

# Quantitative and Modeling Methods for Health Informatics

SDS6210 – Informatics for Health

Comprehensive Presentations (Q1-Q4)

---

## Group 5

January 20, 2026

University of Nairobi  
Department of Mathematics

# Unit and Group Information

Unit	SDS6210 – Informatics for Health
Programme	MSc Public Health Data Science
Institution	University of Nairobi
Department	Department of Mathematics
<b>Group 5 Members</b>	
Cavin Otieno (SDS6/46982/2024)	
Laura Nabalayo Kundu (SDS6/47543/2024)	
John Andrew (SDS6/47659/2024)	

# Overview: Quantitative and Modeling Module

This comprehensive presentation covers four fundamental quantitative methods essential for health informatics and data science applications in healthcare settings.

# **Q1: Maximum Likelihood Estimation and Logistic Regression**

---

# **Q1: Maximum Likelihood Estimation and Logistic Regression**

---

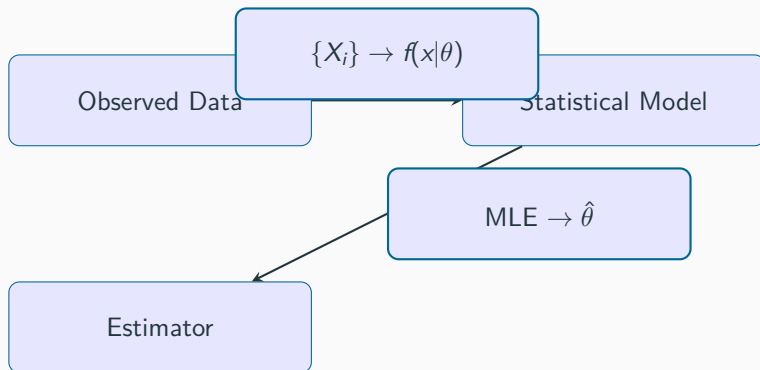
**Foundational Concepts of MLE**

# The Problem of Statistical Estimation

Statistical inference involves estimating unknown parameters from observed data:

## The Estimation Problem

Given a random sample  $X_1, X_2, \dots, X_n$  from a distribution with probability density function  $f(x|\theta)$ , where  $\theta$  is an unknown parameter, we seek to estimate  $\theta$  based on the observed data.



# Definition of Maximum Likelihood Estimation

Maximum Likelihood Estimation (MLE) is a fundamental method for estimating statistical parameters:

## Fundamental Principle

The likelihood function measures the plausibility of different parameter values given the observed data. The MLE is the parameter value that makes the observed data most probable (most likely).

## Definition (Likelihood Function)

Given a sample  $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$  from a distribution with pdf  $f(x|\theta)$ , the likelihood function is:

$$L(\theta) = L(\theta|x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i|\theta)$$

## Maximum Likelihood Estimator

The MLE  $\hat{\theta}$  is the value of  $\theta$  that maximizes the likelihood function:

## Example: Bernoulli and Normal Distributions

MLE yields familiar estimators for common distributions:

### Bernoulli Distribution

For  $X_i \sim \text{Bernoulli}(p)$ , the MLE is  $\hat{p} = \frac{y}{n}$  where  $y = \sum x_i$ .

The likelihood is  $L(p) = p^y(1-p)^{n-y}$ , maximizing this gives  $\hat{p} = \frac{y}{n}$ .

### Normal Distribution

For  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ :

$$\hat{\mu} = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

Note:  $\hat{\sigma}^2$  is biased (uses  $n$  instead of  $n-1$ ).



# The Log-Likelihood Function

The log-likelihood transformation simplifies maximization:

## Log-Likelihood Definition

$$\ell(\theta) = \log L(\theta) = \log \left( \prod_{i=1}^n f(x_i|\theta) \right) = \sum_{i=1}^n \log f(x_i|\theta)$$

## Benefits

- Transforms products to sums
- Converts multiplication to addition
- Provides numerical stability
- Simplifies differentiation

For Bernoulli:  $\ell(p) = y \log p + (n - y) \log(1 - p)$

# Properties of MLE Estimators

MLE possesses several desirable statistical properties:

## Consistency

As  $n \rightarrow \infty$ ,  $\hat{\theta}_n \xrightarrow{P} \theta$ , meaning the estimator converges to the true parameter value.

## Asymptotic Normality

For large  $n$ ,  $\hat{\theta} \approx \mathcal{N}(\theta, \text{Var}(\hat{\theta}))$ .

## Asymptotic Efficiency

The MLE achieves the Cramér-Rao lower bound asymptotically.

## Invariance Property

If  $\hat{\theta}$  is the MLE of  $\theta$ , then  $g(\hat{\theta})$  is the MLE of  $g(\theta)$ .

# **Q1: Maximum Likelihood Estimation and Logistic Regression**

---

## **Logistic Regression Model**

# Why Logistic Regression?

When the outcome variable is binary ( $Y \in \{0, 1\}$ ), linear regression is problematic:

## Problems with Linear Probability Model

$$P(Y = 1|X) = \beta_0 + \beta_1 X_1 + \cdots + \beta_k X_k$$

- Predicted probabilities can be  $< 0$  or  $> 1$
- Heteroscedasticity: Variance depends on  $X$
- Error terms are not normally distributed

The logistic function constrains predictions to  $(0, 1)$  while modeling the log-odds linearly.

