# Establishing Physical Activity Guidelines

An Application of Shape Constrained Regression

Eli Kravitz

2018

## Introduction

- This talk is based in part on Development and Testing of an Integrated Score for Physical Behaviors by Keadle et al.

- Scan QR code to take you to the paper



**SCAN ME**

## Motivation

- Common in epidemiology to take complicated multidimensional behavior and reduce it to a single interpretable number between 0 and 100

    - often called a composite score or an index
    - Ex: Healthy Eating Index

- Use this number to predict risk of disease

- Currently, no scoring system exist for physical activity.

# Healthy Eating Index (HEI)

| Component | Maximum points | Standard for maximum score | Standard for minimum score of zero |
|---|---|---|---|
| Total Fruit (includes 100% juice) | 5 | ≥0.8 cup equiv. per 1,000 kcal | No Fruit |
| Whole Fruit (not juice) | 5 | ≥0.4 cup equiv. per 1,000 kcal | No Whole Fruit |
| Total Vegetables | 5 | ≥1.1 cup equiv. per 1,000 kcal | No Vegetables |
| Dark Green and Orange Vegetables and Legumes[2] | 5 | ≥0.4 cup equiv. per 1,000 kcal | No Dark Green or Orange Vegetables or Legumes |
| Total Grains | 5 | ≥3.0 oz equiv. per 1,000 kcal | No Grains |
| Whole Grains | 5 | ≥1.5 oz equiv. per 1,000 kcal | No Whole Grains |
| Milk[3] | 10 | ≥1.3 cup equiv. per 1,000 kcal | No Milk |
| Meat and Beans | 10 | ≥2.5 oz equiv. per 1,000 kcal | No Meat or Beans |
| Oils[4] | 10 | ≥12 grams per 1,000 kcal | No Oil |
| Saturated Fat | 10 | ≤7% of energy[5] | ≥15% of energy |
| Sodium | 10 | ≤0.7 gram per 1,000 kcal[5] | ≥2.0 grams per 1,000 kcal |
| Calories from Solid Fats, Alcoholic beverages, and Added Sugars (SoFAAS) | 20 | ≤20% of energy | ≥50% of energy |

[1]Intakes between the minimum and maximum levels are scored proportionately, except for Saturated Fat and Sodium (see note 5).
[2]Legumes counted as vegetables only after Meat and Beans standard is met.
[3]Includes all milk products, such as fluid milk, yogurt, and cheese, and soy beverages.
[4]Includes nonhydrogenated vegetable oils and oils in fish, nuts, and seeds.
[5]Saturated Fat and Sodium get a score of 8 for the intake levels that reflect the 2005 Dietary Guidelines, <10% of calories from saturated fat and 1.1 grams of sodium/1,000 kcal, respectively.

## Physical Activtiy Data

- We work with the NIH-AARP Study of Diet and Health

  - Adults between 50-71 year old

  - The self-report questionnaire asked how much time per week was spent in 16 different physical behaviors during the past 12 months.

  - The physical behaviors were categorized into 5 exercise or sport activities, two sitting behaviors and sleep.

## Physical Activity Components

- Activity is broken down as:
  1. Light household activity: cooking, cleaning, laundry, dusting
  2. Moderate-vigorous household activity:vacuuming, sweeping weeding, raking, home repairs painting
  3. Moderate activity: walking for exercise, walking for other daily activities, playing golf
  4. Vigorous Exercise: playing tennis, swimming laps, bicycling, jogging
  5. Weight training
  6. Sitting watching television
  7. Other sitting: reading, knitting, using a computer
  8. Sleeping: at night or napping during the day

## Inital Model

- How do the 8 physical activity components and additional covariates relate to overall health?

- Use a Generalized Additive Model (GAM) with **survival** as the outcomes of interest:

$$Pr(Y = 1|X, Z) = H\{\sum_{j=1}^{d} f_j(X_{ij}) + Z^T\theta\},$$

where $H(\cdot)$ is the logistic distribution function, $X_{ij}$ are physical activity measurements, $f_j(\cdot)$ are unspecified smooth functions, and $Z_i$ are covariates,

- Pros: Flexible model, doesn't require any parametric assumptions, no one has done this in the physical activity literature.

- Our collaborator gave us some additional information, some of which is common sense and some is unexpected:

    - Activity is always good for you but the effect levels off
    - No sleep is bad for you, sleep is good for your health until some change point which indicates poor health
    - Sitting is bad for your health but fine in moderation

- We can use this information to apply shape constraints to the individual functions, $f_j$, in the GAM.

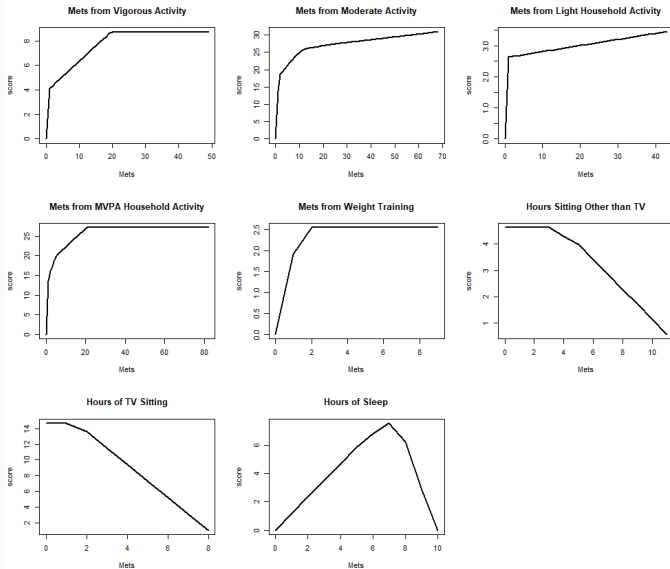- Two methods for doing this: SCAR (Yining Chen and Richard Samworth) and SCAM (Natalya Pya and Simon Wood)
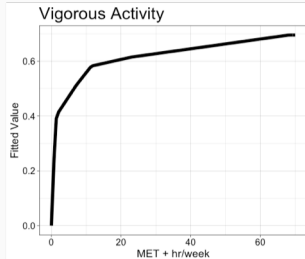
## Shape Constraints

