# Understanding Risk Factors for Mortality Among Older Individuals

(A Comprehensive Analysis using Classical, Ensemble and Deep Learning Models)

Srimanta Ghosh

University Roll No.: E1773U234013

**Supervisors:**

Dr. Rahul Ghosal[1]

Dr. Sudhakar Sahoo[2]

[1]Assistant Professor, Arnold School of Public Health
University of South Carolina

[2]Associate Professor
Institute of Mathematics and Applications, Bhubaneswar

# Outline

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## Outline

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## The Problem

- As populations around the world continue to age, understanding the determinants of healthy aging and longevity becomes increasingly vital.

- While physical activity has been shown to reduce mortality risk, the complex interplay between lifestyle behaviors, demographic characteristics, and health conditions remains poorly understood.

- Leveraging advancements in statistical modeling and machine learning, this project aims to create a comprehensive framework for predicting mortality risk—enabling more precise identification of high-risk individuals and informing public health strategies for aging populations.

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## The Problem

- As populations around the world continue to age, understanding the determinants of healthy aging and longevity becomes increasingly vital.

- While physical activity has been shown to reduce mortality risk, the complex interplay between lifestyle behaviors, demographic characteristics, and health conditions remains poorly understood.

- Leveraging advancements in statistical modeling and machine learning, this project aims to create a comprehensive framework for predicting mortality risk—enabling more precise identification of high-risk individuals and informing public health strategies for aging populations.

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## The Problem

- As populations around the world continue to age, understanding the determinants of healthy aging and longevity becomes increasingly vital.

- While physical activity has been shown to reduce mortality risk, the complex interplay between lifestyle behaviors, demographic characteristics, and health conditions remains poorly understood.

- Leveraging advancements in statistical modeling and machine learning, this project aims to create a comprehensive framework for predicting mortality risk—enabling more precise identification of high-risk individuals and informing public health strategies for aging populations.

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

# Outline

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

Research Questions

- How are physical activity, demographic factors, and health indicators related to mortality in older adults?

- What are the key risk factors that contribute to mortality in older adults?

- How do classical statistical models, ensemble methods, and deep learning approaches compare in their effectiveness at predicting mortality risk?

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## Research Questions

- How are physical activity, demographic factors, and health indicators related to mortality in older adults?

- What are the key risk factors that contribute to mortality in older adults?

- How do classical statistical models, ensemble methods, and deep learning approaches compare in their effectiveness at predicting mortality risk?

Introduction
The Data
The Models
Conclusion and Future Work

The Problem
Research Questions

## Research Questions

- How are physical activity, demographic factors, and health indicators related to mortality in older adults?

- What are the key risk factors that contribute to mortality in older adults?

- How do classical statistical models, ensemble methods, and deep learning approaches compare in their effectiveness at predicting mortality risk?

Introduction
The Data
The Models
Conclusion and Future Work

Data Source
The Survival Data
Preprocessing (Physical Activity Data)

# Outline

Introduction
**The Data**
The Models
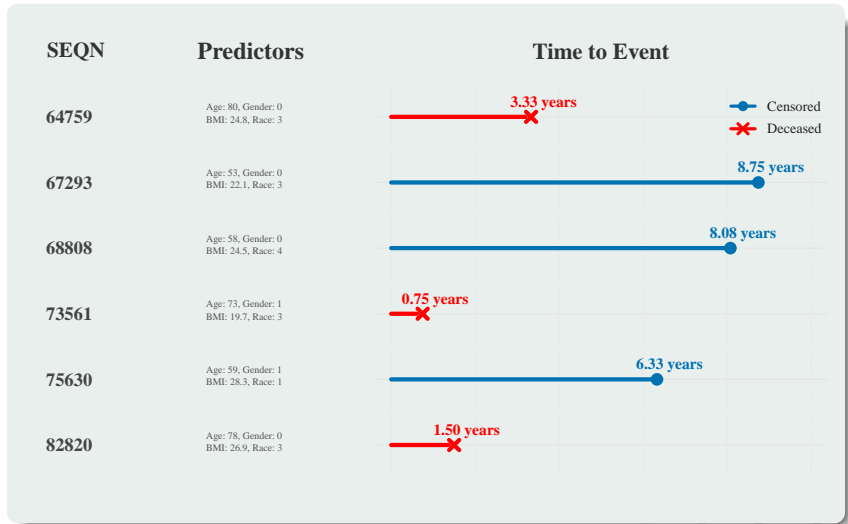Conclusion and Future Work

Data Source
The Survival Data
Preprocessing (Physical Activity Data)

## Data Source

The study will use data from the National Health and Nutrition Examination Survey (NHANES) 2011-2014, "*https://wwwn.cdc.gov/nchs/nhanes/Default.aspx*", which includes demographic, lifestyle, and health-related variables. The mortality information is linked to the National Death Index (NDI).