# The Logistic Function

The logistic (sigmoid) function transforms linear predictors to probabilities:

## Logistic (Sigmoid) Function

$$\pi(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z}$$

where  $z$  is the linear predictor and  $\pi(z) \in (0, 1)$  for all  $z \in \mathbb{R}$ .

## Key Properties

- $\pi(-z) = 1 - \pi(z)$  (symmetry)
- $\pi(0) = 0.5$  (median effective dose)
- As  $z \rightarrow \infty$ ,  $\pi(z) \rightarrow 1$

# The Logistic Regression Model

Logistic regression models the log-odds as a linear function of predictors:

## Probability Model

For a binary outcome  $Y \in \{0, 1\}$  and covariates  $X_1, X_2, \dots, X_k$ :

$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}}$$

## Log-Odds Form (Logit)

Taking the log-odds (logit transformation):

$$\log \left( \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

This linearizes the model while constraining probabilities to  $(0, 1)$ .

# Assumptions of Logistic Regression

## Model Assumptions

- Binary outcome (Bernoulli distribution)
- Linearity in the log-odds
- Independence of observations
- No perfect multicollinearity

## GLM Framework

Logistic regression is a Generalized Linear Model (GLM) with:

- **Distribution:** Binomial/Bernoulli
- **Link function:** Logit (log-odds)

# **Q1: Maximum Likelihood Estimation and Logistic Regression**

---

**Likelihood Derivation for Logistic  
Regression**



# Data Structure and Likelihood Function

Let the observed data consist of  $n$  independent observations:

## Notation

For each observation  $i = 1, \dots, n$ :

- $Y_i \in \{0, 1\}$ : Binary outcome
- $X_i = (1, X_{i1}, \dots, X_{ik})$ : Vector of predictors
- $\beta = (\beta_0, \beta_1, \dots, \beta_k)^T$ : Vector of coefficients

## Likelihood Function

$$L(\beta) = \prod_{i=1}^n P(Y_i = y_i | X_i, \beta) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

where  $\pi_i = \frac{1}{1 + e^{-X_i^T \beta}}$ .

# The Log-Likelihood Function

Taking the log of the likelihood simplifies the product to a sum:

## Log-Likelihood Derivation

$$\begin{aligned}\ell(\beta) &= \log L(\beta) \\ &= \sum_{i=1}^n \log [P(Y_i = y_i | X_i, \beta)] \\ &= \sum_{i=1}^n [y_i \log \pi_i + (1 - y_i) \log(1 - \pi_i)] \\ &= \sum_{i=1}^n [y_i X_i^T \beta - \log(1 + \exp(X_i^T \beta))]\end{aligned}$$

## Goal

Find  $\hat{\beta}$  that maximizes  $\ell(\beta)$ :

$$\hat{\beta} = \arg \max \ell(\beta)$$

# The Score Function and Hessian

To find the MLE, we set the score function to zero:

## Score Vector (First Derivative)

$$U(\beta) = \frac{\partial \ell(\beta)}{\partial \beta} = \sum_{i=1}^n X_{ij}(y_i - \pi_i)$$

The score equations are  $U_j(\beta) = 0$ ,  $j = 0, 1, \dots, k$ .

## Hessian Matrix (Second Derivative)

$$\mathbf{H}(\beta) = -\mathbf{X}^T \mathbf{W} \mathbf{X}$$

where  $\mathbf{W} = \text{diag}\{\pi_1(1 - \pi_1), \dots, \pi_n(1 - \pi_n)\}$  is the weight matrix.

# **Q1: Maximum Likelihood Estimation and Logistic Regression**

---

**Numerical Estimation and  
Interpretation**

# Newton-Raphson and IRLS Algorithm

Unlike linear regression, logistic regression has no closed-form solution:

## Newton-Raphson Update

$$\beta^{(t+1)} = \beta^{(t)} - \mathbf{H}^{-1}(\beta^{(t)}) U(\beta^{(t)})$$

$$\beta^{(t+1)} = \beta^{(t)} + (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}^{(t)})$$

## Fisher Scoring (IRLS)

The update equation resembles weighted least squares:

$$\beta^{(t+1)} = (\mathbf{X}^T \mathbf{W}^{(t)} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}^{(t)} \mathbf{z}^{(t)}$$

where  $z_i^{(t)} = X_i^T \beta^{(t)} + \frac{y_i - \pi_i^{(t)}}{\pi_i^{(t)}(1 - \pi_i^{(t)})}$ .

Typically converges in 5-10 iterations.

# Interpreting Estimated Coefficients

The coefficient  $\beta_j$  represents the change in log-odds per unit change in  $X_j$ :

## Odds Ratio Interpretation

For a one-unit increase in predictor  $X_j$ :

$$\text{OR} = e^{\beta_j} = \frac{\text{odds}(Y = 1 | X_j + 1)}{\text{odds}(Y = 1 | X_j)}$$

$$\beta_j > 0 \implies e^{\beta_j} > 1: \text{Increased odds of outcome}$$

$$\beta_j = 0 \implies e^{\beta_j} = 1: \text{No association}$$

$$\beta_j < 0 \implies e^{\beta_j} < 1: \text{Decreased odds of outcome}$$

# Applications in Public Health

Logistic regression is widely used for disease risk modeling:

Predictor	$\hat{\beta}$	OR	95% CI
Age (per 10 years)	0.72	2.05	(1.85, 2.28)
Male (vs. Female)	0.45	1.57	(1.32, 1.86)
SBP (per 20 mmHg)	0.38	1.46	(1.31, 1.63)
Current Smoker	0.89	2.44	(2.05, 2.90)
Diabetes	1.12	3.06	(2.51, 3.73)

Used in Framingham Risk Score and numerous clinical prediction models.

