



Feature Engineering

Data Science & Machine Learning Team

02/22/2021

Andrew Wheeler, PhD

andrew.wheeler@hms.com

Agenda

- Machine Learning and Causality
 - Importance of business domain knowledge
- Functional form
 - Polynomial terms, non-linear effects, step functions
 - Encoding categorical variables
- Examples in Python using linear regression

Machine Learning and Causality

- To get a better prediction, need a basic understanding of the causal mechanisms behind the phenomenon
 - Need to feed the machine the correct data in the correct format, or it will not generate valid predictions in practice.
- Example business problems at HMS
 - Subrogation – younger people are more likely to have car accidents
 - Payment Integrity – some insurance claim types are more discretionary, more likely to be upcoded
 - Technical Denial – some claims have higher complexity, and so are likely to be missing critical information
- Each requires unique solutions – same model architecture would not work for all three projects

What is Feature Engineering?

- A model has inputs used to predict an output, e.g.:

$$y = f(a, b, c)$$

- Feature engineering involves:

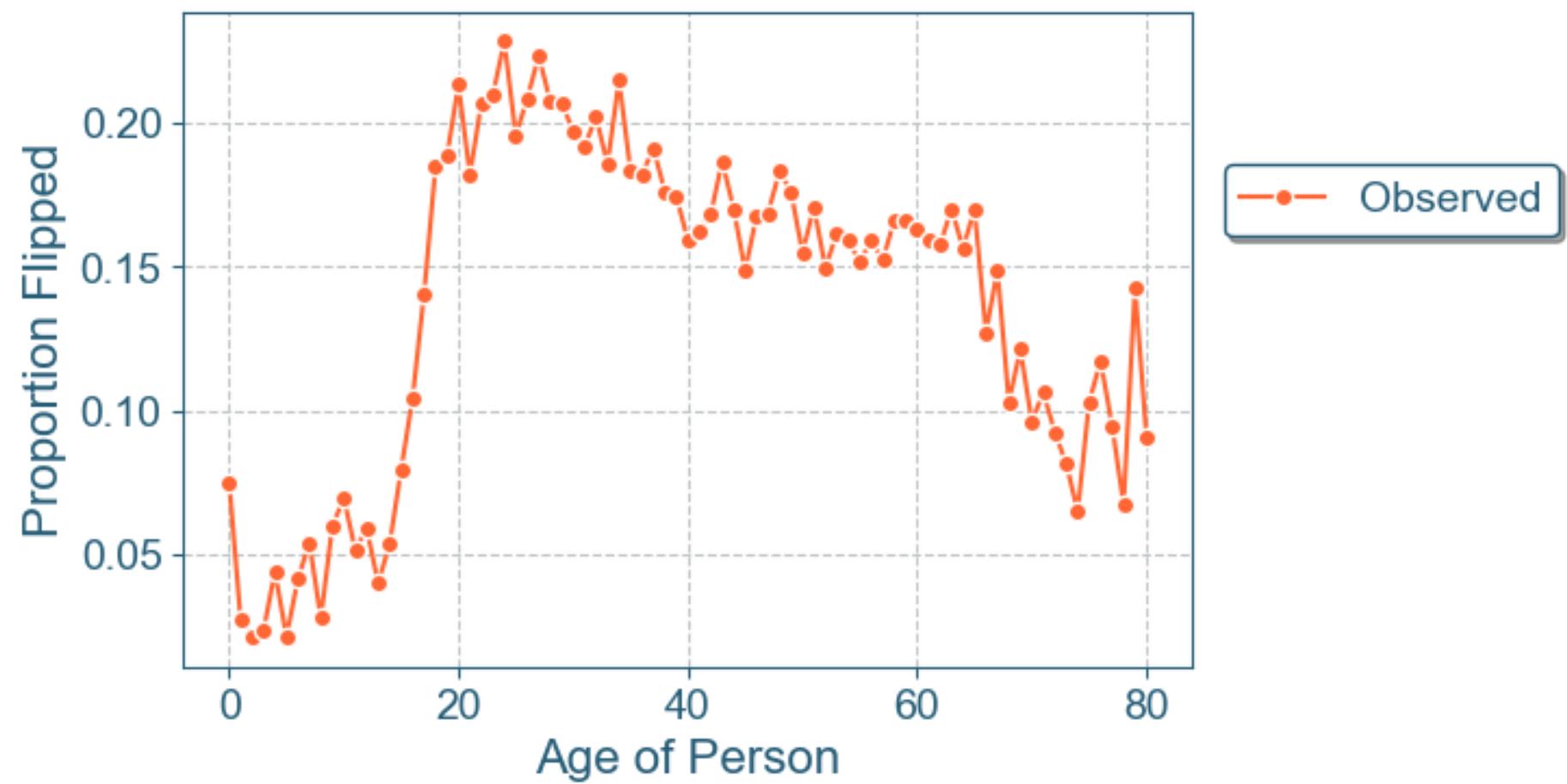
- What features to include in the model, e.g. (a,b,c) instead of (a,e,f)?
- Transformations of variables, e.g. $\log(y) = \beta_1 \cdot \sqrt{a}$
- Representing categorical variables in a model, e.g.

