

# Session 4

## One Hot Encoding



بیشتر مدل های یادگیری ماشین با داده ها عددی کار می کنند.

؟ اگر بخواهیم بررسی داده های Categorical کار کنیم، چه طور باید این مشکل را حل کنیم ؟

خب ساده است باید داده هایمان را به عدد تبدیل کنیم، اما چه طور ؟

1. Integer Encoding

2. One-Hot Encoding

"Red" : 1

"Yellow" : 2

"Green" : 3

Color
Red
Red
Yellow
Green
Yellow

Integer Encoding  
→

Color
1
1
2
3
2

ظاهر آنکه همه چیز به نظر خوب می‌رسد، اما واقعاً اینجوری نیست؟ **ordered relationship**

اگر ما به هر Category یک عدد بدهیم و مدل ما ارزش آن عدد را در حسابات دخیل کند در حالی که

عدد بزرگتر نشان دهنده‌ی برتری یک Category نسبت به دیگری نیست.

؟: برای رفع این مشکل چه کنیم؟ می‌توانیم از **One-Hot Encoding** استفاده کنیم.

Color		Red	Yellow	Green
Red		1	0	0
Red	One-Hot Encoding →	1	0	0
Yellow		0	1	0
Green		0	0	1
Yellow		0	1	0

۲: مشکلات One-Hot Encoding چیست؟

۱. افزایش تعداد ویژگی‌ها

۲. Dummy Variable Trap

Red	Yellow	Green		Red	Yellow
1	0	0		1	0
1	0	0		1	0
0	1	0	→	0	1
0	0	1		0	0
0	1	0		0	1