

# CEE 616: Probabilistic Machine Learning Foundations: Optimization

**Jimi Oke**

UMassAmherst  

---

College of Engineering

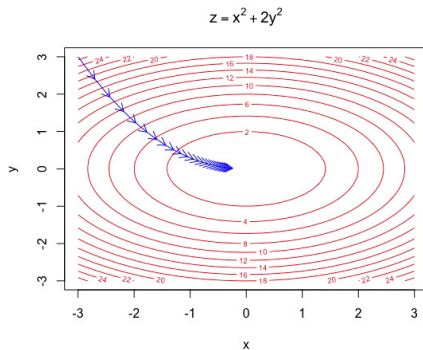
Thu, Sep 18, 2025

# Outline

- ① Introduction
- ② First-order methods
- ③ Second-order methods
- ④ Application: MLE
- ⑤ Constrained optimization
- ⑥ Outlook

# Optimization

Optimization is the body of mathematics that deals with the theory and algorithms for characterizing the maximum/minimum values of functions.



We will consider two widely-used approaches in machine learning:

- First-order methods (e.g. gradient descent)
- Second-order methods (e.g. Newton's method)

# Definitions

- Minimum of convex function

$$\theta^* = \arg \min_{\theta} \mathcal{L}(\theta) : \frac{d\mathcal{L}(\theta^*)}{d\theta} = 0 \quad (1)$$

- Derivative of a function:

$$\mathcal{L}'(\theta) = \frac{d\mathcal{L}(\theta)}{d\theta} \quad (2)$$

- Gradient of a function (vector of partial derivatives):

$$\nabla_{\theta} \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial \mathcal{L}(\theta)}{\partial \theta_1} \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial \mathcal{L}(\theta)}{\partial \theta_n} \end{bmatrix} \quad (3)$$

# Further definitions

- Second-order derivative:

$$\mathcal{L}''(\theta) = \frac{d^2 \mathcal{L}(\theta)}{d\theta^2} \quad (4)$$

- Hessian (matrix of partial second derivatives):

$$\mathbf{H}_{\theta}(\mathcal{L}(\theta)) = \nabla_{\theta}^2 \mathcal{L}(\theta) = \begin{bmatrix} \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_1^2} & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_1 \partial \theta_2} & \cdots & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_1 \partial \theta_n} \\ \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_2 \partial \theta_1} & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_2^2} & \cdots & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_2 \partial \theta_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_n \partial \theta_1} & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_n \partial \theta_2} & \cdots & \frac{\partial^2 \mathcal{L}(\theta)}{\partial \theta_n^2} \end{bmatrix} \quad (5)$$

# Optimality conditions

For continuous, twice-differentiable functions:

- If  $\theta^*$  is a local minimum, then  $\mathbf{g}^* = \mathbf{0}$  and  $\mathbf{H}^*$  must be psd (necessary condition)
- If  $\mathbf{g}^* = \mathbf{0}$  and  $\mathbf{H}^*$  is pd, then  $\theta^*$  is a local optimum (sufficient condition)

## General approach

- Begin with an initial value  $\theta_0$
- At each iteration  $t$ , update  $\theta_{t+1}$
- Terminate when  $\mathcal{L}(\theta_{t+1}) - \mathcal{L}(\theta_t) = \epsilon$

# Convex optimization

If the objective is convex and defined over a convex set, then every local minimum is a global minimum. pe

- Convex set: For any  $\mathbf{x}, \mathbf{x}' \in \mathcal{S}$ :

$$\lambda \mathbf{x} + (1 - \lambda) \mathbf{x}' \in \mathcal{S}, \quad \forall \lambda \in [0, 1] \quad (6)$$

- Convex function: For any  $\mathbf{x}, \mathbf{y} \in \mathcal{S}$  and for any  $0 \leq \lambda \leq 1$ :

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \quad (7)$$

## Quadratic form

Given  $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$ :

- $f$  is convex if  $\mathbf{A}$  is psd
- $f$  is strictly convex if  $\mathbf{A}$  is pd

