

Optimization Techniques for Machine Learning

AMLZC326 · #16 Recap

Anshid Aboobacker

COURSE AT A GLANCE

M 1: Analytic Geometry (L1–2)

Vector spaces, norms, inner products, orthogonality, Gram-Schmidt, spectral theorem

M 2: Linear Optimization (L3–4)

LPP, simplex method, duality, sensitivity analysis

M 3: Vector Calculus (L5–6)

Gradients, Jacobian, Hessian, Taylor expansion, critical points

M 4: Continuous Optimization (L7–10)

Gradient descent, SGD, AdaGrad, RMSProp, Adam, Lagrangians, KKT

M 5: ML Models (L11–15)

PCA, SVD, Gaussian mixture models, EM algorithm, SVM, ARIMA, Kalman filter

LEARNING OBJECTIVES

By the end of this recap lecture you should be able to:

- Recall the key definitions, theorems, and algorithms from each of the 5 modules
- Connect the mathematical tools (analytic geometry, calculus, duality) to the ML models they enable
- Identify which optimisation technique applies to which ML problem
- See the through-line: how every module builds on the previous one toward EM, SVM, and Kalman

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization
- 3 Vector Calculus
- 4 Continuous Optimization & Duality
- 5 Principal Component Analysis
- 6 Machine Learning Models

VECTOR SPACES & NORMS

Vector Space $(V, +, \cdot)$

A set V closed under vector addition and scalar multiplication, satisfying 8 axioms (commutativity, associativity, distributivity, existence of zero and additive inverse, etc.)

Norm $\|\cdot\| : V \rightarrow \mathbb{R}$

- **Positivity:** $\|\mathbf{x}\| \geq 0$; $\|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$
- **Homogeneity:** $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$
- **Triangle inequality:** $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

Euclidean Norm

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

INNER PRODUCTS & ORTHOGONALITY

Inner Product $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$

Symmetric, bilinear, positive definite. Standard: $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$

Cauchy-Schwarz Inequality

$$|\langle \mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{u}\| \|\mathbf{v}\|$$

Angle Between Vectors

$$\cos \theta = \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

$$\mathbf{u} \perp \mathbf{v} \Leftrightarrow \langle \mathbf{u}, \mathbf{v} \rangle = 0$$

Gram-Schmidt Process

$$\mathbf{v}_k = \mathbf{u}_k - \sum_{j=1}^{k-1} \frac{\langle \mathbf{u}_k, \mathbf{v}_j \rangle}{\|\mathbf{v}_j\|^2} \mathbf{v}_j$$

Gives orthogonal basis;
normalise $\hat{\mathbf{v}}_k = \mathbf{v}_k / \|\mathbf{v}_k\|$ for
orthonormal basis.

SPECTRAL THEOREM

Eigenvalue Problem

$A\mathbf{x} = \lambda\mathbf{x}$, where λ is an eigenvalue and $\mathbf{x} \neq \mathbf{0}$ is the eigenvector.
Characteristic polynomial: $\det(A - \lambda I) = 0$

Spectral Theorem

Every real **symmetric** matrix $A \in \mathbb{R}^{n \times n}$ has real eigenvalues and mutually orthogonal eigenvectors:

$$A = Q\Lambda Q^T, \quad Q^T Q = I, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

- $\text{tr}(A) = \sum_i \lambda_i$ $\det(A) = \prod_i \lambda_i$
- $A \succ 0$ (positive definite) \Leftrightarrow all $\lambda_i > 0$

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization**
- 3 Vector Calculus
- 4 Continuous Optimization & Duality
- 5 Principal Component Analysis
- 6 Machine Learning Models

LINEAR PROGRAMMING PROBLEM

Standard Form $\min_{\mathbf{x}} \mathbf{c}^T \mathbf{x}$ s.t. $A\mathbf{x} = \mathbf{b}$, $\mathbf{x} \geq \mathbf{0}$

Inequality $\mathbf{a}^T \mathbf{x} \leq b$ converted via slack variable $s \geq 0$: $\mathbf{a}^T \mathbf{x} + s = b$

- Feasible region is a **convex polytope**
- **Corner-point theorem:** optimal solution occurs at a vertex (BFS)
- **BFS:** set $n-m$ non-basic variables to zero, solve for m basic variables

Simplex Method

Start at a BFS; check reduced costs \bar{c}_j . If all $\bar{c}_j \geq 0$: optimal (minimisation). Otherwise pivot via **minimum ratio test** (leaving variable) and most negative \bar{c}_j (entering variable).