## Q1 Summary: Key Equations

**Logistic Model:** 
$$P(Y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}}$$

**Log-Odds:** 
$$\log \left( \frac{P(Y = 1|X)}{1 - P(Y = 1|X)} \right) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

**Log-Likelihood:** 
$$\ell(\beta) = \sum_{i=1}^n \left[ y_i X_i^T \beta - \log(1 + e^{X_i^T \beta}) \right]$$

**Score Function:** 
$$U_j(\beta) = \sum_{i=1}^n X_{ij}(y_i - \pi_i)$$

**Odds Ratio:** 
$$OR = e^{\beta_j}$$



## **Q2: Cox Proportional Hazards Model and Hazard Ratios**

---

## **Q2: Cox Proportional Hazards Model and Hazard Ratios**

---

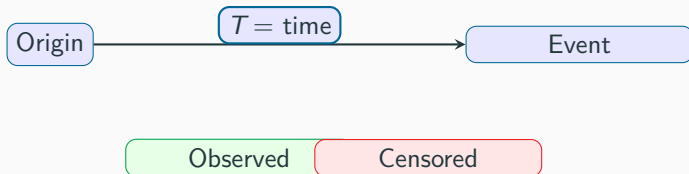
**Fundamentals of Survival Analysis**

# Survival Data: Time-to-Event Analysis

Survival analysis studies the time until the occurrence of a specific event:

## Definition of Survival Time

Let  $T$  denote the survival time, defined as the time from a defined origin (e.g., diagnosis, treatment initiation) to the occurrence of a specified event (e.g., death, disease recurrence, recovery).



# Censoring: The Central Challenge

In survival analysis, complete event times are often not observed:

## Types of Censoring

### Right Censoring

The event has not occurred by the end of the study period. We know  $T > c$ .

### Left Censoring

The event occurred before the study began.  
We know  $T < c$ .

## Data Representation

For each subject  $i = 1, \dots, n$ :

- $t_i$ : The observed time (minimum of survival time and censoring time)
- $\delta_i$ : The event indicator ( $\delta_i = 1$  if event observed, 0 if censored)

# The Survival and Hazard Functions

Two fundamental functions in survival analysis:

## Definition (Survival Function)

$$S(t) = P(T > t)$$

This represents the probability that the event has not occurred by time  $t$ .

## Definition (Hazard Function)

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t \mid T \geq t)}{\Delta t}$$

The hazard function describes the instantaneous risk of the event at time  $t$ , conditional on survival to that time.

## Relationship

$$h(t) = -\frac{d}{dt} \log S(t), \quad S(t) = \exp \left( - \int_0^t h(u) du \right)$$

## **Q2: Cox Proportional Hazards Model and Hazard Ratios**

---

### **The Cox Model Formulation**

# The Cox Proportional Hazards Model

The Cox model extends regression to time-to-event data:

## The Cox Model Equation

$$h(t|X) = h_0(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_k X_k)$$

## Model Components

- $h(t|X)$ : Hazard at time  $t$  given covariates  $X$
- $h_0(t)$ : Baseline hazard function
- $X = (X_1, \dots, X_k)$ : Vector of covariates

## Model Classification

### Semi-parametric:

- $h_0(t)$  is non-parametric (unspecified)
- Covariate effects are parametric

# The Baseline Hazard Function

The baseline hazard  $h_0(t)$  represents the hazard for an individual with all covariates equal to zero:

## Properties of $h_0(t)$

- Non-parametric: No distributional assumption required
- Unknown function that can take any non-negative shape
- Acts as a time-varying intercept in the model

## Log-Linear Formulation

$$\log h(t|X) = \log h_0(t) + \beta^T X$$

The coefficient  $\beta_j$  represents the change in the log-hazard associated with a one-unit increase in  $X_j$ .



# Proportional Hazards Assumption

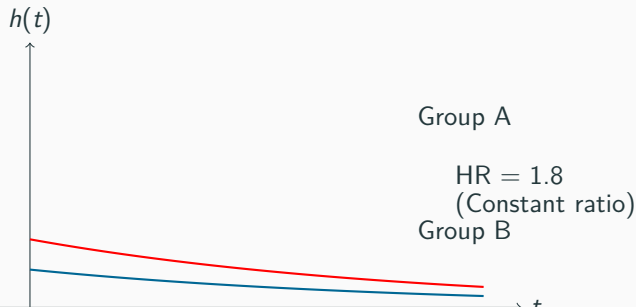
The term "proportional hazards" refers to the constant ratio of hazards:

## Proportionality Assumption

For any two individuals with covariates  $X_A$  and  $X_B$ :

$$\frac{h(t|X_A)}{h(t|X_B)} = \exp(\beta^T(X_A - X_B))$$

This ratio is constant for all time  $t$ .



