

Continuous variables (features)

X continuous feature \mathbb{R}

categorical variable

'dog', 'cat', 'gerbil'

→ one-hot encoding

animal	dog	cat	gerbil	const
'dog'	1	0	0	1
'dog'	1	0	0	1
'cat'	0	1	0	1
'gerbil'	0	0	1	1
'cat'	0	1	0	1

dog + cat + gerbil

→ columns can't
linearly independent

gender height weight

M	F	height
1	0	
1	0	
1	0	
1	0	
0	1	
0	1	
0	1	

ordinal variable

education	→ no HS	0
	HS	1
	some collg	2
	collge	3
	grad degree	4

linear regression → categorical variables
one-hot encoding

tree-based models → handle ordinal variables

Classification

↳ binary classification

target data Y can be thought of
as all 1s and 0s

↳ given some design matrix X ,
I want to be able to predict
the probability that $Y_i = 1$ P_i

$$P(Y_i = 0) = 1 - P_i$$

machine learning: ① regression problems → continuous
output

② classification

↳ binary (categorical)
output

evaluate regression:

$$\frac{1}{n} \sum (Y_i - \hat{Y}_i)^2$$

True value classification:

T

F

True positive

False positive

T

Positive

F

False negative

True

① Accuracy \rightarrow the proportion of prediction that are correct

$$\frac{\#TP + \#TN}{N}$$

② Precision \rightarrow if we predict positive, how often it is correct

$$\frac{\#TP}{\# \text{ predicted positives}}$$

③ recall (sensitivity) \rightarrow what percent of positives do we catch?

$$\frac{\#TP}{\# \text{ positives}}$$

$$\textcircled{4} \text{ specificity} = \frac{\#TN}{\# \text{ negatives}}$$

Professor Dahn hates all of these!

⑤ Brier score

$$\sum_{i=1}^n (\hat{p}_i - y_i)^2 \Rightarrow \text{squared error using probabilities}$$

so of y_i 's are 1

②.1. of my data set is 1's

↳ lowest predicted probability
I would set

0.3

$$\hat{y}_i = 0$$