ARTIFICIAL VARIABLES & DUALITY

Finding an Initial BFS

Big-M method: add artificial variable a_i with penalty M to objective.

Two-phase method:

Phase I — minimise $\sum a_i$ to find a BFS.

Phase II — optimise original objective.

Shadow price y_i^* : rate of change of optimal objective per unit increase in b_i .

Primal & Dual

$$(P) \min \mathbf{c}^\top \mathbf{x}, \mathbf{Ax} \geq \mathbf{b}, \mathbf{x} \geq \mathbf{0}$$

$$(D) \max \mathbf{b}^\top \mathbf{y}, \mathbf{A}^\top \mathbf{y} \leq \mathbf{c}, \mathbf{y} \geq \mathbf{0}$$

Complementary Slackness:

At optimality: $x_i s_i = 0$,

$$y_j t_j = 0$$

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization
- 3 Vector Calculus**
- 4 Continuous Optimization & Duality
- 5 Principal Component Analysis
- 6 Machine Learning Models

GRADIENTS & TAYLOR EXPANSION

Gradient

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \partial f / \partial x_1 \\ \vdots \\ \partial f / \partial x_n \end{bmatrix} \in \mathbb{R}^n$$

Points in direction of steepest ascent.

Directional Derivative: $D_{\mathbf{u}}f(\mathbf{x}) = \nabla f(\mathbf{x})^T \mathbf{u} = \|\nabla f\| \cos \theta$

Jacobian: $[J_f]_{ij} = \partial f_i / \partial x_j$ for $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$

Hessian: $H_{ij} = \partial^2 f / \partial x_i \partial x_j$ (symmetric for smooth f)

Taylor Expansion (2nd order)

Univariate:

$$f(x+h) \approx f(x) + f'(x)h + \frac{1}{2}f''(x)h^2$$

Multivariate:

$$f(\mathbf{x} + \mathbf{h}) \approx f(\mathbf{x}) + \nabla f^T \mathbf{h} + \frac{1}{2} \mathbf{h}^T \mathbf{H} \mathbf{h}$$

CRITICAL POINTS & EXTREMA CLASSIFICATION

First-Order Necessary Condition

\mathbf{x}^* is a local minimum $\Rightarrow \nabla f(\mathbf{x}^*) = \mathbf{0}$

Second-Order Conditions (via Hessian $H(\mathbf{x}^*)$)

$H \succ 0$ (positive definite) \Rightarrow local **minimum**

$H \prec 0$ (negative definite) \Rightarrow local **maximum**

H indefinite \Rightarrow **saddle point**

Multivariate Chain Rule

For $f(\mathbf{g}(t))$ with $\mathbf{g} : \mathbb{R} \rightarrow \mathbb{R}^n$:

$$\frac{d}{dt}f(\mathbf{g}(t)) = \nabla f(\mathbf{g}(t))^\top \frac{d\mathbf{g}}{dt} = \sum_{i=1}^n \frac{\partial f}{\partial x_i} \frac{dg_i}{dt}$$

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization
- 3 Vector Calculus
- 4 Continuous Optimization & Duality**
- 5 Principal Component Analysis
- 6 Machine Learning Models

GRADIENT DESCENT & LAGRANGE MULTIPLIERS

Gradient Descent (Unconstrained)

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

Moves in direction of steepest descent; η_k is the learning rate.

Constrained Optimisation & Lagrangian

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad \text{s.t.} \quad g_i(\mathbf{x}) \leq 0, \quad h_j(\mathbf{x}) = 0$$

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x}) + \sum_j \mu_j h_j(\mathbf{x})$$

Stationarity Condition

Unconstrained minimum: $\nabla f(\mathbf{x}^*) = \mathbf{0}$.