## **Q2: Cox Proportional Hazards Model and Hazard Ratios**

---

**Hazard Ratios: Definition and  
Interpretation**

# The Hazard Ratio as a Measure of Effect

The hazard ratio (HR) is the primary measure of association in survival analysis:

## Definition (Hazard Ratio)

The hazard ratio compares the hazards of two groups:

$$\text{HR} = \frac{h(t|X_A)}{h(t|X_B)} = \exp(\beta^T(X_A - X_B))$$

The baseline hazard cancels out, leaving only the covariate effects.

## Key Result

$$\text{HR} = \exp(\beta_j)$$

The exponentiated coefficient directly gives the hazard ratio for a one-unit change.

## Interpretation of Hazard Ratios

$\beta > 0 \implies \text{HR} > 1$ : The covariate increases the hazard (risk factor)

$\beta = 0 \implies \text{HR} = 1$ : No association

$\beta < 0 \implies \text{HR} < 1$ : The covariate decreases the hazard (protective factor)

## Numerical Interpretation Examples

$\beta$	HR = $\exp(\beta)$	Interpretation	95% CI for HR
0	1.00	No effect	(0.82, 1.22)
0.10	1.11	11% increased hazard	(0.91, 1.35)
0.50	1.65	65% increased hazard	(1.35, 2.02)
1.00	2.72	172% increased hazard	(2.23, 3.32)
-0.50	0.61	39% decreased hazard	(0.50, 0.74)
-1.00	0.37	63% decreased hazard	(0.30, 0.45)

## Continuous Covariates: Unit Changes

For a continuous predictor, the HR represents the hazard ratio per unit increase:

### Example: Age and Cardiovascular Disease

If  $\beta_{\text{age}} = 0.08$  (per year):

- Per 1 year:  $\text{HR} = \exp(0.08) = 1.083$  (8.3% increased hazard)
- Per 10 years:  $\text{HR} = \exp(0.80) = 2.23$  (123% increased hazard)

## Q2 Summary: Key Equations

**Hazard Function:** 
$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

**Survival Function:** 
$$S(t) = P(T > t)$$

**Cox Model:** 
$$h(t|X) = h_0(t) \exp(\beta_1 X_1 + \cdots + \beta_k X_k)$$

**Hazard Ratio:** 
$$HR = \exp(\beta_j)$$

**Data:** Observed time  $t_i = \min(T_i, C_i)$ ,  $\delta_i = I(T_i \leq C_i)$

## **Q3: Principal Component Analysis for Dimensionality Reduction**

---



# **Q3: Principal Component Analysis for Dimensionality Reduction**

---

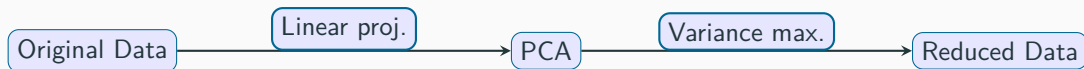
## **Foundational Concepts of PCA**

# What is Principal Component Analysis?

Principal Component Analysis (PCA) is a linear dimensionality reduction technique:

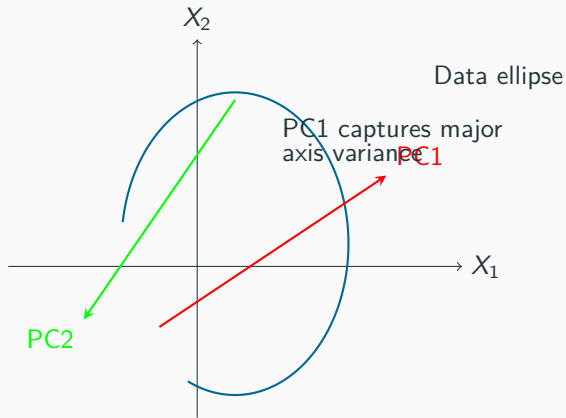
## Definition

PCA is an unsupervised learning technique that transforms a set of correlated variables into a smaller set of uncorrelated variables called **principal components** (PCs), while retaining the maximum possible variance in the original data.



# Geometric Interpretation of PCA

PCA finds the directions of maximum variance in the data:



## Key Principle

PC1 captures the maximum variance direction; PC2 is orthogonal to PC1 and captures remaining maximum variance.

# Mathematical Formulation of PCA

The goal of PCA is to find a linear transformation that maximizes variance:

## Data Matrix

Let  $\mathbf{X}$  be an  $n \times p$  data matrix where:

- $n$ : Number of observations (patients)
- $p$ : Number of variables (15 clinical indicators)

## Principal Components as Linear Combinations

Each principal component  $PC_j$  is a linear combination:

$$PC_j = a_{j1}X_1 + a_{j2}X_2 + \cdots + a_{jp}X_p$$

where  $\mathbf{a}_j = (a_{j1}, \dots, a_{jp})$  is the loading vector (eigenvector).