$$y = \beta_1(b = \text{Aetna}) + \beta_2(b = \text{Amerihealth})$$

- Combinations of all of these variables together, e.g. interaction effects

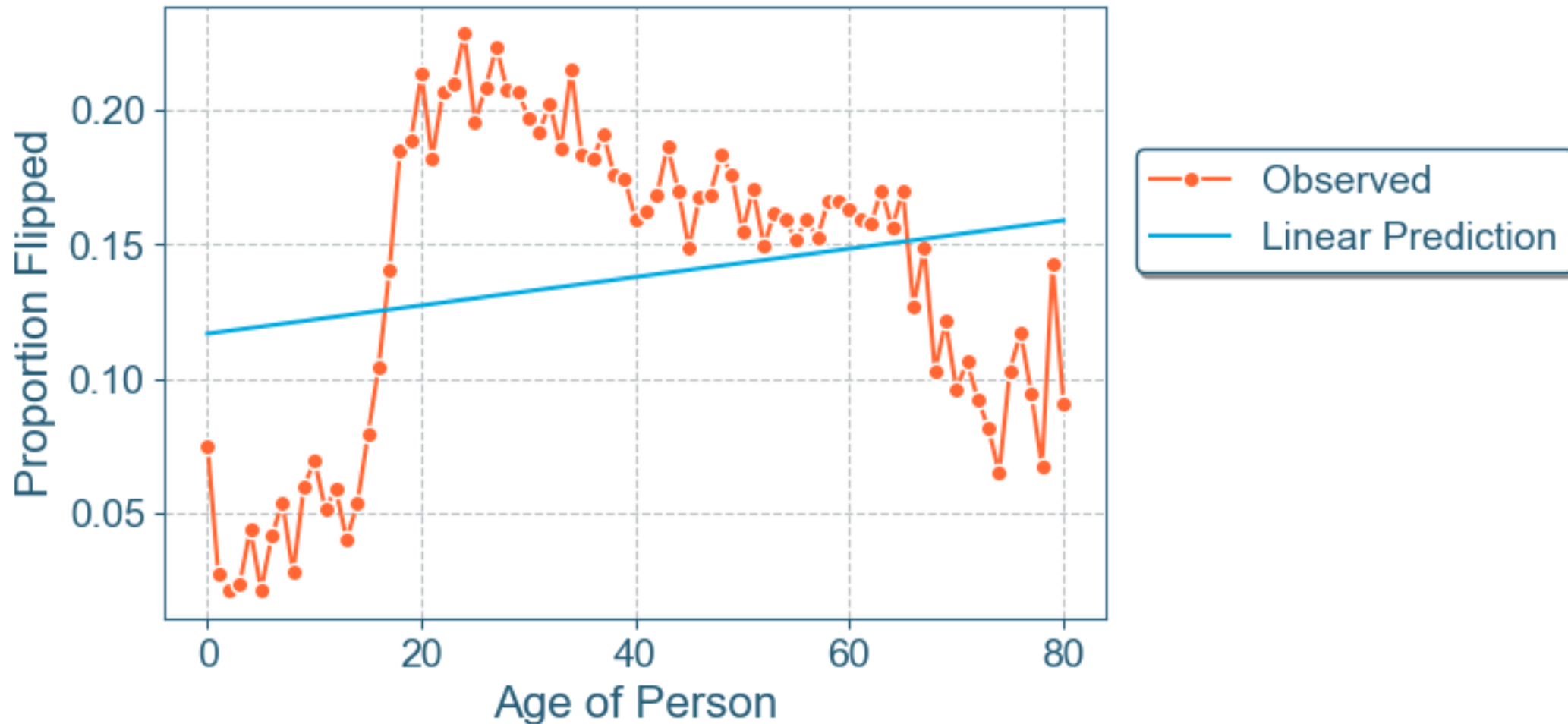
$$y = \beta_1\sqrt{a} + \beta_2(b = \text{Aetna}) + \beta_3(b = \text{Amerihealth}) + \beta_4(\sqrt{a} \cdot [b = \text{Aetna}])$$

Example: Subrogation (via Accent)