Constrained optimum: $\nabla_{\mathbf{x}} \mathcal{L} = \mathbf{0}$.

ADAPTIVE GRADIENT METHODS

Stochastic Gradient Descent (SGD)

Use mini-batch gradient instead of full dataset: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \mathbf{g}_t$

AdaGrad

$$G_t = G_{t-1} + g_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{G_t + \varepsilon}} \mathbf{g}_t$$

Adapts per-parameter; shrinks step for frequent features. G_t grows monotonically (drawback).

RMSProp

$$E[g^2]_t = \rho E[g^2]_{t-1} + (1-\rho) g_t^2$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta}{\sqrt{E[g^2]_t + \varepsilon}} \mathbf{g}_t$$

Fixes diminishing learning rate via exponential moving average.

ADAM OPTIMIZER

Adaptive Moment Estimation

Combines **momentum** (1st moment) with **adaptive learning rate** (2nd moment):

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t \quad (1\text{st moment})$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \quad (2\text{nd moment})$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (\text{bias correction})$$

$$w_{t+1} = w_t - \frac{\eta \hat{m}_t}{\sqrt{\hat{v}_t} + \varepsilon} \quad (\text{update})$$

Defaults: $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 10^{-8}$. Bias correction compensates for zero initialisation of m_0 , v_0 .

LAGRANGIAN DUALITY & KKT CONDITIONS

Primal & Dual

$$p^* = \min_{\mathbf{x}} \max_{\lambda \geq 0} \mathcal{L}(\mathbf{x}, \lambda)$$

$$d^* = \max_{\lambda \geq 0} \min_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda)$$

Weak Duality (always holds):

$$d^* \leq p^*$$

Strong Duality: $d^* = p^*$ under

Slater's condition — $\exists \mathbf{x}$ strictly feasible ($g_i(\mathbf{x}) < 0$), f and g_i convex.

KKT Conditions

(necessary & sufficient for convex)

- 1 **Stationarity:** $\nabla_{\mathbf{x}} \mathcal{L} = 0$
- 2 **Primal feasibility:**
 $g_i(\mathbf{x}^*) \leq 0$
- 3 **Dual feasibility:** $\lambda_i^* \geq 0$
- 4 **Comp. slackness:**
 $\lambda_i^* g_i(\mathbf{x}^*) = 0$

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization
- 3 Vector Calculus
- 4 Continuous Optimization & Duality
- 5 Principal Component Analysis**
- 6 Machine Learning Models

PCA: MAXIMUM VARIANCE PERSPECTIVE

Goal: Find M orthogonal directions $\mathbf{b}_1, \dots, \mathbf{b}_M \in \mathbb{R}^D$ capturing maximum variance in centred data $\{\mathbf{x}_n\}_{n=1}^N$.

Covariance Matrix & Eigenvalue Problem

$$S = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \in \mathbb{R}^{D \times D}$$

Maximise $\mathbf{b}^\top S \mathbf{b}$ subject to $\|\mathbf{b}\| = 1 \Rightarrow S \mathbf{b} = \lambda \mathbf{b}$

- **Principal components:** eigenvectors of S , sorted by $\lambda_1 \geq \lambda_2 \geq \dots$
- **Projected data:** $\mathbf{y}_n = B^\top \mathbf{x}_n$, $B = [\mathbf{b}_1, \dots, \mathbf{b}_M]$
- **Explained variance** of k -th PC: $\lambda_k / \sum_j \lambda_j$

PCA VIA SVD & LOW-RANK APPROXIMATION

Singular Value Decomposition

$$X = U\Sigma V^T, \quad U^T U = I, \quad V^T V = I$$
$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p), \quad \sigma_1 \geq \dots \geq \sigma_p \geq 0$$

Principal components \equiv columns of V ; $\lambda_i \propto \sigma_i^2$.