# The Problem of Dimensionality in Clinical Research

Modern clinical studies often collect numerous measurements per patient:

High correlation among clinical indicators

Small sample sizes relative to number of variables

Multicollinearity complicates regression modeling

Difficult interpretation of 15+ correlated variables

# Benefits of Dimensionality Reduction

Reducing 15 indicators to 4 principal components offers several advantages:

Dimensionality reduction:  $15 \rightarrow 4$  variables (73% reduction)

Decorrelation: Components are uncorrelated (eliminates multicollinearity)

Noise reduction: Low-variance components often capture measurement error

Interpretability: Each PC represents a clinical dimension

## Step 1: Data Standardization

Before performing PCA, variables must be standardized:

### Standardization Formula

For each variable  $X_j$ :

$$Z_{ij} = \frac{X_{ij} - \bar{X}_j}{s_j}$$

where  $\bar{X}_j$  is the mean and  $s_j$  is the standard deviation.

### Problem of Scale Differences

Clinical indicators often have different scales (e.g., BP in mmHg, cholesterol in mg/dL). Without standardization, larger-scale variables would dominate.

### Result

After standardization, each variable has mean 0 and variance 1.

## Step 2: Computing the Covariance/Correlation Matrix

The covariance matrix captures the relationships between variables:

### Sample Correlation Matrix

For standardized data  $\mathbf{Z}$ :

$$\mathbf{R} = \frac{1}{n-1} \mathbf{Z}^T \mathbf{Z}$$

where  $\mathbf{R}$  is a  $p \times p$  symmetric matrix with 1s on the diagonal.

### Example Correlations in Clinical Data

Variable Pair	Interpretation	Correlation
Systolic BP - Diastolic BP	Strong positive	0.78
HDL - Triglycerides	Moderate negative	-0.52
Creatinine - eGFR	Strong negative	-0.85



## Step 3: Eigenvalue Decomposition

The core of PCA is the eigenvalue decomposition of the correlation matrix:

### Eigenvalue Equation

Find eigenvalues  $\lambda_j$  and eigenvectors  $\mathbf{v}_j$  such that:

$$\mathbf{R}\mathbf{v}_j = \lambda_j\mathbf{v}_j$$

where  $\mathbf{v}_j^T\mathbf{v}_j = 1$  and  $\mathbf{v}_j^T\mathbf{v}_k = 0$  for  $j \neq k$ .

### Eigen Decomposition

$$\mathbf{R} = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^T$$

where  $\mathbf{V}$  contains eigenvectors and  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_p)$ .

### Properties

Eigenvalues  $\lambda_j \geq 0$  (since  $\mathbf{R}$  is positive semi-definite).

## Eigenvalues as Variance Explained

Each eigenvalue represents the variance explained by its corresponding principal component:

### Variance Interpretation

For standardized data:  $\lambda_j = \text{Var}(PC_j)$  and  $\sum_{j=1}^p \lambda_j = p$ .

### Example Eigenvalue Spectrum

Component	Eigenvalue ( $\lambda_j$ )	Variance %	Cumulative %
PC1	4.52	30.1	30.1
PC2	2.87	19.1	49.2
PC3	2.14	14.3	63.5
PC4	1.68	11.2	74.7

## Step-by-Step Algorithm Summary

**Step 1:** Standardize data:  $Z_{ij} = (X_{ij} - \bar{X}_j)/s_j$

**Step 2:** Compute correlation matrix:  $\mathbf{R} = \frac{1}{n-1} \mathbf{Z}^T \mathbf{Z}$

**Step 3:** Eigenvalue decomposition:  $\mathbf{R} \mathbf{v}_j = \lambda_j \mathbf{v}_j$

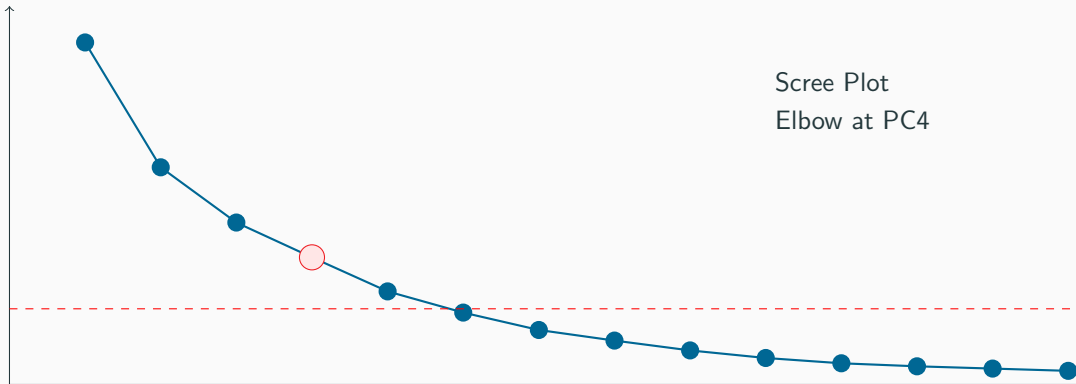
**Step 4:** Order by eigenvalues:  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{15}$

**Step 5:** Select top  $k$  components and compute scores:  $\mathbf{PC} = \mathbf{Z} \mathbf{V}_k$

# Scree Plot for Component Selection

The scree plot visualizes eigenvalues to guide component selection:

Eigenvalue



Scree Plot  
Elbow at PC4

## Justifying 4 Principal Components

Several criteria support selecting 4 components from 15:

### **Kaiser Criterion**

Retain components with eigenvalue  $> 1$ : PC1 ( $\lambda = 4.52$ ), PC2 ( $\lambda = 2.87$ ), PC3 ( $\lambda = 2.14$ ), PC4 ( $\lambda = 1.68$ ) all exceed 1.

