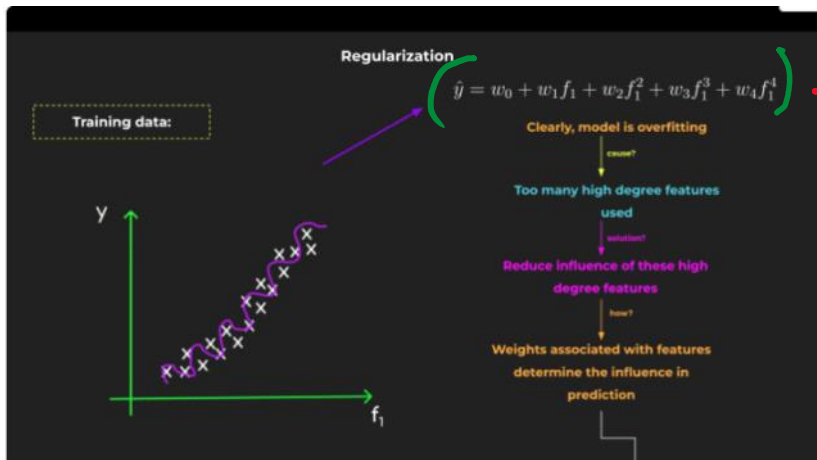


Instead of training and testing Multiple Polynomial Models with different degrees, is there a better solution?



→ Maybe f_1^2 is not required
 $w_0 + w_1 f_1 + w_3 f_1^3 + w_4 f_1^4$ May be enough

Maybe a higher degree like f_1^4 May not be required
 $w_0 + w_1 f_1 + w_2 f_1^2 + w_3 f_1^3$ May be enough

Is there a way we can reduce the influence or impact of degrees on features to avoid overfitting?

For example :- $y = 2 + 3.1x_1 + 4.2x_1^2 + 2.2x_1^3 + \left(\overset{\uparrow \text{Very low influence}}{0.01} x_1^4 \right)$
 ↓
 observe and compare with other Co-efficients

Any technique to help us achieve the above?

{Regularization !} \Rightarrow as part of training process

Instead of Minimizing just the loss

function using $\min L$, add another term

to the objective function $\Rightarrow \sum W$

$$L = SSE = \sum (y^{\wedge} - y_i)^2 + \sum W^2 \rightarrow \text{co-efficients / impact}$$

New objective :- minimize $\left\{ L + \sum_{i=1}^d W_i^2 \right\}$

Why do we add $\sum_{i=1}^d W_i^2$?

Along with loss function, We also want to minimize impact of W (weights) as much as possible !

Outcome :- Influence of **Unwanted weights** significantly reduces !!

In Linear Regression, What do we want to minimize?

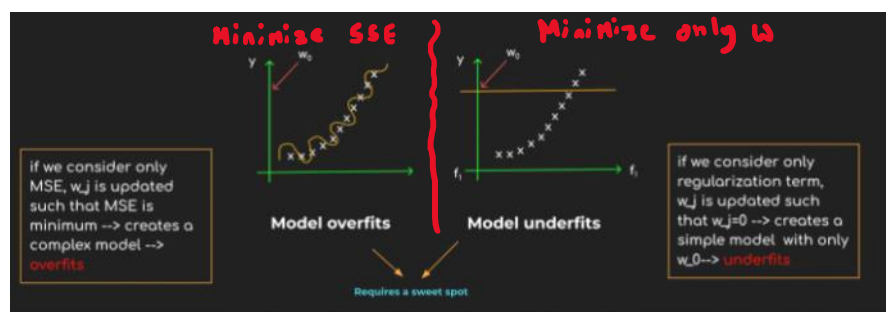
$$SSE = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

Downside of using this approach alone is overfitting

To reduce impact of a feature F_1 or F_1^2 or F_1^3 ,
What do we need to minimize?

$\{w_i\}$

Downside of using this approach alone is Underfitting



Giving birth to J

$$loss = \min_{w_j} \frac{1}{n} \left[\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 + \sum_{j=1}^d w_j^2 \right]$$

$$loss = \min_{w_j} \frac{1}{n} \left[\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 + \sum_{j=1}^d w_j^2 \right]$$

Updates w_j such that MSE is least

Makes w_j such that its value ≈ 0 , except w_{w_0}

All okay but we want some control over how much impact regularization has right?

Introducing λ (lambda), also called strength

$$\text{loss} = \min_{w_j} \frac{1}{n} \left[\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 + \sum_{j=1}^d w_j^2 \right]$$

becomes

$$\text{loss} = \min_{w_j} \frac{1}{n} \left[\sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2 + \lambda \sum_{j=1}^d w_j^2 \right]$$

$\lambda \geq 0$ Acts as a tuner

if $\lambda = 0$, no regularization, overfitting

if $\lambda = 1000000$, weights will heavily reduce \Rightarrow underfitting

But be careful!

Too high λ causes all weights to become 0,
which causes underfitting

Too low λ causes all weights to stay as they are, causing overfitting

$$\lambda \sum_{j=0}^M W_j^2$$

→ together, this is known as
{penalty.}

Interview → Regularization is the process
of adding a penalty term to your loss
Function to avoid overfitting

$$-\lambda \sum_{j=0}^M W_j^2$$

W^2 → Goal is to Minimize
Weights of Features that cause
overfitting.

