i.$$



- ▶ The primal problem is $\min f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0, 1 \leq i \leq m$. Optimization is performed over the primal variables \mathbf{x} .
- ▶ The associated Lagrangian dual problem is $\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \mathfrak{D}(\boldsymbol{\lambda})$ subject to $\boldsymbol{\lambda} \geq 0$ where $\boldsymbol{\lambda}$ are dual variables.
- ▶ $\mathfrak{D}(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$.
- ▶ The following minimax inequality holds over two arguments \mathbf{x}, \mathbf{y} : $\max_{\mathbf{y}} \min_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{y}) \leq \min_{\mathbf{x}} \max_{\mathbf{y}} \phi(\mathbf{x}, \mathbf{y})$.



- ▶ Why is this inequality true?
- ▶ Assume that \mathbf{x}, \mathbf{y} : $\max_{\mathbf{y}} \min_{\mathbf{x}} \phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}_A, \mathbf{y}_A)$ and $\min_{\mathbf{x}} \max_{\mathbf{y}} \phi(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x}_B, \mathbf{y}_B)$.
- ▶ Fixing \mathbf{y} at \mathbf{y}_A we see that the inner operation on the left hand side of the minimax inequality is a min operation over \mathbf{x} and returns \mathbf{x}_A . Thus we have $\phi(\mathbf{x}_A, \mathbf{y}_A) \leq \phi(\mathbf{x}_B, \mathbf{y}_A)$.
- ▶ Fixing \mathbf{x} at \mathbf{x}_B we see that the inner operation on the right hand side of the minimax inequality is a max operation over \mathbf{y} and returns \mathbf{y}_B . Thus we have $\phi(\mathbf{x}_B, \mathbf{y}_B) \geq \phi(\mathbf{x}_B, \mathbf{y}_A)$.
- ▶ From the above we conclude that $\phi(\mathbf{x}_B, \mathbf{y}_B) \geq \phi(\mathbf{x}_A, \mathbf{y}_A)$.



- ▶ The difference between $J(\mathbf{x})$ and the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is that the indicator function is relaxed to a linear function.
- ▶ When $\boldsymbol{\lambda} \geq 0$, the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is a lower bound on $J(\mathbf{x})$.
- ▶ The maximum of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ with respect to $\boldsymbol{\lambda}$ is $J(\mathbf{x})$ - if the point \mathbf{x} satisfies all the constraints $g_i(\mathbf{x}) \leq 0$, then the maximum of the Lagrangian is obtained at $\boldsymbol{\lambda} = 0$ and it is equal to $J(\mathbf{x})$. If one or more constraints is violated such that $g_i(\mathbf{x}) > 0$, then the associated Lagrangian coefficient λ_i can be taken to be ∞ so as to equal $J(\mathbf{x})$.



- ▶ From the previous slide, we have $J(\mathbf{x}) = \max_{\lambda \geq 0} \mathcal{L}(\mathbf{x}, \lambda)$.
- ▶ Our original constrained optimization problem boiled down to minimizing $J(\mathbf{x})$, in other words we are looking at $\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\lambda \geq 0} \mathcal{L}(\mathbf{x}, \lambda)$
- ▶ Using the minimax inequality we see that $\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\lambda \geq 0} \mathcal{L}(\mathbf{x}, \lambda) \geq \max_{\lambda \geq 0} \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \lambda)$.
- ▶ This is known as weak duality. The inner part of the right hand side of the inequality is $\mathcal{D}(\lambda)$, and the inequality above is the reason for setting up the associated Lagrangian dual problem for the original constrained optimization problem.



- ▶ In contrast to the original formulation $\mathfrak{D}(\lambda) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \lambda)$ is an unconstrained optimization problem for a given value of λ .
- ▶ We observe that $\mathfrak{D}(\lambda) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \lambda)$ is a point-wise minimum of affine functions and hence $\mathfrak{D}(\lambda)$ is concave even though $f()$ and $g()$ may be nonconvex.
- ▶ We have obtained a Lagrangian formulation for a constrained optimization problem where the constraints are inequalities. What happens when some constraints are equalities?



- ▶ Suppose the problem is $\min_{\mathbf{x}} f(\mathbf{x})$ subject to $g_i(\mathbf{x}) \leq 0$ for all $1 \leq i \leq m$ and $h_j(\mathbf{x}) = 0$ for $1 \leq j \leq n$.
- ▶ We model the equality constraint $h_j(\mathbf{x}) = 0$ with two inequality constraints $h_j(\mathbf{x}) \geq 0$ and $h_j(\mathbf{x}) \leq 0$.
- ▶ The resulting Lagrange multipliers are then unconstrained.
- ▶ The Lagrange multipliers for the original inequality constraints are non-negative while those corresponding to the equality constraints are unconstrained.



- ▶ We are interested in a class of optimization problems where we can guarantee global optimality.
- ▶ When $f()$, the objective function, is a convex function and $g()$ and $h()$ are convex functions, we have a convex optimization problem.
- ▶ In this setting we have strong duality - the optimal solution of the primal problem is equal to the optimal solution of the dual problem.
- ▶ What is a convex function?