Introduction
**The Data**
The Models
Conclusion and Future Work

Data Source
**The Survival Data**
Preprocessing (Physical Activity Data)

## Outline

Introduction
**The Data**
The Models
Conclusion and Future Work

Data Source
The Survival Data
Preprocessing (Physical Activity Data)

# The Survival Data



| SEQN | Predictors | Time to Event |
|------|-----------|---------------|

| 64759 | Age: 80, Gender: 0<br>BMI: 24.8, Race: 3 | 3.33 years — Deceased |
| 67293 | Age: 53, Gender: 0<br>BMI: 22.1, Race: 3 | 8.75 years — Censored |
| 68808 | Age: 58, Gender: 0<br>BMI: 24.5, Race: 4 | 8.08 years |
| 73561 | Age: 73, Gender: 1<br>BMI: 19.7, Race: 3 | 0.75 years |
| 75630 | Age: 59, Gender: 1<br>BMI: 28.3, Race: 1 | 6.33 years |
| 82820 | Age: 78, Gender: 0<br>BMI: 26.9, Race: 3 | 1.50 years |

Legend: ● — Censored ✕ — Deceased

Introduction
**The Data**
The Models
Conclusion and Future Work

Data Source
The Survival Data
**Preprocessing (Physical Activity Data)**

# Outline

Introduction
**The Data**
The Models
Conclusion and Future Work

Data Source
The Survival Data
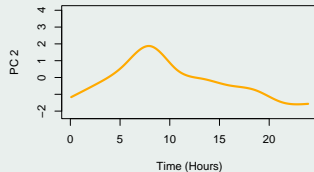**Preprocessing (Physical Activity Data)**

## Preprocessing

### FPCA

$$X_i(t) = \mu(t) + \sum_{k=1}^{K} \xi_{ik}\phi_k(t) + \epsilon_i(t)$$



**Mean Physical Acitivity**



**Morning–Active vs. Evening–Active**

Introduction
The Data
**The Models**
Conclusion and Future Work

**Structures**
Output
C-Index Comparison

# Outline

Introduction
The Data
**The Models**
Conclusion and Future Work

**Structures**
Output
C-Index Comparison

## Models

| Model | Type | Loss Function |
|-------|------|---------------|
| Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right) + \lambda \|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f, \beta) = -\sum_{i:\delta_i=1} \left( \eta_i - \log \sum_{j \in R(T_i)} \exp(\eta_i) \right) + \lambda \sum_{j=1}^{p} \int \left( f_j''(x) \right)^2 dx$ where $\eta_i = \sum_{j=1}^{p} f_j(X_{ij}) + \sum_{k=1}^{q} \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E} \sum_{i:\delta_i=1} \left[ h_\theta(x_i) - \log \sum_{j \in R(T_i)} \exp(h_\theta(x_j)) \right] + \lambda \|\theta\|_2^2$ |

| Model | Type | Split Rule |
|-------|------|-----------|
| RSF | Ensemble | Log-rank test statistic: $L(x, c) = \sum_{i=1}^{N} \dfrac{\left( d_{i,1} - \frac{Y_{i,1} d_i}{Y_i} \right)}{\sqrt{\sum_{i=1}^{N} \frac{Y_{i,1}}{Y_i} \left( 1 - \frac{Y_{i,1}}{Y_i} \right) \frac{(Y_i - d_i)}{(Y_i - 1)} d_i}}$ |