| Activity | Constraint |
|---|---|
| Vigorous Activity | Concave Increasing |
| Moderate Activity | Concave Increasing |
| Light Household Activity | Concave Increasing |
| MVPA Household Activity | Concave Increasing |
| Weight Training | Concave Increasing |
| Hours Sitting Other than TV | Concave Decreasing |
| Hours of TV Sitting | Concave Decreasing |
| Hours of Sleep | Concave |

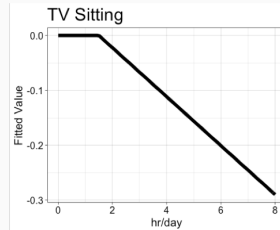Each constraint is justified in the physical activity literature and make sense to our collaborators

# Shape Constraints

# Shape Constraints



**Figure 1: Aerobic Activity**
Concave, increasing



**Figure 2:**
**Sedentary Behavior**:Concave,
decreasing

## Making a composite score

- Want our results to be on same scale as other composite scores

- We proportionally rescale the physical activity fitted values, $\sum_j f_j(x_{ij})$ to take values between 0 and 100.

## Composite Score

| Component | Contribution to Total | Criteria for Maximum |
|---|:---:|:---:|
| Vigorous Activity | 10 | >20 MET-hrs/wk |
| Moderate Activity | 32 | >50 MET-hrs/wk |
| Light Household Activity | 3 | >3 MET-hrs/wk |
| MVPA Household Activity | 25 | >20 MET-hrs/wk |
| Weight Training | 3 | >2 MET-hrs/wk |
| Sitting Other than TV | 5 | 0 hours |
| Hours of TV Sitting | 14 | 0 hours |
| Hours of Sleep | 8 | 7.5 hours |
| Total | 100 | |

## Validation Concerns

- We've built a 0-to-100-score. Now we ask is our score meaningful?

    - Is a higher score (more physical activity) predictive of longer survival time?

- Ideal Scenario: Build a score for one population $\rightarrow$ validate on a different population

- Validation with a second dataset is not possible.

## Validation Concerns

- Denote the new physical activity scores as $\mathcal{S}$, the covariates as $Z$, and time of death as $t$:

$$h(t) = h_0(t)exp(\beta\mathcal{S} + Z^T\theta)$$

- Is $\beta$ estimated correctly? Is the variability of $\hat{\beta}$ estimated correctly?

- The hypothesis test $H_0 : \beta = 0$ is very important, indicates if $\mathcal{S}$ is useful.

  - Do not want to incorrectly reject $H_0 \rightarrow$ conclude $\mathcal{S}$ is predictive of mortality when it really is not.

## Sample Splitting

- Possible Solution: Split the data in half into a "training" set and "test" set.
    - Can be highly variable, a single "lucky" split can cause us to underestimate standard error and improperly estimate coefficient (simulation error)
- Idea: Repeatedly split the data in half into training and test. Aggregate all the estimates
    - Explored by Meinshausen et. all (2008).
    - Variants of this technique are used in high dimensional variable selection literature, Wasserman and Roeder (2008), Lei et. all (2016)

This should reduce simulation error.

## Sample Splitting Procedure

1. Split data in half, $b = 1, \ldots, B$ times to create $B$ partitions of the data. Denote these partitions as $D_{in}^{(b)}$ and $D_{out}^{(b)}$

2. Fit shape constrained GAM on $D_{in}^{(b)}$.

3. Use the fitted value of the GAM to assign a score, $\mathcal{S}$ to every person in $D_{out}^{(b)}$. Run a cox regression using $D_{out}^{(b)}$ to get $\widehat{\beta}^{(b)}$ and p-value, $P^{(b)}$

## New Estimates

- New estimate of $\beta$: $\widehat{\beta} = B^{-1} \sum_{b=1}^{B} \widehat{\beta}^{(b)}$

- Standard error: sample standard deviation of $\{\widehat{\beta}^{(b)} | b = 1, \ldots, B\}$

- New p-value is more complicated:

    - Define $Q(\gamma) = q_\gamma(\{P^{(b)}/\gamma; b = 1, \ldots, B\})$, where $q_\gamma(\cdot)$ is the $\gamma^{th}$ quantile function
    - Search for the best quantile:
      $P = min\{1, (1 - log\gamma_{min}) \inf_{\gamma \in (\gamma_{min}, 1)} Q(\gamma)\}$

## New Estimates

- Our estimate, $\widehat{\beta}$ is consistent and the Type-I error is controled.

- Justification is messy and technical. Involves a lot of frequentist asymptotics.

- Excluded it from this talk, talk to me after if you want more information.

## Results

|  | coef | exp(coef) | se(coef) | Conf Int. | p-value |
|---|---|---|---|---|---|
| PA Score Quintile 1 | · | · | · | · | |
| PA Score Quintile 2 | -0.27 | 0.76 | 0.021 | (0.73, 0.80) | 0 |
| PA Score Quintile 3 | -0.39 | 0.68 | 0.023 | (0.65, 0.71) | 0 |
| PA Score Quintile 4 | -0.51 | 0.60 | 0.024 | (0.57, 0.63) | 0 |
| PA Score Quintile 5 | -0.61 | 0.54 | 0.027 | (0.52, 0.57) | 0 |