Eckart-Young Theorem

Best rank- M approximation: $\tilde{X}_M = \sum_{i=1}^M \sigma_i \mathbf{u}_i \mathbf{v}_i^T$

Power iteration: $\mathbf{v}_{k+1} = A\mathbf{v}_k / \|A\mathbf{v}_k\|$ — finds dominant eigenvector iteratively

TABLE OF CONTENTS

- 1 Analytic Geometry
- 2 Linear Optimization
- 3 Vector Calculus
- 4 Continuous Optimization & Duality
- 5 Principal Component Analysis
- 6 Machine Learning Models**

GAUSSIAN MIXTURE MODELS

GMM

$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad \sum_k \pi_k = 1, \quad \pi_k \geq 0$$

Expectation-Maximization (EM)

Iteratively maximise $\sum_n \log p(\mathbf{x}_n)$ with latent cluster assignments.

E-step — compute responsibilities:

$$r_{nk} = \frac{\pi_k \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_j \pi_j \mathcal{N}(\mathbf{x}_n \mid \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}$$

M-step — update parameters:

$$\boldsymbol{\mu}_k^{\text{new}} = \frac{\sum r_{nk} \mathbf{x}_n}{\sum_n r_{nk}}, \quad \pi_k^{\text{new}} = \frac{\sum r_{nk}}{N}$$

SUPPORT VECTOR MACHINE (SVM)

Hard Margin SVM

Linearly separable. Boundary:

$$\mathbf{w}^\top \mathbf{x} + b = 0, \text{ margin} = 2/\|\mathbf{w}\|$$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1$$

Support vectors: points on
 $y_i(\mathbf{w}^\top \mathbf{x}_i + b) = 1$

The SVM dual (via KKT, Sec. 4) reveals that only support vectors determine the decision boundary.

Hinge Loss: $\ell_i = \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$ (zero for correctly classified points beyond margin)

Soft Margin SVM

Non-separable. Slack $\xi_i \geq 0$:

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i$$

$$\text{s.t. } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i, \xi_i \geq 0$$

C trades margin vs. violations.

AUTOREGRESSIVE & MOVING AVERAGE MODELS

Stationarity

Mean, variance, and autocovariance constant over time. Non-stationary series made stationary by **differencing** (d times).

AR(p) Model

$$y_t = \mu + \sum_{i=1}^p \phi_i y_{t-i} + \varepsilon_t$$

Uses p past observations.

- Use **PACF** to identify p (cuts off after lag p); use **ACF** to identify q (cuts off after lag q)
- SMA, WMA, EMA smooth the series and remove noise

MA(q) Model

$$y_t = \mu + \varepsilon_t + \sum_{j=1}^q \theta_j \varepsilon_{t-j}$$

Uses q past errors.

ARIMA, SARIMA & KALMAN FILTER

ARIMA(p, d, q)

Difference d times for stationarity, then fit ARMA(p, q).

SARIMA(p, d, q) \times (P, D, Q, s)

Adds seasonal AR and MA terms at period s .

Kalman Filter (state-space model)

Predict:

$$\hat{\mathbf{x}}_{t|t-1} = F\hat{\mathbf{x}}_{t-1|t-1}$$

$$P_{t|t-1} = FP_{t-1|t-1}F^T + Q$$

Update:

$$K_t = P_{t|t-1}C^T(CP_{t|t-1}C^T + R)^{-1}$$

$$\hat{\mathbf{x}}_{t|t} = \hat{\mathbf{x}}_{t|t-1} + K_t(\mathbf{y}_t - C\hat{\mathbf{x}}_{t|t-1})$$

K_t : Kalman gain — optimally weights prediction vs. observation.

KEY TAKEAWAYS

- **Module thread:** Geometry \rightarrow LP \rightarrow Calculus \rightarrow Optimisation \rightarrow ML Models — each module builds on the previous
- **KKT conditions** (M4) are the bridge: they power the SVM dual, LASSO, and neural network training
- **PCA** (M5) uses both the Spectral Theorem (M1) and SVD — showing the full mathematical circle
- **Iterative optimisers everywhere:** EM (GMM), ARIMA fitting, and Kalman filter are all gradient-descent-style iterative schemes applied to structured models

Thank you :)