- ▶ First we need to know what is a convex set. A set C is a convex set if for any $x, y \in C$, $\theta x + (1 - \theta)y \in C$.
- ▶ For any two points lying in the convex set, a line joining them lies entirely in the convex set.
- ▶ Let a function $f : \mathbb{R}^d \rightarrow R$ be a function whose domain is a convex set C .
- ▶ The function is a convex function if for any $\mathbf{x}, \mathbf{y} \in C$,
$$f(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta) f(\mathbf{y})$$
- ▶ Another way of looking at a convex function is to use the gradient: for any two points \mathbf{x} and \mathbf{y} , we have
$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla_{\mathbf{x}} f(\mathbf{x})(\mathbf{y} - \mathbf{x}).$$

- ▶ The negative entropy, a useful function in Machine Learning, is convex: $f(x) = x \log_2 x$ for $x > 0$.
- ▶ First let us check if $f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y)$. Take $x = 2$, $y = 4$, and $\theta = 0.5$ to get
 $f(0.5 * 2 + 0.5 * 4) = f(3) = 3 \log_2 3 \approx 4.75$. Then
 $\theta f(2) + (1 - \theta)f(4) = 0.5 * 2 \log_2 2 + 0.5 * 4 \log_2 4 = \log_2 32 = 5$.
Therefore the convexity criterion is satisfied for these two points.
- ▶ Let us now use the gradient criterion. We have
 $\nabla f(x) = \log_2 x + x \frac{1}{x \log_e 2}$. Calculating $f(2) + \nabla f(2) * (4 - 2)$ gives $2 \log_2 2 + (\log_2 2 + \frac{1}{\log_e 2} * 2) \approx 6.9$. We see that
 $f(4) = 4 \log_2 4 = 8$ which shows that the gradient criterion is also satisfied.



- ▶ Let us look at a convex optimization problem where the objective function and constraints are all linear.
- ▶ Such a convex optimization problem is called a linear programming problem.
- ▶ We can express a linear programming problem as $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ subject to $\mathbf{Ax} \leq \mathbf{b}$ where $\mathbf{A} \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^{m \times 1}$.
- ▶ The Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is given by $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^T \mathbf{x} + \boldsymbol{\lambda}^T (\mathbf{Ax} - \mathbf{b})$ where $\boldsymbol{\lambda} \in \mathbb{R}^m$ is the vector of non-negative Lagrangian multipliers.
- ▶ We can rewrite the Lagrangian as $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{x} - \boldsymbol{\lambda}^T \mathbf{b}$.



- ▶ Taking the derivative of the Lagrangian with respect to \mathbf{x} and setting it to zero we get $\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} = 0$.
- ▶ Since $\mathcal{D}(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$, plugging in the above equation gives $\mathcal{D}(\boldsymbol{\lambda}) = -\boldsymbol{\lambda}^T \mathbf{b}$.
- ▶ We would like to maximize $\mathcal{D}(\boldsymbol{\lambda})$, subject to the constraint $\boldsymbol{\lambda} \geq 0$.
- ▶ Thus we end up with the following problem:

$$\begin{aligned} \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} & -\boldsymbol{\lambda}^T \mathbf{b} \\ \text{subject to} & \mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda} = 0 \\ & \boldsymbol{\lambda} \geq 0 \end{aligned}$$



- ▶ We can solve the original primal linear program or the dual one - the optimum in each case is the same.
- ▶ The primal linear program is in d variables but the dual is in m variables, where m is the number of constraints in the original primal program.
- ▶ We choose to solve the primal or dual based on which of m or d is smaller.



- ▶ We now consider the case of a quadratic objective function subject to affine constraints:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x} \text{ subject to}$$
$$\mathbf{A} \mathbf{x} \leq \mathbf{b}$$

- ▶ Here $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$, $\mathbf{c} \in \mathbb{R}^d$



- ▶ The Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is given by $\frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} + \mathbf{c}^T \mathbf{x} + \boldsymbol{\lambda}^T (\mathbf{A}\mathbf{x} - \mathbf{b})$.
- ▶ Rearranging the above we have $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q}\mathbf{x} + (\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{x} - \boldsymbol{\lambda}^T \mathbf{b}$
- ▶ Taking the derivative of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ and setting it equal to zero gives $\mathbf{Q}\mathbf{x} + (\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda}) = \mathbf{0}$.
- ▶ If we take \mathbf{Q} to be invertible, we have $\mathbf{x} = \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda})$.
- ▶ Plugging this value of \mathbf{x} into $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ gives us $\mathcal{D}(\boldsymbol{\lambda}) = -\frac{1}{2}(\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda}) - \boldsymbol{\lambda}^T \mathbf{b}$.
- ▶ This gives us the dual optimization problem:
 $\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} -\frac{1}{2}(\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda})^T \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^T \boldsymbol{\lambda}) - \boldsymbol{\lambda}^T \mathbf{b}$ subject to $\boldsymbol{\lambda} \geq \mathbf{0}$.