Introduction
The Data
The Models
Conclusion and Future Work

Structures
Output
C-Index Comparison

## Models

| Model | Type | Loss Function |
|-------|------|---------------|
| **Cox** | **Classical** | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right) + \lambda \|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f,\beta) = -\sum_{i:\delta_i=1} \left( \eta_i - \log \sum_{j \in R(T_i)} \exp(\eta_j) \right) + \lambda \sum_{j=1}^p \int \left( f_j''(x) \right)^2 dx$ where $\eta_i = \sum_{j=1}^p f_j(X_{ij}) + \sum_{k=1}^q \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E} \sum_{i:\delta_i=1} \left[ h_\theta(x_i) - \log \sum_{j \in R(T_i)} \exp(h_\theta(x_j)) \right] + \lambda \|\theta\|_2^2$ |

| Model | Type | Split Rule |
|-------|------|-----------|
| RSF | Ensemble | Log-rank test statistic: $L(x,c) = \sum_{i=1}^N \dfrac{\left( d_{i,1} - \frac{Y_{i,1} d_i}{Y_i} \right)}{\sqrt{\sum_{i=1}^N \frac{Y_{i,1}}{Y_i} \left( 1 - \frac{Y_{i,1}}{Y_i} \right) \frac{(Y_i - d_i)}{(Y_i - 1)} d_i}}$ |

Introduction
The Data
The Models
Conclusion and Future Work

Structures
Output
C-Index Comparison

## Models

| Model | Type | Loss Function |
|-------|------|---------------|
| Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right) + \lambda \|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f, \beta) = -\sum_{i:\delta_i=1} \left( \eta_i - \log \sum_{j \in R(T_i)} \exp(\eta_j) \right) + \lambda \sum_{j=1}^{p} \int \left( f_j''(x) \right)^2 dx$ where $\eta_i = \sum_{j=1}^{p} f_j(X_{ij}) + \sum_{k=1}^{q} \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E} \sum_{i:\delta_i=1} \left[ h_\theta(x_i) - \log \sum_{j \in R(T_i)} \exp(h_\theta(x_j)) \right] + \lambda \|\theta\|_2^2$ |

| Model | Type | Split Rule |
|-------|------|------------|
| RSF | Ensemble | Log-rank test statistic: $L(x, c) = \sum_{i=1}^{N} \dfrac{\left( d_{i,1} - \frac{Y_{i,1} d_i}{Y_i} \right)}{\sqrt{\sum_{i=1}^{N} \frac{Y_{i,1}}{Y_i} \left( 1 - \frac{Y_{i,1}}{Y_i} \right) \frac{(Y_i - d_i)}{(Y_i - 1)} d_i}}$ |

Introduction
The Data
The Models
Conclusion and Future Work

Structures
Output
C-Index Comparison

## Models

| Model | Type | Loss Function |
|-------|------|---------------|
| Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right) + \lambda \|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f, \beta) = -\sum_{i:\delta_i=1} \left( \eta_i - \log \sum_{j \in R(T_i)} \exp(\eta_i) \right) + \lambda \sum_{j=1}^{p} \int \left( f_j''(x) \right)^2 dx$ where $\eta_i = \sum_{j=1}^{p} f_j(X_{ij}) + \sum_{k=1}^{q} \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E} \sum_{i:\delta_i=1} \left[ h_\theta(x_i) - \log \sum_{j \in R(T_i)} \exp(h_\theta(x_j)) \right] + \lambda \|\theta\|_2^2$ |

| Model | Type | Split Rule |
|-------|------|------------|
| RSF | Ensemble | Log-rank test statistic: $L(x, c) = \sum_{i=1}^{N} \dfrac{\left( d_{i,1} - \frac{Y_{i,1} d_i}{Y_i} \right)}{\sqrt{\sum_{i=1}^{N} \frac{Y_{i,1}}{Y_i} \left( 1 - \frac{Y_{i,1}}{Y_i} \right) \frac{(Y_i - d_i)}{(Y_i - 1)} d_i}}$ |