# Subgradients

Given a convex function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\mathbf{g} \in \mathbb{R}^n$  is a subgradient of  $f$  at  $\mathbf{x} \in \text{dom}(f)$  if:

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^\top (\mathbf{z} - \mathbf{x}) \quad \forall \mathbf{z} \in \text{dom}(f) \quad (8)$$

- The set of subgradients is called the **subdifferential**:  $\partial f(\mathbf{x})$
- $f$  is subdifferentiable at  $\mathbf{x}$  if it has at least one subgradient at that point

## Subdifferential of ReLU

The rectified linear unit function (ReLU) is given by:

$$\text{ReLU}(z) = \max(0, z), \quad z \in \mathbb{R} \quad (9)$$

Its subdifferential is:

$$\partial \text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ [0, 1] & \text{if } z = 0 \\ 1 & \text{if } z > 0 \end{cases} \quad (10)$$



# First-order methods

These methods only require first-order derivatives in order to compute an update:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \rho_t \mathbf{d}_t \quad (11)$$

where:

- $\rho_t$ : step size/learning rate
- $\mathbf{d}_t$ : descent direction (e.g. negative gradient)
- $t$ : index of iteration

The classic first-order method is **gradient descent**

# Descent direction

A vector  $\mathbf{d}$  is considered a descent direction if a nonzero step size  $\rho$  moved in its direction yields a decrease in the objective:

$$\mathcal{L}(\boldsymbol{\theta} + \rho \mathbf{d}) < \mathcal{L}(\boldsymbol{\theta}) \quad \forall 0 < \rho < \rho_{\max} \quad (12)$$

## Steepest descent

The direction of steepest descent is given by the negative gradient:

$$-\mathbf{g}_t = -\nabla \mathcal{L}(\boldsymbol{\theta}_t) \quad (13)$$

# Step size

The choice of step size impacts the convergence of an optimization routine. It determines how far along in the descent direction the variable is updated at each step.

- **Constant step size:**  $\rho$  is fixed throughout
- **Adaptive step size:**  $\rho_t$  is modified at each iteration to satisfy certain conditions (e.g. line search methods)
- Sequence  $\{\rho_t\}$  in an optimization algorithm is known as the learning rate schedule

# Example: Gradient descent

Given the function

$$f(x) = (x - 1)^4 - 3x + 4$$

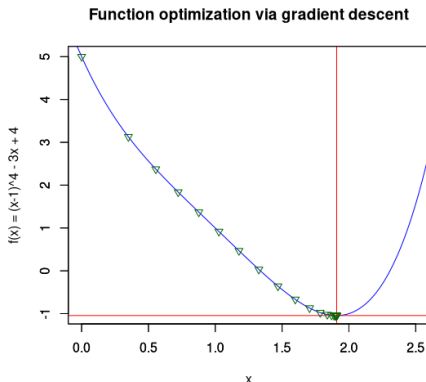
find the optimal point  $(x^*, f(x^*))$ .

First, we find the **gradient** and define the **update** step:

$$\begin{aligned} f'(x) &= 4(x - 1)^3 - 3 \\ x_{k+1} &= x_k - \lambda [4(x_k - 1)^3 - 3] \end{aligned}$$

We choose  $\lambda = 0.05$  and a random starting point between 0 and 1.

# Example: Gradient descent (cont.)



Using the `gradient-descent.ipynb` notebook, we find the optimal point as **(1.908, -1.044)**. How does the learning rate impact the solution process?