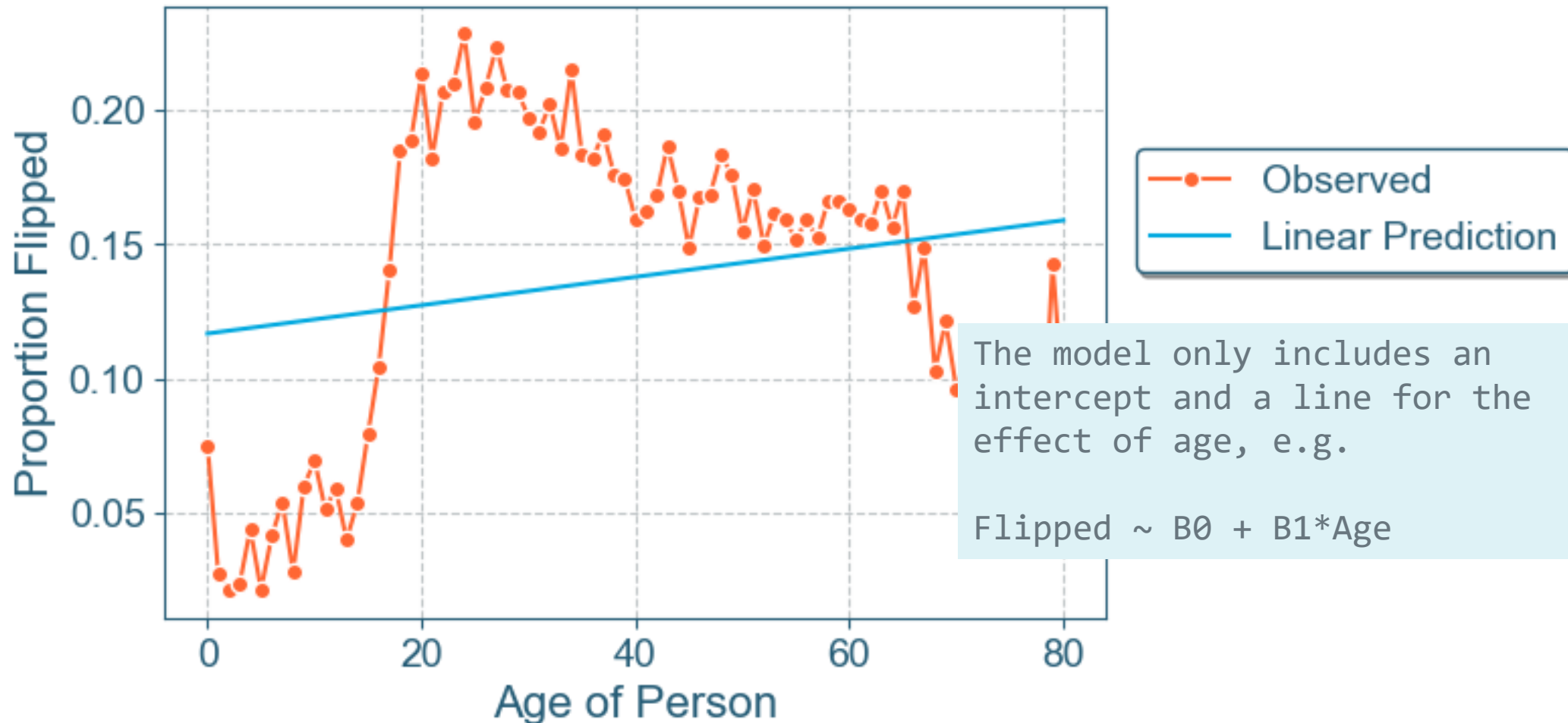


	Age	PercentFlipped
Index		
0	0	0.075314
1	1	0.027778
2	2	0.021786
3	3	0.023641
4	4	0.044118
...
76	76	0.117188
77	77	0.094862
78	78	0.067568
79	79	0.142857
80	80	0.090909

