

Cervical cancer risk classification

JiEun Song

Overview

Cervical cancer is the fourth most common type of cancer among women in the world. More than 10,000 new cases are diagnosed each year in the US. Although each year cervical cancer kills about 4,000 women in the US, it is the most preventable type of cancer. Cervical cancer is highly preventable because many risk factors are known for initiation of the cancer. My aim is to examine the dataset, and use a classification model to predict whether a patient is diagnosed with cervical cancer or not.

I found the dataset on Kaggle, which has information of 850 patients with 37 columns with risk factors. Each patients are classified as cancer positive or negative.

Features

Number of partners	STDs:vaginal condylomatosis
First sexual intercourse	STDs:vulvo-perineal condylomatosis
Number of pregnancies	STDs:syphilis
Smokes	STDs:pelvic inflammatory disease
Smokes (years)	STDs:genical herpes
Smokes (packs/year)	STDs:molluscum contagiosum
Hormonal contraceptives	STDs:AIDS
Hormonal contraceptives (years)	STDs:HIV
IUD	STDs:Hepatitis B
IUD (years)	STDs:HPV
STDs	STDs:number of diagnosis
STDs (number)	STDs:time since first diagnosis
STDs:condylomatosis	STDs:time since last diagnosis
STDs:cervical condylomatosis	

Known unknowns

I have studied cancer biology, especially brain cancer, and brain cancer is distinct from other cancers, so I am not familiar with cervical cancer. I know that HPV infection is a major cause of cervical cancer, but am not sure prediction power of each features.