# Exact line search

Finds optimal  $\rho_t$  by:

$$\rho_t^* = \arg \min_{\rho > 0} \phi_t(\rho) = \arg \min_{\rho > 0} \mathcal{L}(\boldsymbol{\theta}_t + \rho \mathbf{d}_t) \quad (14)$$

If  $\mathcal{L}(\boldsymbol{\theta})$  is quadratic, then we can write:

$$\phi_t(\rho) = \mathcal{L}(\boldsymbol{\theta}_t + \rho \mathbf{d}_t) = \frac{1}{2}(\boldsymbol{\theta} + \rho \mathbf{d})^\top \mathbf{A}(\boldsymbol{\theta} + \rho \mathbf{d}) + \mathbf{b}^\top(\boldsymbol{\theta} + \rho \mathbf{d}) + c \quad (15)$$

Solving for  $\phi'_t(\rho) = 0$  gives the optimal update as:

$$\rho_t^* = -\frac{\mathbf{d}_t^\top (\mathbf{A}\boldsymbol{\theta}_t + \mathbf{b})}{\mathbf{d}_t^\top \mathbf{A} \mathbf{d}_t} \quad (16)$$

This approach introduces additional computational expense

# Armijo backtracking

Finds a suitable  $\rho_t^*$  by iteratively shrinking  $\rho_t$  by  $\beta \in (0, 1)$ , i.e.

$$\rho_t^{k+1} = \beta \rho_t^k \quad (17)$$

until the **Armijo-Goldstein** condition is satisfied:

$$\mathcal{L}(\theta_t + \rho \mathbf{d}_t) \leq \mathcal{L}(\theta_t) + c \rho \mathbf{d}_t^\top \nabla \mathcal{L}(\theta_t) \quad (18)$$

where  $c \in [0, 1]$  is typically chosen as  $10^{-4}$

# Momentum

Momentum methods are used to improve convergence. The momentum update is given by:

$$\mathbf{m}_{t+1} = \beta \mathbf{m}_t + \mathbf{g}_t \quad (19)$$

where  $\beta \in (0, 1)$ .

The momentum can also be expressed as an exponentially weighted average of past gradients:

$$\mathbf{m}_{t+1} = \sum_{\tau=0}^t \beta^\tau \mathbf{g}_{t-\tau} \quad (20)$$

The update is given by:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \rho_t \mathbf{m}_{t+1} \quad (21)$$



# Nesterov momentum

In the Nesterov method, the momentum update computes the gradient at the new location, which can speed up convergence:

$$\mathbf{m}_{t+1} = \beta \mathbf{m}_t - \rho_t \nabla \mathcal{L}(\boldsymbol{\theta}_t + \beta \mathbf{m}_t) \quad (22)$$

The update is given by:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \mathbf{m}_{t+1} \quad (23)$$

This method is also called **Nesterov accelerated gradient**

# Stochastic gradient descent (SGD)

In stochastic optimization, goal is to minimize:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z})}[\mathcal{L}(\boldsymbol{\theta}, \mathbf{z})] \quad (24)$$

where  $\mathbf{z}_t \sim q$  could be a training sample drawn from a set. Thus, the **stochastic gradient descent** update is given by:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \rho_t \nabla \mathcal{L}(\boldsymbol{\theta}_t, \mathbf{z}_t) = \boldsymbol{\theta}_t - \rho_t \mathbf{g}_t \quad (25)$$

- Finite sum:  $\mathcal{L}(\boldsymbol{\theta}_t) = \frac{1}{N} \sum_{n=1}^N \mathcal{L}_n(\boldsymbol{\theta}_t)$

# Newton's method

Second-order methods incorporate curvature via the second derivative to speed up convergence.

The basic method of this form is Newton's method:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \rho_t \mathbf{H}^{-1} \mathbf{g}_t \quad (26)$$

where:

$$\mathbf{H}_t := \nabla^2 \mathcal{L}(\boldsymbol{\theta}_t) = \mathbf{H}(\boldsymbol{\theta}_t) \quad (27)$$

# Derivation of Newton's method

Taylor-expand  $\mathcal{L}(\boldsymbol{\theta})$  around  $\boldsymbol{\theta}_t$ :

$$\mathcal{L}_{\text{quad}} = \mathcal{L}(\boldsymbol{\theta}_t) + \mathbf{g}_t^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_t) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_t)^\top \mathbf{H}_t (\boldsymbol{\theta} - \boldsymbol{\theta}_t) \quad (28)$$

