

Regression

In unserer Analyse haben wir ausschließlich lineare Regressionsmodelle verwendet, um die Daten zu modellieren. Dabei wurden die Variablen Age, Height, BMI, Weight, FoodConsumption und Activity berücksichtigt. Andere verfügbare Daten, wie beispielsweise jene, die lediglich aus den Datensätzen (1) und (2) bestanden, wurden nicht einbezogen, da sie für die Modellierung als nicht aussagekräftig eingeschätzt wurden.

Training eines LR-Modells

Wir haben in diesem Abschnitt das Ziel, die Daten zu modellieren und dadurch zu analysieren.

1.py	Age	Height	Weight	BMI	FoodConsumption	Activity
count	1610.000000	1610.000000	1610.000000	1610.000000	1610.000000	1610.000000
mean	33.115528	167.741615	73.504745	25.990435	5.824845	6.614907
std	9.835076	7.979873	17.096973	5.149009	1.707041	1.881387
min	18.000000	150.000000	38.760000	17.000000	1.000000	1.000000
25%	25.000000	161.000000	59.800000	21.700000	5.000000	5.000000
50%	32.000000	168.000000	70.400000	27.500000	6.000000	7.000000
75%	41.000000	174.000000	85.180000	27.500000	7.000000	8.000000
max	54.000000	193.000000	130.370000	35.000000	11.000000	12.000000

BMI: Ein BMI über 24,9 zählt als übergewichtig. Somit lässt sich ablesen, dass mehr als die Hälfte übergewichtig sind. Außerdem liegt der Durchschnittswert bei 25,99 und somit im Bereich des Übergewichts.

Activitiy: Bei Activity steht eine hohe Zahl für mehr Aktivität. Somit lässt sich sagen, dass mehr als die Hälfte Aktiv tätig sind.

FoodConsumption: Bei FoodConsumption steht eine niedrige Zahl für „gut“. Somit lässt sich sagen, dass Durchschnittlich mehr als die Hälfte zu viel Essen konsumiert.

Interpretation eines LR-Modells

Aus der Analysesicht kann man entweder mit einem Modell beginnen, das alle Features enthält, und es dann schrittweise verfeinern, oder man **beginnt mit einem künstlich einfachen Modell als Vergleichspunkt. Annahme: Geringere Aktivität führt zu höherem BMI.**

In diesem Modell werden nur „Activity“, „FoodConsumption“ und „Age“ betrachtet. Hierbei bleibt jeder p-Wert relevant. Jedoch hat sich R-squared deutlich verschlechtert, aber dennoch ist die Korrelation „stark“. Gemessen an R-squared ist also Bild 2 > Bild 3 > Bild 1.

ANOVA results						
	df_resid	ssr	df_diff	ss_diff		Pr(>F)
0	1608.0	42619.793893	0.0	NaN		NaN
1	1606.0	24396.163048	2.0	18223.630845	37579.391518	0.0
2	1604.0	388.919334	2.0	24007.243714	49505.920005	0.0

Modellvergleich mit ANOVA

Vergleich von Bild 3 mit Bild 1:

Der Vergleich in der zweiten Zeile zeigt einen signifikanten p-Wert ($\text{Pr}(>F) = 0.0$), was bedeutet, dass Bild 3 eine signifikante Verbesserung gegenüber Bild 1 darstellt.

Vergleich von Bild 2 mit Bild 3:

Der Vergleich in der dritten Zeile hat ebenfalls einen signifikanten p-Wert ($\text{Pr}(>F) = 0.0$), was darauf hinweist, dass Bild 2 eine signifikante Verbesserung gegenüber Bild 3 darstellt.

Schlussfolgerung:

Da sowohl der Vergleich von Bild 3 mit Bild 1 als auch der Vergleich von Bild 2 mit Bild 3 signifikant ist, können wir sagen: Bild 2 ist besser als Bild 3 somit auch besser als Bild 1.