Example: Subrogation (via Accent) – Linear

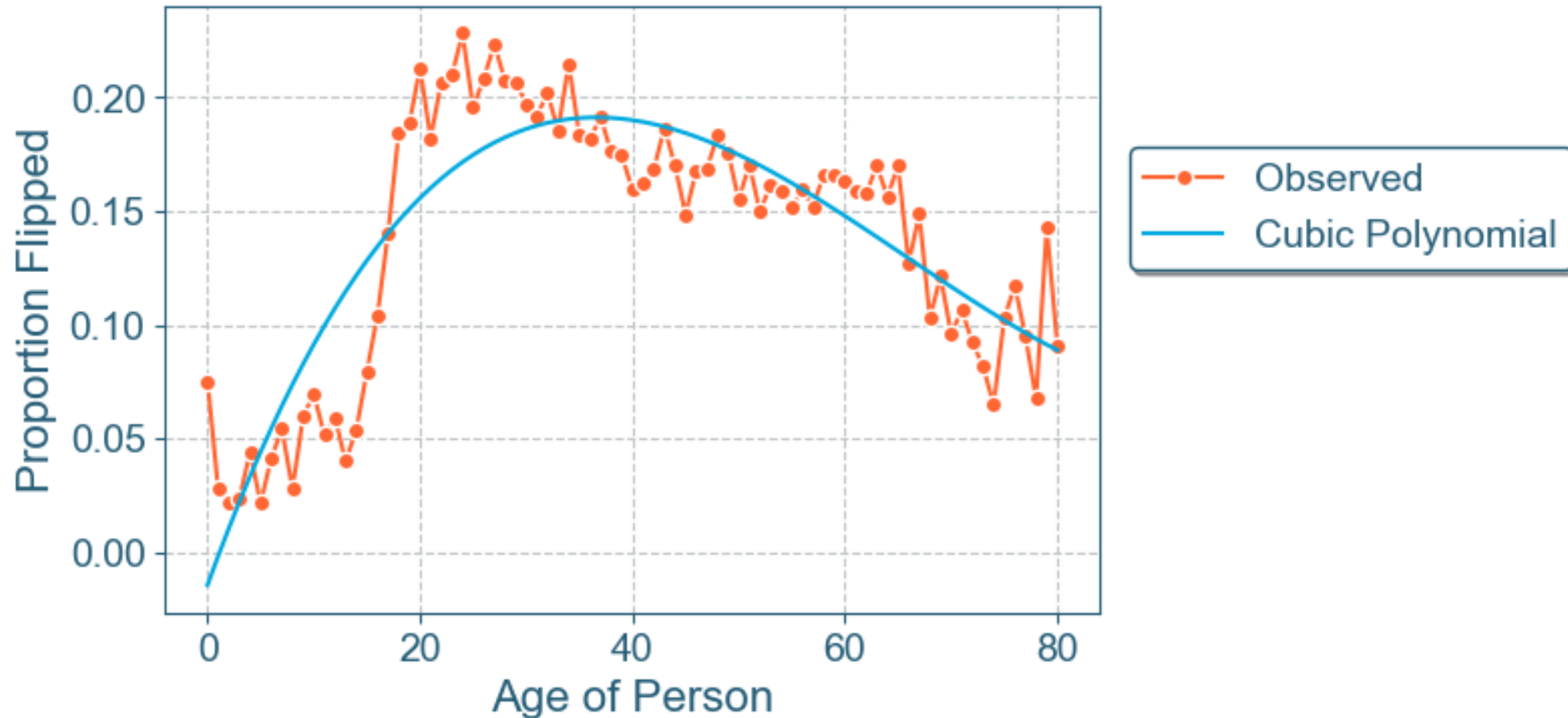


Example: Subrogation (via Accent) – Linear

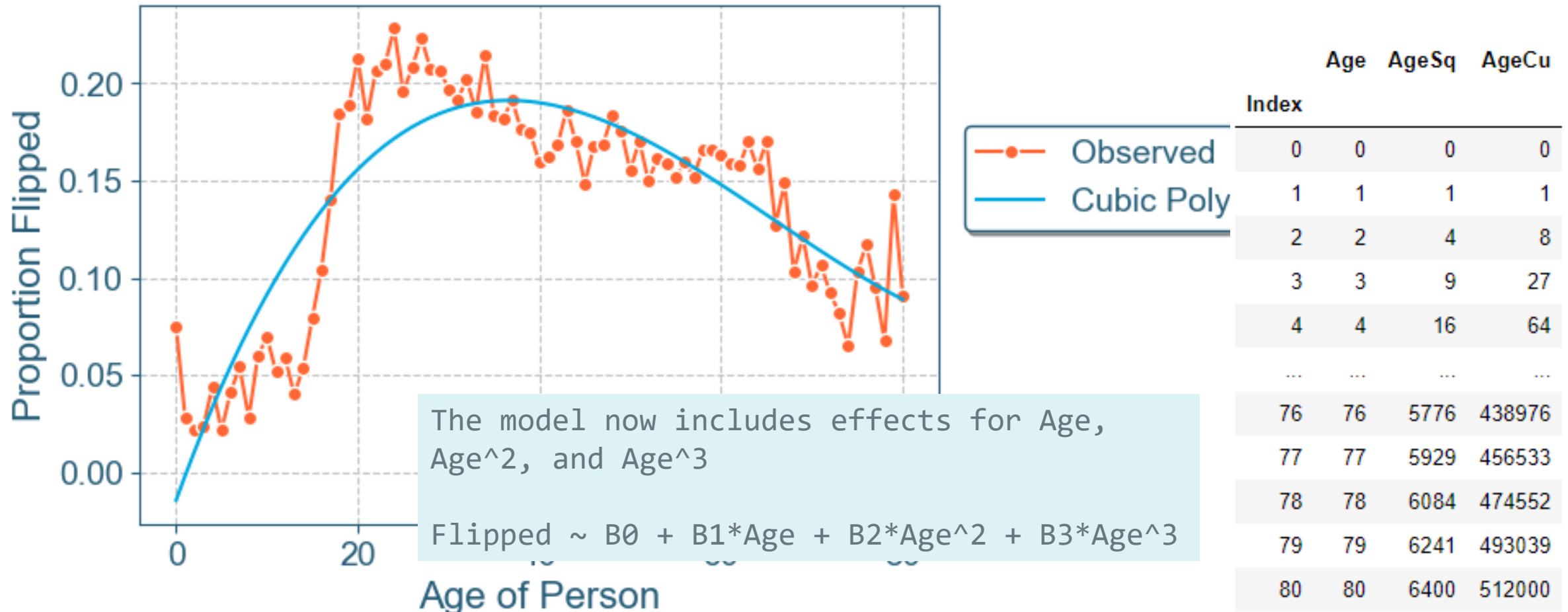


Age	
Index	
0	0
1	1
2	2
3	3
4	4
...	...
76	76
77	77
78	78
79	79
80	80

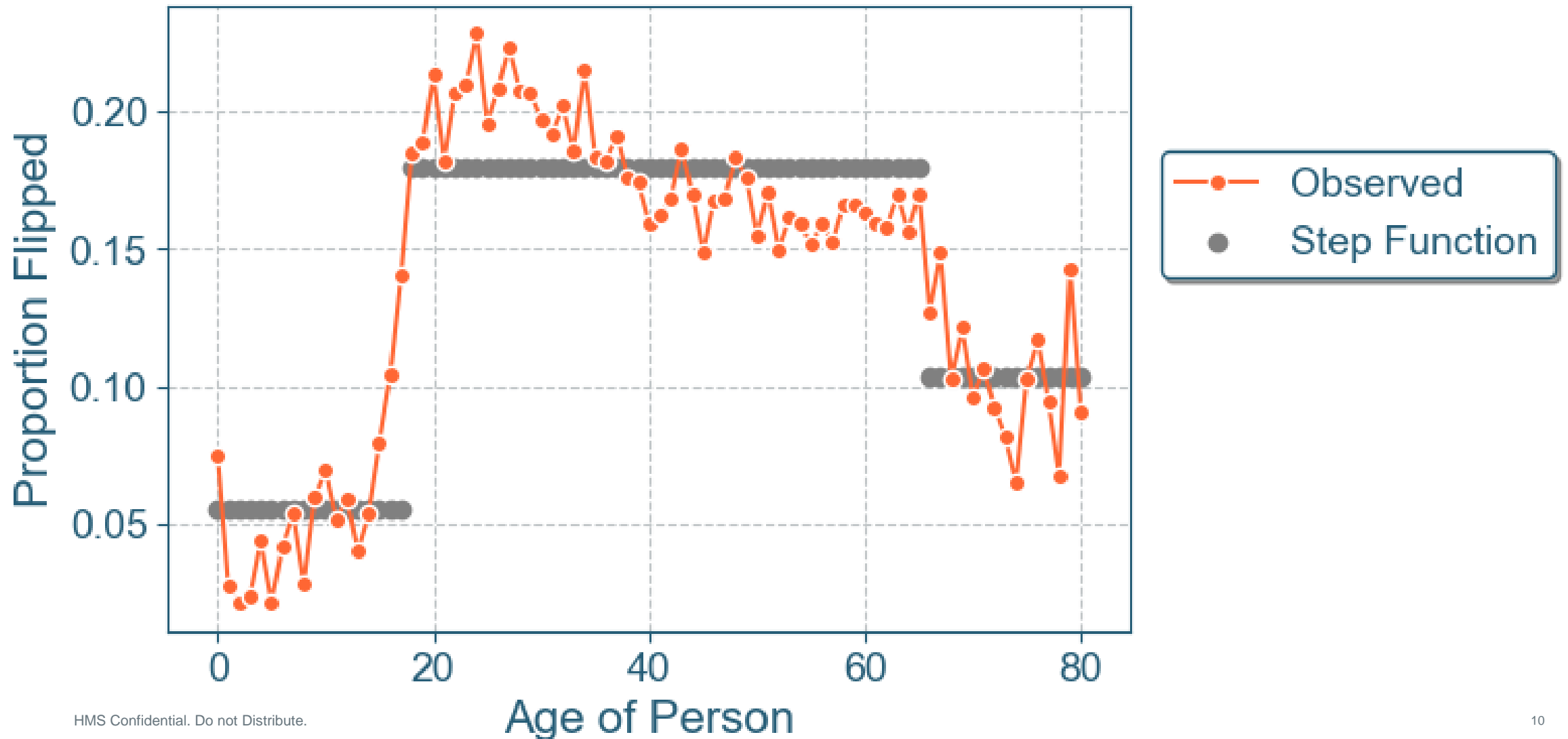
Example: Subrogation (via Accent) – Polynomial