Taking the derivative of  $\mathcal{L}_{\text{quad}}$  and setting it to zero to find the minimum, we obtain:

$$\mathbf{g}_t + \mathbf{H}_t (\boldsymbol{\theta} - \boldsymbol{\theta}_t) = \mathbf{0} \quad (29)$$

$$\boldsymbol{\theta} - \boldsymbol{\theta}_t = -\mathbf{H}_t^{-1} \mathbf{g}_t \quad (30)$$

$$\boldsymbol{\theta} = \boldsymbol{\theta}_t - \mathbf{H}_t^{-1} \mathbf{g}_t \quad (31)$$

Thus, we set the descent direction  $\mathbf{d}_t$  as  $-\mathbf{H}_t^{-1} \mathbf{g}_t$

# BFGS

Broyden-Fletcher-Goldfarb-Shanno method. Approximates  $\mathbf{H}_t \approx \mathbf{B}_t$ :

$$\mathbf{B}_{t+1} = \mathbf{B}_t + \frac{\mathbf{y}_t \mathbf{y}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} - \frac{(\mathbf{B}_t \mathbf{s}_t)(\mathbf{B}_t \mathbf{s}_t)^\top}{\mathbf{s}_t^\top \mathbf{B}_t \mathbf{s}_t} \quad (32)$$

$$\mathbf{s}_t = \boldsymbol{\theta}_t - \boldsymbol{\theta}_{t-1} \quad (33)$$

$$\mathbf{y}_t = \mathbf{g}_t - \mathbf{g}_{t-1} \quad (34)$$

- $\mathbf{B}_0$  is typically initialized as  $\mathbf{I}$  (positive definite)
- To ensure  $\mathbf{B}_{t+1}$  remains pd,  $\rho$  must satisfy the Wolfe conditions:

$$\mathcal{L}(\boldsymbol{\theta}_t + \rho \mathbf{d}_t) \leq \mathcal{L}(\boldsymbol{\theta}_t) + c_1 \rho \mathbf{d}_t^\top \nabla \mathcal{L}(\boldsymbol{\theta}_t) \quad (35)$$

$$\mathcal{L}(\boldsymbol{\theta}_t + \rho \mathbf{d}_t) \geq c_2 \rho \mathbf{d}_t^\top \nabla \mathcal{L}(\boldsymbol{\theta}_t) \quad (36)$$

where  $0 < c_1 < c_2 < 1$

# Limited memory BFGS

In practice, BFGS directly provides the inverse Hessian approximation:

$$\mathbf{C}_{t+1} = \left( \mathbf{I} - \frac{\mathbf{s}_t \mathbf{y}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} \right) \mathbf{C}_t \left( \mathbf{I} - \frac{\mathbf{y}_t \mathbf{s}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} \right) + \frac{\mathbf{s}_t \mathbf{s}_t^\top}{\mathbf{y}_t^\top \mathbf{s}_t} \quad (37)$$

where  $\mathbf{C}_t \approx \mathbf{H}_t^{-1}$ .

- The update can be reduced to a recurrence relation that depends explicitly on  $\mathbf{C}_0$  and the history of  $\mathbf{s}_t$  and  $\mathbf{y}_t$
- For computational efficiency, the  $M$  most recent  $\mathbf{s}_t, \mathbf{y}_t$  may be used instead
- $M$  is typically  $\in [5, 20]$
- This approach is termed L-BFGS (limited memory BFGS)

# Maximum likelihood estimation: Key definitions

Given a set of  $n$  independent observations  $x_1, x_2, \dots, x_n$  from a random sample, the **likelihood function** is:

$$L(x_1, x_2, \dots, x_n; \theta) = f(x_1; \theta)f(x_2; \theta) \cdots f(x_n; \theta) \quad (38)$$

The **maximum likelihood estimator (MLE)**,  $\hat{\theta}$ , is the value of  $\theta$  that maximizes the [log-]likelihood function:

$$\frac{\partial \ln L(x_1, x_2, \dots, x_n; \theta)}{\partial \theta} = 0 \quad (39)$$

For multiple parameters, the likelihood function is:

$$L(x_1, \dots, x_n; \theta_1, \dots, \theta_m) = \prod_{i=1}^n f(x_i; \theta_i; \theta_1, \dots, \theta_m) \quad (40)$$

And the MLE's would be found by simultaneously solving the partial derivatives set to 0 for each parameter.