Introduction
The Data
**The Models**
Conclusion and Future Work

**Structures**
Output
C-Index Comparison

## Models

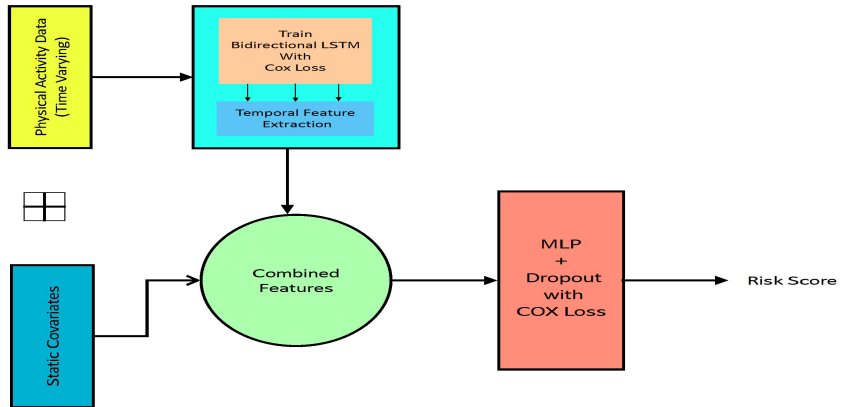| Model | Type | Loss Function |
|-------|------|---------------|
| Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1} \left( x_i^\top \beta - \log \sum_{j \in R(T_i)} \exp(x_j^\top \beta) \right) + \lambda \|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f, \beta) = -\sum_{i:\delta_i=1} \left( \eta_i - \log \sum_{j \in R(T_i)} \exp(\eta_i) \right) + \lambda \sum_{j=1}^{p} \int \left( f_j''(x) \right)^2 dx$ where $\eta_i = \sum_{j=1}^{p} f_j(X_{ij}) + \sum_{k=1}^{q} \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E} \sum_{i:\delta_i=1} \left[ h_\theta(x_i) - \log \sum_{j \in R(T_i)} \exp(h_\theta(x_j)) \right] + \lambda \|\theta\|_2^2$ |

| Model | Type | Split Rule |
|-------|------|------------|
| RSF | Ensemble | Log-rank test statistic: $L(x, c) = \sum_{i=1}^{N} \dfrac{\left( d_{i,1} - \frac{Y_{i,1} d_i}{Y_i} \right)}{\sqrt{\sum_{i=1}^{N} \frac{Y_{i,1}}{Y_i} \left( 1 - \frac{Y_{i,1}}{Y_i} \right) \frac{(Y_i - d_i)}{(Y_i - 1)} d_i}}$ |

Introduction
The Data
The Models
Conclusion and Future Work

Structures
Output
C-Index Comparison

## Models

| Model | Type | Loss Function |
|---|---|---|
| Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1}\left(x_i^\top \beta - \log\sum_{j\in R(T_i)}\exp(x_j^\top\beta)\right)$ |
| Penalized Cox | Classical | $\mathcal{L}(\beta) = -\sum_{i:\delta_i=1}\left(x_i^\top \beta - \log\sum_{j\in R(T_i)}\exp(x_j^\top\beta)\right) + \lambda\|\beta\|_1$ |
| GAM Cox | Classical | $\mathcal{L}(f,\beta) = -\sum_{i:\delta_i=1}\left(\eta_i - \log\sum_{j\in R(T_i)}\exp(\eta_i)\right) + \lambda\sum_{j=1}^p\int\left(f_j''(x)\right)^2 dx$ <br> where $\eta_i = \sum_{j=1}^p f_j(X_{ij}) + \sum_{k=1}^q \beta_k Z_{ik}$ |
| DeepSurv | Deep Learning | $\mathcal{L}(\theta) = -\frac{1}{N_E}\sum_{i:\delta_i=1}\left[h_\theta(x_i) - \log\sum_{j\in R(T_i)}\exp(h_\theta(x_j))\right] + \lambda\|\theta\|_2^2$ |

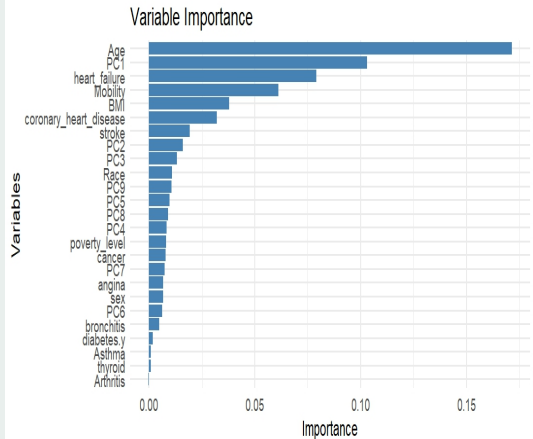| Model | Type | Split Rule |
|---|---|---|
| RSF | Ensemble | Log-rank test statistic: <br> $L(x,c) = \sum_{i=1}^N \dfrac{\left(d_{i,1} - \frac{Y_{i,1}d_i}{Y_i}\right)}{\sqrt{\sum_{i=1}^N \frac{Y_{i,1}}{Y_i}\left(1 - \frac{Y_{i,1}}{Y_i}\right)\frac{(Y_i - d_i)}{(Y_i - 1)}d_i}}$ |