Example: Subrogation (via Accent) – Polynomial



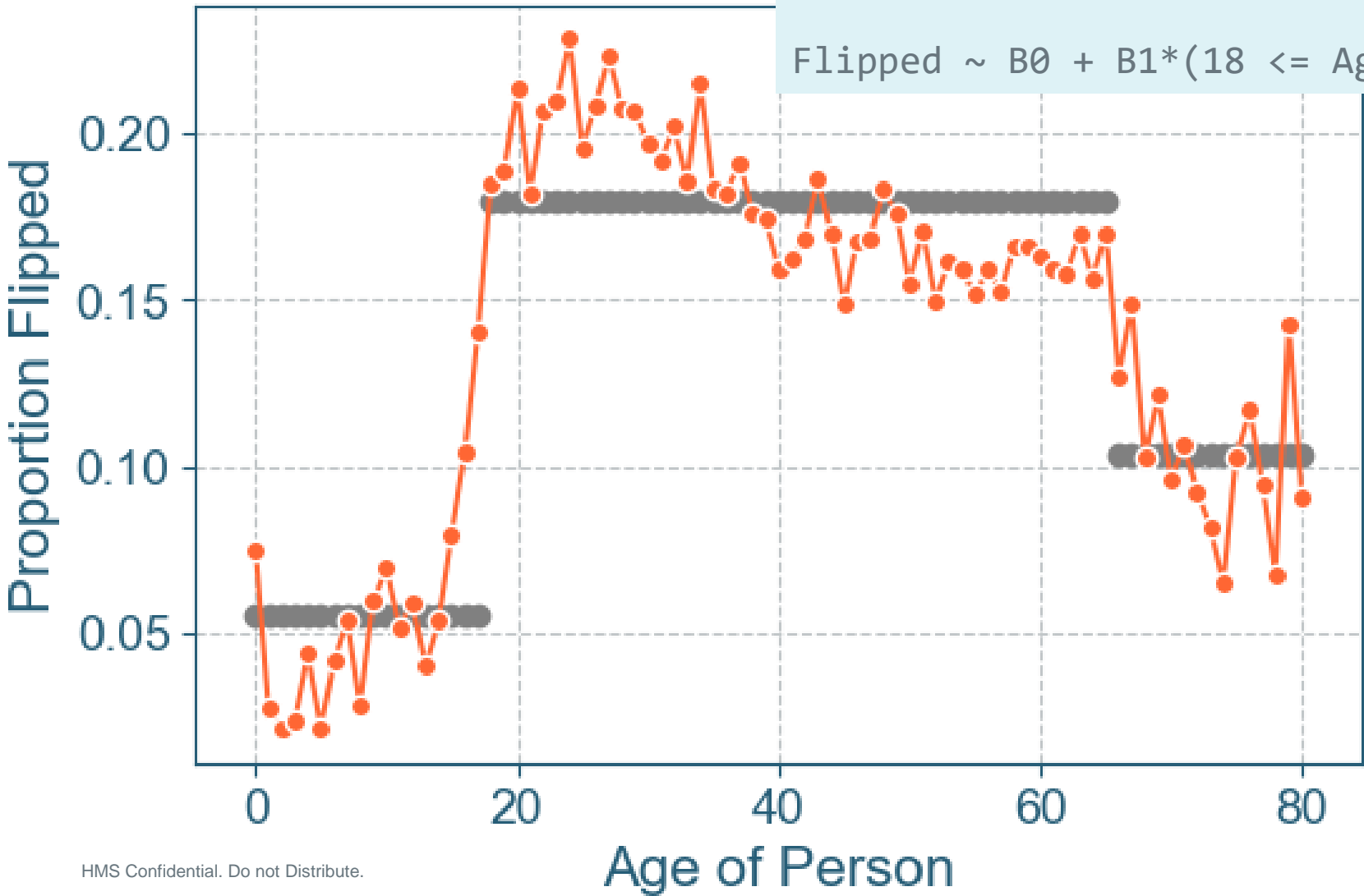
Example: Subrogation (via Accent) – Steps