# Maximum likelihood estimation (MLE) in logistic regression

The log-likelihood function for the binomial logistic regression case is

$$\ell(\beta) = \sum_i [y_i (\beta_0 + \beta_1 x_i) - \log (1 + e^{\beta_0 + \beta_1 x_i})] \quad (41)$$

The optimal  $\hat{\beta}$  which maximizes  $\ell(\beta)$  is the **maximum likelihood estimate**.

Also recall the derivative of  $\ell$ :

$$\nabla_{\beta} \ell = \begin{pmatrix} \frac{\partial \ell}{\partial \beta_0} \\ \frac{\partial \ell}{\partial \beta_1} \end{pmatrix} = \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (42)$$

We can use either Newton-Raphson or gradient *ascent* to *maximize*  $\ell$ .



# Gradient ascent for MLE in logistic regression

This approach only requires the first derivative:

$$\beta_{k+1} = \beta_k + \lambda \nabla \ell(\beta_k) \quad (43)$$

Thus, to find  $\hat{\beta}$  we iterate using:

$$\begin{pmatrix} \beta_{0,k} \\ \beta_{1,k} \end{pmatrix} = \begin{pmatrix} \beta_{0,k} \\ \beta_{1,k} \end{pmatrix} + \lambda \begin{pmatrix} \sum_i [y_i - p(x_i)] \\ \sum_i [x_i (y_i - p(x_i))] \end{pmatrix} \quad (44)$$

Because the log-likelihood is *concave*, and thus a *maximization* problem, we *ascend* the function and thus *add* the scaled derivative.

- As we can see, the gradient ascent method does not require a second derivative
- However, it may require more iterations to converge than Newton-Raphson

# Newton-Raphson approach for MLE in logistic regression

The optimal point  $\hat{\beta}$  is given by the root of the equation  $\nabla_{\beta}\ell = 0$ .

Applying Newton-Raphson, the update step is:

$$\beta_{k+1} = \beta_k - \mathbf{H}_{\beta_k}^{-1}(\ell) \nabla_{\beta_k} \ell(\beta_k) \quad (45)$$

The operator  $\mathbf{H}^{-1}$  represents the inverse **Hessian** (second derivative) matrix of  $\ell$  with respect to  $\beta$ :

$$\mathbf{H}_{\beta_k}(\beta_k) = \nabla_{\beta_k}^2 \ell = \begin{pmatrix} \frac{\partial^2 \ell(\beta)}{\partial \beta_0^2} & \frac{\partial^2 \ell(\beta)}{\partial \beta_0 \beta_1} \\ \frac{\partial^2 \ell(\beta)}{\partial \beta_1 \beta_0} & \frac{\partial^2 \ell(\beta)}{\partial \beta_1^2} \end{pmatrix} \quad (46)$$

Note that the (45) is just the matrix representation of the 1-D case:

$$\beta_{k+1} = \beta_k - \frac{\ell'(\beta_k)}{\ell''(\beta_k)} \quad (47)$$

# Newton-Raphson approach for MLE (cont.)

We can work out each component of the second derivative:

$$\frac{\partial^2 \ell(\beta)}{\partial \beta_0^2} = - \sum_i p(x_i)(1 - p(x_i)) \quad (48)$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta_0 \beta_1} = - \sum_i x_i p(x_i)(1 - p(x_i)) \quad (49)$$

$$\frac{\partial^2 \ell(\beta)}{\partial \beta_1^2} = - \sum_i x_i^2 p(x_i)(1 - p(x_i)) \quad (50)$$