Introduction
The Data
The Models
Conclusion and Future Work

Structures
Output
C-Index Comparison

# BiLSTM-Deepsurv Architecture

Introduction
The Data
**The Models**
Conclusion and Future Work

Structures
**Output**
C-Index Comparison

# Outline

Introduction
The Data
**The Models**
Conclusion and Future Work

Structures
**Output**
C-Index Comparison

## Output

| Variable | exp(coef.) | Pr(>|z|) |
|----------|-----------|----------|
| Age | 1.06 | < 2e-16 |
| BMI | 0.96 | 4.67e-07 |
| Mobility2 | 0.58 | 2.03e-08 |
| povertylevel | 0.92 | 0.00517 |
| heartfailure1 | 1.87 | 7.23e-08 |
| PC1 | 0.86 | 3.04e-14 |

Introduction
The Data
**The Models**
Conclusion and Future Work

Structures
Output
C-Index Comparison

# C-Index Comparison

$$C = \frac{\sum \mathbf{1}(h(X_i) > h(X_j)) \cdot \mathbf{1}(T_i < T_j)}{\sum \mathbf{1}(T_i < T_j)}$$

where $h(X)$ is the model's predicted risk score, and $T$ is the observed survival time.

| Model | Average C-Index |
|-------|----------------|
| Cox | 0.77421 |
| GAM Cox | 0.077763 |
| Penalized Cox | 0.77424 |
| RSF | 0.79895 |
| Deepsurv | 0.78173 |
| BiLSTM Deepsurv | 0.76041 |

## Conclusion and Future Work

- If your primary goal is achieving high accuracy, Random Survival Forest (RSF) is the best-performing model.

- If your primary goal is interpretability, the Generalized Additive Cox model is one of the best reliable choices.

- DeepSurv is a strong deep learning-based alternative that shows good result in prediction.

- The BiLSTM-DeepSurv model appears promising, but it still requires improvement. With more computational resources and further tuning, it may perform better in future experiments.

## Conclusion and Future Work

- If your primary goal is achieving high accuracy, Random Survival Forest (RSF) is the best-performing model.

- If your primary goal is interpretability, the Generalized Additive Cox model is one of the best reliable choices.

- DeepSurv is a strong deep learning-based alternative that shows good result in prediction.

- The BiLSTM-DeepSurv model appears promising, but it still requires improvement. With more computational resources and further tuning, it may perform better in future experiments.

## Conclusion and Future Work

- If your primary goal is achieving high accuracy, Random Survival Forest (RSF) is the best-performing model.
- If your primary goal is interpretability, the Generalized Additive Cox model is one of the best reliable choices.
- DeepSurv is a strong deep learning-based alternative that shows good result in prediction.
- The BiLSTM-DeepSurv model appears promising, but it still requires improvement. With more computational resources and further tuning, it may perform better in future experiments.

## Conclusion and Future Work

- If your primary goal is achieving high accuracy, Random Survival Forest (RSF) is the best-performing model.
- If your primary goal is interpretability, the Generalized Additive Cox model is one of the best reliable choices.
- DeepSurv is a strong deep learning-based alternative that shows good result in prediction.
- The BiLSTM-DeepSurv model appears promising, but it still requires improvement. With more computational resources and further tuning, it may perform better in future experiments.

# References

📄 Rahul Ghosal, Marcos Matabuena, and Sujit K Ghosh.
Functional time transformation model with applications to digital health.
*Computational Statistics & Data Analysis*, page 108131, 2025.

📄 Hemant Ishwaran and Udaya B Kogalur.
Random survival forests for r.
*R news*, 7(2):25–31, 2007.

📄 Jared L Katzman, Uri Shaham, Alexander Cloninger, Jonathan Bates, Tingting Jiang, and Yuval Kluger.
Deepsurv: personalized treatment recommender system using a cox proportional hazards deep neural network.
*BMC medical research methodology*, 18:1–12, 2018.