### **Cumulative Variance**

4 components explain 74.7% of total variance, exceeding the commonly recommended 70-75% threshold.

### **Elbow Criterion**

Visual inspection of the scree plot shows an elbow at PC4.

# Interpreting Principal Components in Clinical Context

Each principal component represents a clinical dimension:

Component	Clinical Interpretation	Variance %
PC1	<b>Cardiovascular Risk:</b> Systolic BP, Diastolic BP, Pulse Pressure, LDL, Triglycerides	30.1
PC2	<b>Metabolic Syndrome:</b> BMI, Waist Circumference, Fasting Glucose, HbA1c	19.1
PC3	<b>Inflammatory Status:</b> CRP, ESR, WBC Count	14.3
PC4	<b>Renal Function:</b> Creatinine, eGFR	11.2

Cumulative variance: 74.7% retained.

# Computing Principal Component Scores

Principal component scores are used for subsequent analysis:

## Score Calculation

For each observation  $i$  and component  $j$ :

$$PC_{ij} = z_{i1}v_{1j} + z_{i2}v_{2j} + \cdots + z_{i,15}v_{15,j}$$

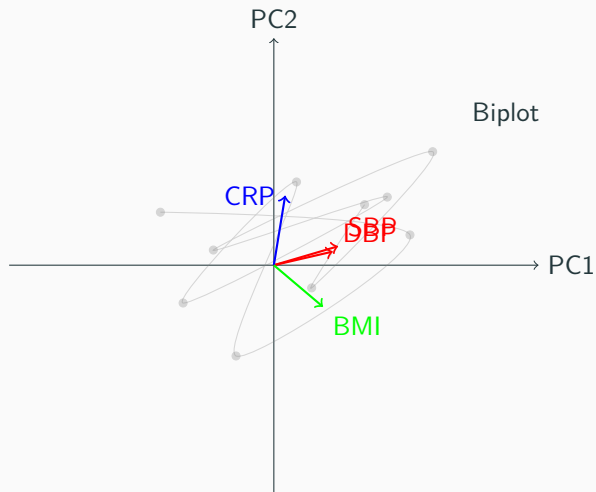
In matrix form:  $\mathbf{PC} = \mathbf{ZV}_4$

## Properties of Scores

- Each column has mean 0
- Each column has variance equal to the corresponding eigenvalue

## Biplot for Joint Interpretation

The biplot visualizes observations and variables simultaneously:



**Biplot Interpretation**



# Limitations of PCA

PCA has important limitations:

**Linearity:** PCA only captures linear relationships

**No class discrimination:** Components may not separate disease groups

**Interpretability:** Component meaning may be unclear

**Information loss:** Some variance is discarded

## Q3 Summary: Key Equations

**Standardization:**  $Z_{ij} = \frac{X_{ij} - \bar{X}_j}{s_j}$

**Correlation Matrix:**  $\mathbf{R} = \frac{1}{n-1} \mathbf{Z}^T \mathbf{Z}$

**Eigen Decomposition:**  $\mathbf{R} \mathbf{v}_j = \lambda_j \mathbf{v}_j$

**Variance Explained:**  $VE_j = \frac{\lambda_j}{\sum_{m=1}^{15} \lambda_m} \times 100\%$

**Component Scores:**  $PC_{ij} = \sum_{k=1}^{15} z_{ik} v_{kj}$

**Cumulative Variance:**  $CVE_4 = 74.7\%$

## **Q4: Machine Learning Workflow in a Hospital Setting**

---

## **Q4: Machine Learning Workflow in a Hospital Setting**

---

**Purpose of ML in Healthcare**

# Why Machine Learning in Healthcare?

Machine learning represents a transformative paradigm in healthcare delivery:

- **Clinical Decision Support:** Real-time risk stratification and early warning systems for sepsis, cardiac arrest, and deterioration.
- **Operational Efficiency:** Resource allocation optimization, patient flow management, and predictive scheduling.
- **Diagnostic Accuracy:** Enhanced image interpretation, pattern recognition in laboratory data, and differential diagnosis assistance.
- **Predictive Analytics:** Hospital readmission prediction, length of stay estimation, and population health management.

The fundamental shift involves moving from reactive medicine to proactive, personalized, and precision-based clinical interventions.

# From Retrospective Analysis to Real-Time Clinical Applications

Machine learning enables a paradigm shift toward real-time predictions:

- **Retrospective Phase:** Secondary analysis of administrative claims, quality improvement initiatives.
- **Proactive Phase:** Prospective risk modeling embedded within electronic health record systems.
- **Real-Time Phase:** Streaming analytics providing continuous risk scores and actionable alerts at the point of care.

## Key Distinction

The transition from descriptive and diagnostic analytics to predictive and prescriptive analytics represents the core value proposition of ML in modern hospital settings.

# Defining the Clinical Problem

The foundation of any successful ML implementation begins with a precisely formulated clinical question:

- **Clinical Validity:** Does the outcome have clinical significance and actionable intervention options?
- **Feasibility Assessment:** Are the required data elements available and of sufficient quality?
- **Stakeholder Alignment:** Engagement with clinical champions, informaticists, and administrative leadership.
- **Success Metrics:** Clear definition of primary and secondary outcomes.