The complete update can then be shown as:

$$\begin{pmatrix} \beta_{0,k+1} \\ \beta_{1,k+1} \end{pmatrix} = \begin{pmatrix} \beta_{0,k} \\ \beta_{1,k} \end{pmatrix} - \left[ \begin{pmatrix} \frac{\partial^2 \ell(\beta)}{\partial \beta_0^2} & \frac{\partial^2 \ell(\beta)}{\partial \beta_0 \partial \beta_1} \\ \frac{\partial^2 \ell(\beta)}{\partial \beta_1 \partial \beta_0} & \frac{\partial^2 \ell(\beta)}{\partial \beta_1^2} \end{pmatrix}^{-1} \begin{pmatrix} \frac{\partial \ell}{\partial \beta_0} \\ \frac{\partial \ell}{\partial \beta_1} \end{pmatrix} \right]_{\beta_k} \quad (51)$$

Alternatively:

$$\beta_{k+1} = \beta_k - \left[ \left( \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \ell(\beta)}{\partial \beta} \right]_{\beta_k} \quad (52)$$

# Constrained optimization

Deals with problems where we seek to minimize an objective:

$$\min_{\theta \in \mathcal{C}} \mathcal{L}(\theta) \quad (53)$$

subject to:

$$h_i(\theta) = 0, \quad i \in \mathcal{E} \quad (54)$$

$$g_j(\theta) \leq 0, \quad j \in \mathcal{I} \quad (55)$$

where:

- $\mathcal{C}$ : constraint/feasible set
- $\mathcal{E}$ : set of equality constraints
- $\mathcal{I}$ : set of inequality constraints

# Lagrange multipliers

Given an optimization problem with  $m$  equality constraints  $\mathbf{h}(\boldsymbol{\theta})$ , we can write the Lagrangian as:

$$L(\boldsymbol{\theta}, \lambda) := \mathcal{L}(\boldsymbol{\theta}) + \sum_{j=1}^m \lambda_j h_j(\boldsymbol{\theta}) \quad (56)$$

where  $\lambda_j$  are the Lagrange multipliers.

Thus, to find  $*$ , we solve:

$$\nabla_{\boldsymbol{\theta}, \lambda} L(\boldsymbol{\theta}, \lambda) = \mathbf{0} \quad (57)$$

- If  $\boldsymbol{\theta} \in \mathbb{R}^D$ , there are  $D + m$  equations with an equal number of unknowns

# Karush-Kuhn-Tucker (KKT) conditions

Given a problem with equality constraints  $\mathbf{h}(\boldsymbol{\theta}) = \mathbf{0}$  and inequality constraints  $\mathbf{g}(\boldsymbol{\theta}) \leq \mathbf{0}$ , the generalized Lagrangian is given by:

$$\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\mu}, \boldsymbol{\lambda}) = \mathcal{L}(\boldsymbol{\theta}) + \sum_i \mu_i g_i(\boldsymbol{\theta}) + \sum_j \lambda_j h_j(\boldsymbol{\theta}) \quad (58)$$

The optimization problem then becomes:

$$\min_{\boldsymbol{\theta}} \max_{\boldsymbol{\mu} \geq \mathbf{0}, \boldsymbol{\lambda}} L(\boldsymbol{\theta}, \boldsymbol{\mu}, \boldsymbol{\lambda}) \quad (59)$$

When  $\mathcal{L}$  and  $g$  are convex, then the KKT conditions are necessary and sufficient for global optimality:

- Feasibility: constraints satisfied
- Stationarity of solution:  $\nabla_{\boldsymbol{\theta}, \boldsymbol{\mu}, \boldsymbol{\lambda}} L = \mathbf{0}$
- Dual feasibility:  $\boldsymbol{\mu} \geq \mathbf{0}$
- Complementary slackness:  $\boldsymbol{\mu} \odot \mathbf{g} = \mathbf{0}$

# Further topics

Some to be covered in Advanced Probabilistic ML:

- Linear programming (simplex algorithm)
- Quadratic programming
- Proximal gradient method
- Bound optimization (majorize-minimize algorithms)
- Expectation maximization (EM); for MLE/MAP estimation

# Reading assignments

- **PMLI 8**
- **PMLCE 5**