Example: Subrogation

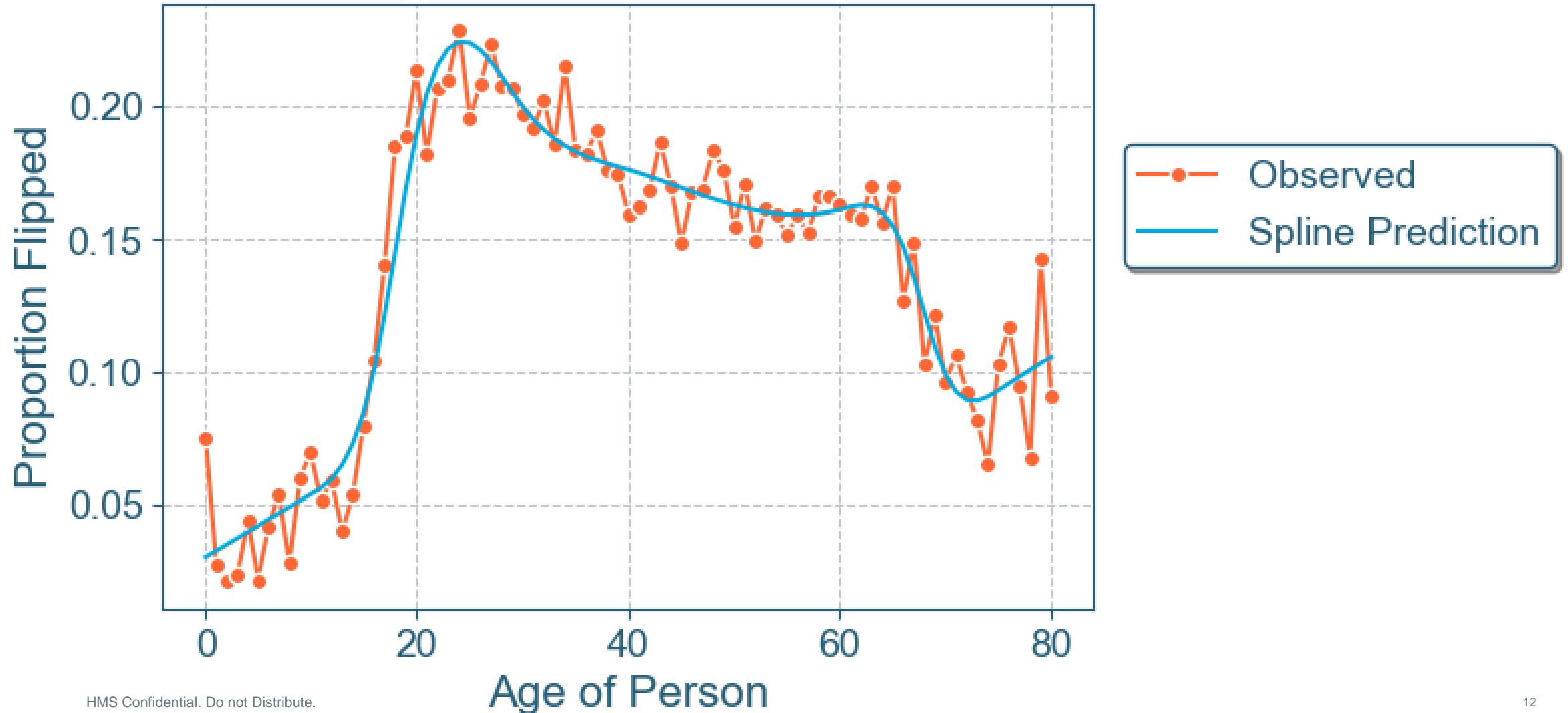
The model now is a series of step functions based on age.

$$\text{Flipped} \sim B_0 + B_1 \cdot (18 \leq \text{Age} \leq 65) + B_2 \cdot (\text{Age} > 65)$$



	Age18_65	AgeOver65
Index		
16	0	0
17	0	0
18	1	0
19	1	0
20	1	0
63	1	0
64	1	0
65	1	0
66	0	1
67	0	1