Common applications include mortality prediction, readmission risk, sepsis early warning, and length of stay estimation.

# Data Sources in Hospital Settings

Hospital data ecosystems comprise multiple interconnected sources:

- **Electronic Health Records (EHR):** Demographics, vital signs, laboratory results, medications, procedures, clinical notes.
- **Administrative Data:** Billing codes, diagnosis codes (ICD-10), procedure codes (CPT), discharge summaries.
- **Real-Time Monitoring Data:** Bedside monitors, infusion pumps, ventilators, and other connected medical devices.
- **External Data Sources:** Social determinants of health, regional epidemiological data.

## Interoperability Standards

HL7 FHIR (Fast Healthcare Interoperability Resources) serves as the primary standard for healthcare data exchange.



# Data Quality and Preprocessing Challenges

Healthcare data presents unique preprocessing challenges:

- **Missing Data:** Systematic patterns of missingness (MNAR, MAR, MCAR) require appropriate imputation strategies.
- **Temporal Heterogeneity:** Changes in clinical practices over time introduce distribution shifts.
- **Irregular Sampling:** Laboratory results and vital signs at irregular intervals.
- **Coding Variations:** ICD-10 code granularity and evolving code sets require standardization.

# Feature Engineering in Clinical Domains

Feature engineering transforms raw clinical data into meaningful predictors:

- **Temporal Aggregation:** Computing trends, rates of change (e.g., delta neutrophil index, shock index).
- **Domain-Specific Transformations:** Composite scores such as SOFA, NEWS, and qSOFA.
- **Natural Language Processing:** Extraction of clinical entities from unstructured notes.
- **Categorical Encodings:** Hierarchical groupings of diagnoses and procedures.

# Feature Selection Strategies

Feature selection aims to identify the most predictive subset of features:

- **Filter Methods:** Correlation analysis, mutual information, variance thresholds.
- **Wrapper Methods:** Recursive feature elimination (RFE), forward/backward selection.
- **Embedded Methods:** L1 regularization (LASSO), tree-based importance scores.
- **Clinical Validation:** Feature review by subject matter experts.

Balancing statistical significance with clinical interpretability is crucial.

# Model Selection Considerations

The selection of ML algorithms requires careful consideration of trade-offs:

- **Interpretable Models:** Logistic regression, decision trees, score-based models.
- **Black-Box Models:** Gradient boosting machines, random forests, deep neural networks.
- **Time-to-Event Models:** Cox proportional hazards for survival outcomes.

## The Interpretability-Performance Trade-off

While complex models may achieve higher discrimination, interpretable models facilitate clinical acceptance and regulatory review.

# Model Training Best Practices

Rigorous model training protocols ensure reproducibility and generalizability:

- **Train-Validation-Test Split:** Temporal splits for time-series data to prevent data leakage.
- **Cross-Validation:** Stratified k-fold cross-validation for robust performance estimation.
- **Class Imbalance Handling:** SMOTE, class weighting, threshold optimization.
- **Hyperparameter Tuning:** Grid search, random search, or Bayesian optimization.

# Validation Strategies

Comprehensive validation assesses performance across multiple dimensions:

- **Internal Validation:** Bootstrap and cross-validation estimates.
- **Temporal Validation:** Performance on holdout data from a later time period.
- **External Validation:** Testing on independent datasets from different institutions.
- **Subgroup Validation:** Performance assessment across demographic subgroups.

# Performance Metrics for Clinical Models

Evaluation metrics must align with clinical objectives:

- **Discrimination:** Area under the ROC curve (AUC-ROC).
- **Calibration:** Hosmer-Lemeshow test, calibration curves, Brier score.
- **Clinical Utility:** Decision curve analysis and net benefit calculations.
- **Sensitivity and Specificity:** Performance at clinically relevant operating points.

# Deployment Strategies

Transitioning from research prototypes to clinical systems:

- **Shadow Mode Deployment:** Running the model without exposing predictions to end-users.
- **Pilot Implementation:** Limited deployment in specific clinical units with enhanced monitoring.
- **Full Integration:** Production deployment within the EHR with user-facing interfaces.

Technical infrastructure must support real-time inference and high availability.



# Integration into Clinical Workflows

Successful deployment depends on thoughtful integration:

- **User Interface Design:** In-context displays within clinician workflows.
- **Alert Design:** Tiered alert severity, clear explanatory information.
- **Feedback Mechanisms:** Capture of clinician acceptance/rejection for improvement.
- **EHR Integration Standards:** SMART on FHIR applications for interoperability.

# Model Governance Framework

Establishing robust governance structures ensures accountability:

- **Model Registry:** Documentation of model versions, training data, performance metrics.
- **Approval Workflows:** Multidisciplinary review committees.
- **Audit Trails:** Comprehensive logging of predictions and outcomes.
- **Incident Response:** Defined procedures for model-related adverse events.

# Model Monitoring and Drift Detection

Continuous monitoring is essential to detect performance degradation:

- **Data Drift Detection:** Statistical monitoring of input feature distributions.
- **Concept Drift Detection:** Tracking of outcome rates over time.
- **Performance Monitoring:** Ongoing calculation of performance metrics.
- **Alert Thresholding:** Statistical process control methods.

