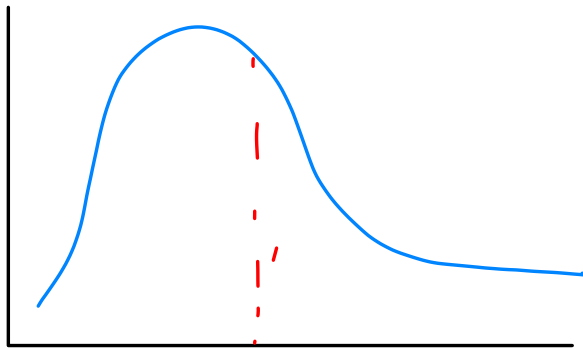


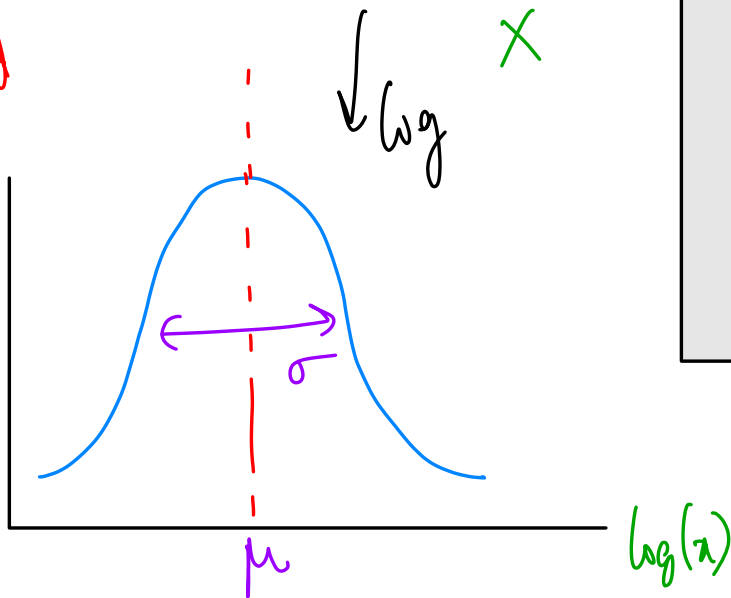
FEATURE ENGINEERING-1



Original



Transformed



LOG NORMAL
DISTRIBUTED

$$E(x) = \mu_x = e^{\left(\mu + \frac{\sigma^2}{2}\right)}$$

$$\text{Var}(x) = (e^{\sigma^2} - 1) \left(e^{2\mu + \sigma^2} \right)$$







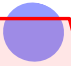




Feature Engineering

Acrofit

Target

Features / attributes

Datapoint/
Records

	Age	Income	Fitness	Miles	Gender	Education	Product
 X_1							KP281
 X_2							KP481

input \rightarrow Model \rightarrow prediction

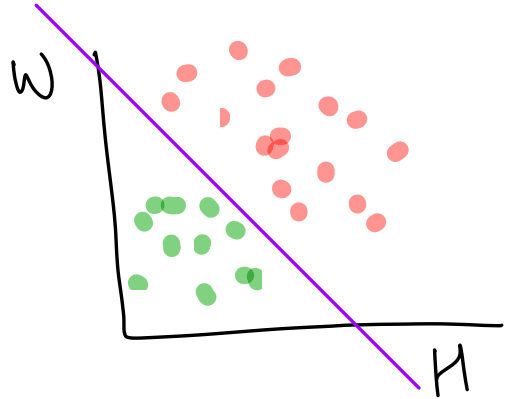
mathematical
function
L

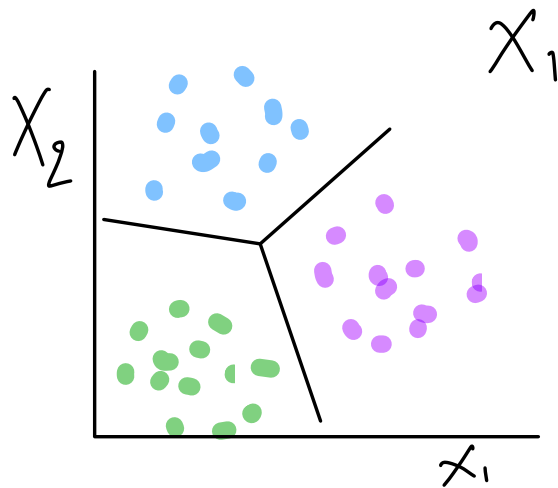
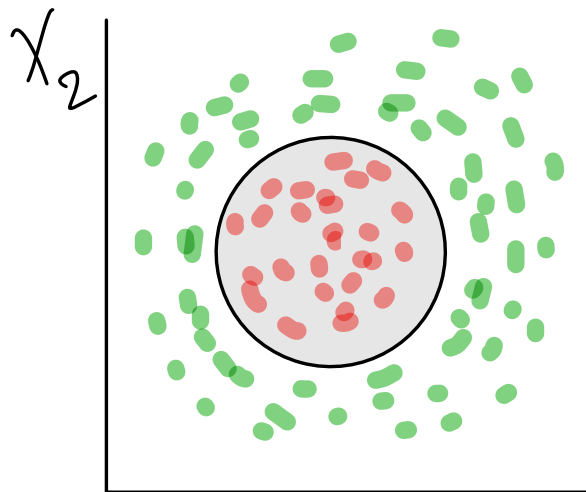
$$y = mx + c$$








Complex complicated

Classification \leftarrow
Regression
Clustering

$$x^2 + y^2 + 2gx + 2fy + c = 0$$





Height	Weight	BMI $\frac{H}{W^2}$	Fitness
			
			
			

* Z test $\rightarrow \mu, \sigma \checkmark$ (Numerical v/s Categories)

* T test (t test - 1 samp, t test - rel, t test - ind, t test - ind - from stats)
 $\sigma \times \rightarrow n < 30$, (Numerical v/s Categories)

* KS Test $\rightarrow \mu_1 \neq \mu_2$, distributions are diff.

* χ^2 \rightarrow Goodness of fit (CATEGORICAL)
Test of Independence (CAT v/s CAT)

* ANOVA \rightarrow More than 2 grps (Numerical)

* KRUSKAL \rightarrow when assumpⁿ of ANOVA fail

* CORRELATION \rightarrow (NUMERICAL v/s NUMERICAL)