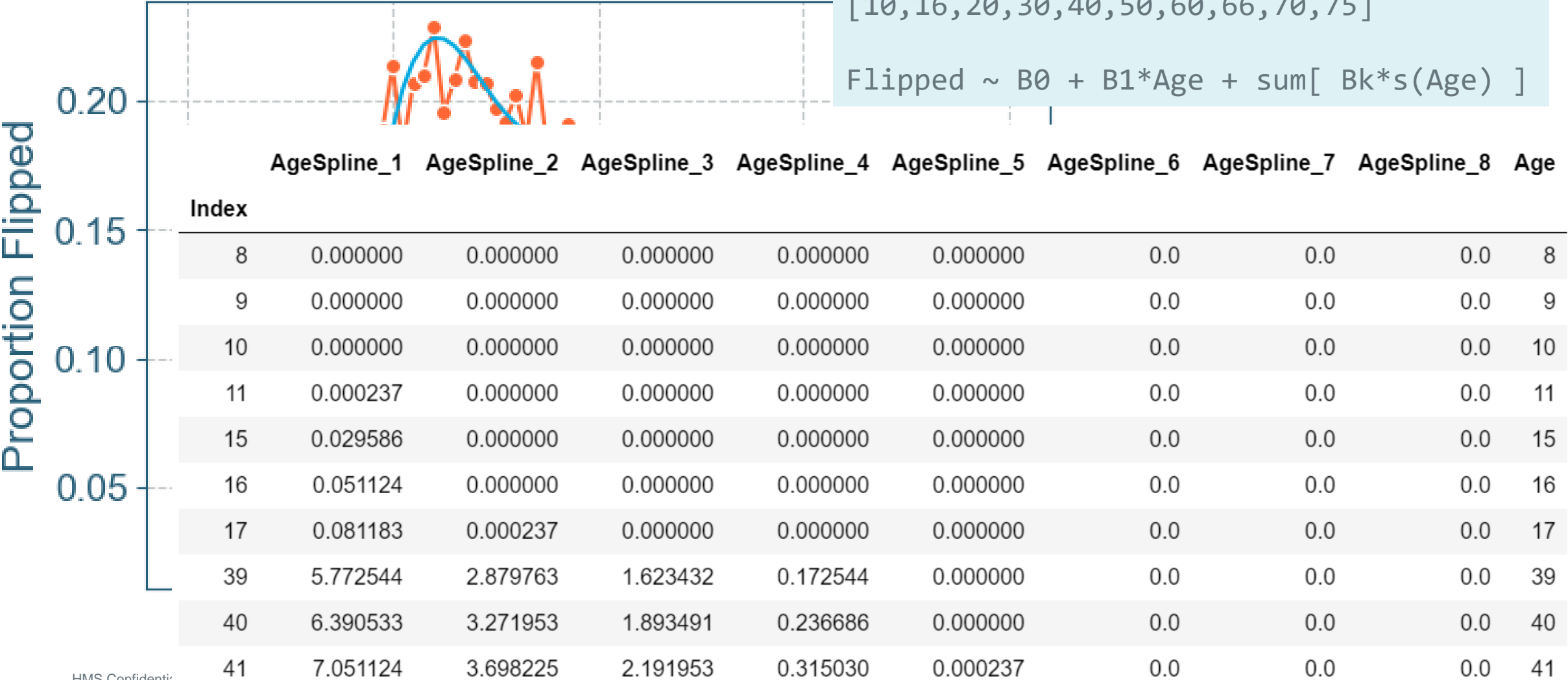
Example: Subrogation (via Accent) – Splines



Example: Subrogation (via A

Splines are more complicated functions, but similar to polynomials. This example includes knots at [10,16,20,30,40,50,60,66,70,75]

$$\text{Flipped} \sim B_0 + B_1 \cdot \text{Age} + \sum [B_k \cdot s(\text{Age})]$$



Encoding Categorical Variables

- For low numbers of categories, one-hot encoding, or dummy variables, is the best you can do:

Sex		Male	Female	Unknown
Male	→	1	0	0
Female		0	1	0
Unknown		0	0	1

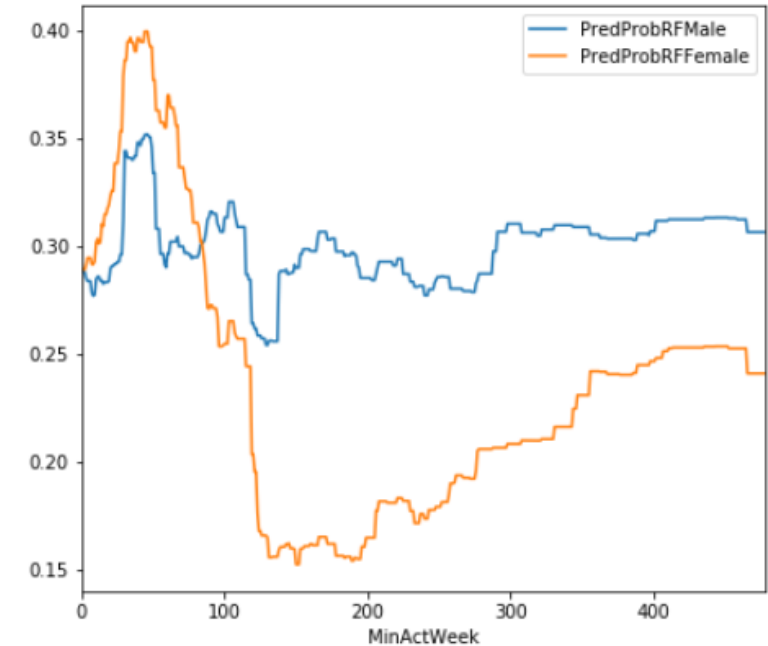
- For many categories, will have another lecture on *High Cardinality*.
 - What counts as many? Depends on data size, typically want a 100+ observations for the smallest category though.
 - For many categories, may be impossible to fit in memory

Other Example Feature Engineering Ideas

- Interaction effects between different variables,
 - $\text{Age} * \text{Male} + \text{Age} * \text{Female}$
 - $\text{Age} * \text{Client1} + \text{Age} * \text{Client2}$, etc.
- Ratio effects, e.g. for loans debt/income ratio, for claims netpaid/length of stay
- Relative to a group, e.g. $[\text{actual netpaid} - \text{average netpaid per DRG}]$ or $[\text{observed length of stay} - \text{average length of stay per DRG}]$

Different ML models and FE

- Forest based models are very good at finding non-linear effects, so don't necessarily need to include non-linear terms (like spline or polynomial variables)
- But they can often improve model fit, especially if they are important variables
- It is ok to include related features, e.g. $\log(\text{netpaid})$ and netpaid , in the same *predictive* model