# Model Updating and Retraining Strategies

Maintaining model relevance requires planned updating protocols:

- **Scheduled Retraining:** Periodic updates based on accumulated new data.
- **Triggered Retraining:** Updates initiated by detected performance degradation.
- **Online Learning:** Continuous model adaptation with safeguards.
- **Version Control:** Careful management of model versions with validation requirements.

# Data Privacy Requirements

Healthcare ML development must comply with stringent privacy regulations:

- **Regulatory Compliance:** Adherence to HIPAA (US), GDPR (Europe).
- **De-identification:** Safe harbor or expert determination methods.
- **Access Controls:** Role-based access, audit logging, principle of least privilege.
- **Consent Considerations:** Patient notification requirements for secondary data use.

# Ethical Considerations in Clinical ML

Responsible development requires proactive attention to fairness:

- **Algorithmic Fairness:** Assessment of model performance across demographic groups.
- **Bias Detection:** Analysis of training data for historical biases.
- **Transparency:** Explainability features and clear communication about limitations.
- **Accountability:** Clear lines of responsibility for model-related decisions.

# Regulatory Framework for Clinical AI

Clinical ML systems may be subject to regulatory oversight:

- **SaMD Classification:** Software as a Medical Device framework from FDA, EU MDR.
- **Risk-Based Validation:** Validation requirements scaled to risk level.
- **Quality Management Systems:** FDA 21 CFR Part 820, ISO 13485.
- **Post-Market Surveillance:** Ongoing monitoring and adverse event reporting.

## Example 1: Sepsis Early Warning System

Sepsis prediction demonstrates the full workflow from problem definition to deployment:

- **Problem:** Early identification of patients at risk for timely antibiotic administration.
- **Data:** Vital signs, laboratory values (lactate, WBC, creatinine), clinical notes.
- **Features:** SOFA score components, infection indicators, temporal trends.
- **Model:** Gradient boosting classifier with interpretability enhancements.
- **Deployment:** Integrated into nursing station dashboards with tiered alerts.

Clinical studies have demonstrated reduced mortality with well-designed systems.



## Example 2: Radiology Image Triage

Computer vision applications in radiology demonstrate ML potential:

- **Problem:** Prioritization of chest radiographs with acute findings (pneumothorax, pneumonia).
- **Data:** Historical chest radiographs with radiologist annotations.
- **Model:** Convolutional neural networks (CNN) for image classification.
- **Deployment:** Integration with PACS systems to flag studies for expedited review.

These systems augment radiologist workflows without replacing human interpretation.

## Q4 Summary: Key Takeaways

Successful implementation of ML in hospital settings requires attention to technical, organizational, and ethical dimensions:

- **Clinical Alignment:** ML solutions must address genuine clinical needs with actionable predictions.
- **Data Foundation:** Robust data quality and preprocessing are prerequisites.
- **Validation Rigor:** Comprehensive validation across temporal, external, and subgroup dimensions.
- **Human-in-the-Loop:** Clinical ML should augment rather than replace human judgment.

### Final Message

Machine learning holds tremendous promise for improving healthcare outcomes, but realizing this potential requires responsible development, rigorous validation, and sustained commitment to patient safety and ethical principles.

## References

---

## References: Q1 - MLE and Logistic Regression

- Hosmer, D.W., Lemeshow, S., & Sturdivant, R.X. (2013). Applied Logistic Regression (3rd ed.). Wiley.
- Agresti, A. (2018). An Introduction to Categorical Data Analysis (3rd ed.). Wiley.
- McCullagh, P. & Nelder, J.A. (1989). Generalized Linear Models (2nd ed.). Chapman & Hall.
- Rothman, K.J., Greenland, S., & Lash, T.L. (2008). Modern Epidemiology (3rd ed.). Lippincott Williams & Wilkins.

## References: Q2 - Cox Proportional Hazards Model

- Cox, D.R. (1972). Regression Models and Life-Tables (with discussion). Journal of the Royal Statistical Society, Series B, 34(2), 187-220.
- Kleinbaum, D.G. and Klein, M. (2012). Survival Analysis: A Self-Learning Text (3rd ed.). Springer.
- Therneau, T.M. and Grambsch, P.M. (2000). Modeling Survival Data: Extending the Cox Model. Springer.
- Hosmer, D.W., Lemeshow, S., and May, S. (2008). Applied Survival Analysis: Regression Modeling of Time to Event Data (2nd ed.). Wiley.

## References: Q3 - PCA

- Jolliffe, I.T. and Cadima, J. (2016). Principal component analysis: A review and recent developments. *Philosophical Transactions of the Royal Society A*, 374(2065), 20150202.
- Ringnér, M. (2008). What is principal component analysis? *Nature Biotechnology*, 26(3), 303-304.

## References: Q4 - ML Workflow

- Hastie, T., Tibshirani, R., and Friedman, J. (2009). The Elements of Statistical Learning (2nd ed.). Springer.
- Rajkomar, A., Oren, E., Chen, K., et al. (2018). Scalable and accurate deep learning with electronic health records. NPJ Digital Medicine, 1(1), 18.
- FDA. (2019). Software as a Medical Device (SaMD): Clinical Evaluation. FDA Guidance Document.