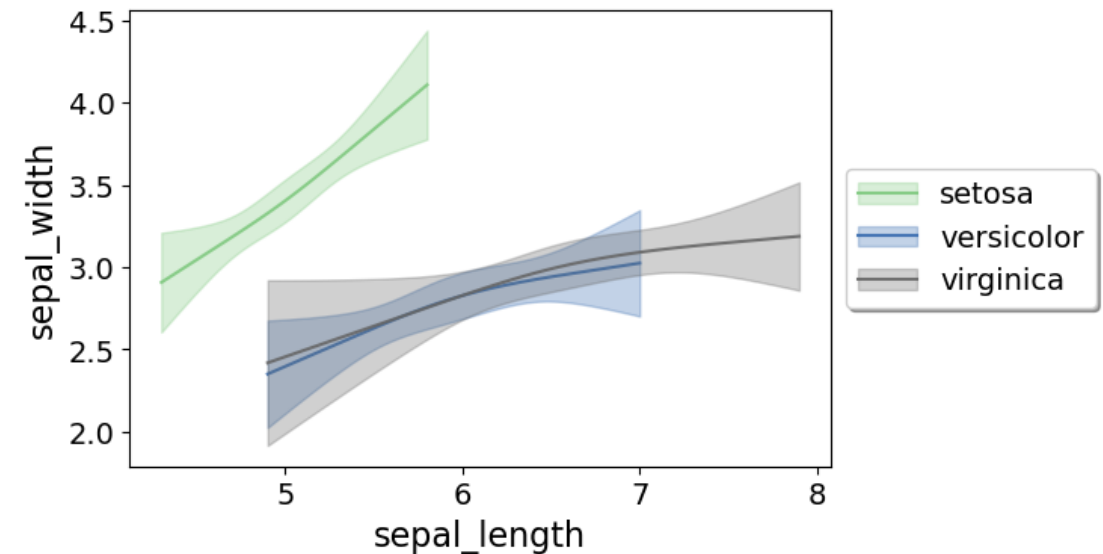
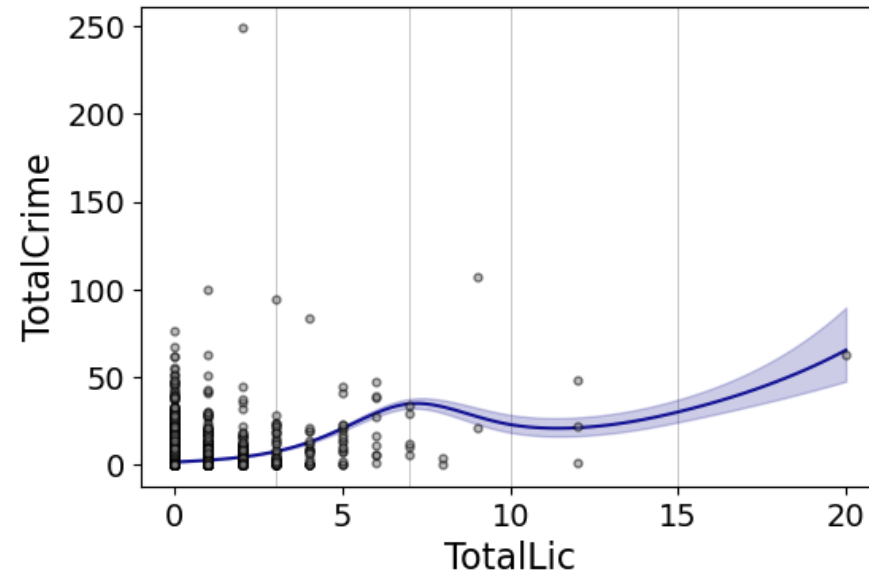
Other Resources

- [Notes on Restricted Cubic Splines](#)
- [Notebook with these examples on Github](#)
- [Smooth.py functions to plot exploratory relationships](#)

So first we specify a set of indicator variables:

$$u_+ = u \text{ if } u > 0$$
$$u_+ = 0 \text{ if } u \leq 0$$

To explain this I need to introduce the full formula for a particular spline variable. So here is that full formula

$$x_i = \left[(x - k_i)_+^3 \right] - \left[(x - k_{K-1})_+^3 \cdot \frac{k_K - k_i}{k_K - k_{K-1}} \right] + \left[(x - k_K)_+^3 \cdot \frac{k_{K-1} - k_i}{k_K - k_{K-1}} \right]$$


Future Topics

- Dealing with a high number of categories in models
- Feature Importance metrics for predictive models
- Partial dependence plots to understand functional form
- Reduced form interpretable machine learning summaries

Questions?

Future Topics

Have requests?
Let me know!

Introduction to Data Science Course Outline

Andrew Wheeler, PhD, andrew.wheeler@hms.com

- Lesson 01: Data Science 101
- Lesson 02: Machine Learning 101
- Lesson 03: Evaluating Predictions
- Lesson 04: Intro Data Transformation in Python
- Lesson 05: Data Visualization 101
- Lesson 06: Feature Engineering
- Lesson 07: Missing Data
- Lesson 08: Big Data and Parallel Computing Intro
- Lesson 09: Dimension Reduction and Unsupervised Learning
- Lesson 10: High Cardinality (Many Categories)
- Lesson 11: Intro to Forecasting
- Lesson 12: Conducting Experiments



Feature Engineering

Data Science & Machine Learning Team

02/22/2021

Andrew Wheeler, PhD

andrew.wheeler@hms.com