This means, some of the weights
become close to 0 but never 0!

$$y = 3.9x_1 + 2.1x_1^2 + \underline{0.01}x_1^3 + 1.1x_2 + \underline{0.001}x_2^2$$

This type of regularization is called L2
or Ridge regularization / Ridge Regression

$$-\lambda \sum_{j=0}^M |W_j|$$

↓
observe
the
change
to
 $|W_j|$
from
 W_j^2

→ Goal is to completely turn off
Weights of Features that cause
overfitting.

This means, some of the weights
will be exactly 0

ex:- $y = 3.9x_1 + 2.1x_1^2 + \underline{0}x_1^3 + 1.1x_2 + \underline{0}x_2^2$

This type of regularization is called L1
 or Lasso regularization

L1 Regularization / Lasso

$$\text{Cost} = \sum_{i=0}^N (y_i - \sum_{j=0}^M x_{ij} W_j)^2 + \lambda \sum_{j=0}^M |W_j|$$

L2 Regularization / Ridge

$$\text{Cost} = \sum_{i=0}^N (y_i - \sum_{j=0}^M x_{ij} W_j)^2 + \lambda \sum_{j=0}^M W_j^2$$

Loss function Regularization Term

So when to use which?

Interesting property of L1 and L2 :

⇒ L1 :

1. Since L1 makes wts of useless features = 0 } ⇒ Feature selection
 Hence used in feature selection
2. As L1 creates sparse weight vector, it reduces the number of non-zero parameter leading to Efficient memory usage and faster computation in high dimensional datasets

⇒ L2 : } ⇒ Minimize impact of high value features

1. The square penalty term w_j^2 encourages smoother parameter values + beneficial in case of noisy data
2. L2 tends to ^{Not} zero out weight value of less useful features ⇒ helpful when we want model to aim for a balanced use of features

What if I'm not sure which to use?

What if I want the best of both

Worlds? Some weights completely turned

off and some very low but not really 0.

Solution: Elastic net regression

Elastic net regularization:

Since both L1 and L2 Regularization has its perks, is there a way to combine them??

$$\left\{ \min_{w_j} \left[\underbrace{MSE}_{L1} + \lambda_1 \sum_{j=1}^d |w_j| + \lambda_2 \sum_{j=1}^d w_j^2 \right] \right\}$$

L1 L2



Elastic Net Regularization

At what temperature to bake a cake?

What brand of shirt to purchase for yourself?

Which size of new pair of shoes to purchase?

These are all hyperparameters in ML.

You try different variations to find the best option. This is called Hyperparameter Tuning

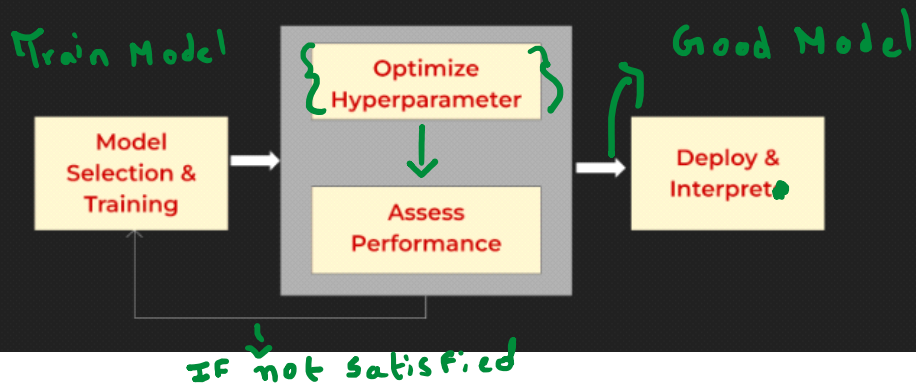
In polynomial Regression, degree (highest power)
You want to fit your Model with.

In Regularization λ (strength of regularization)

Q. How do we choose the hyperparameter values?

— Like in our example, We baked several cakes to get the best one.

We can produce several models using different hyperparameter and can compare them using any method.

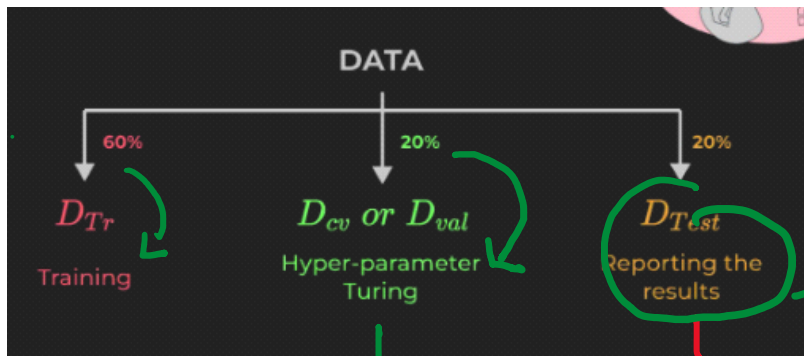


We may take 80% of data to train
and 20% to test.

We perform hyperparameter tuning on Model
by and validate results on test data
by looking at test R_{sqd}

Alternative approach :-

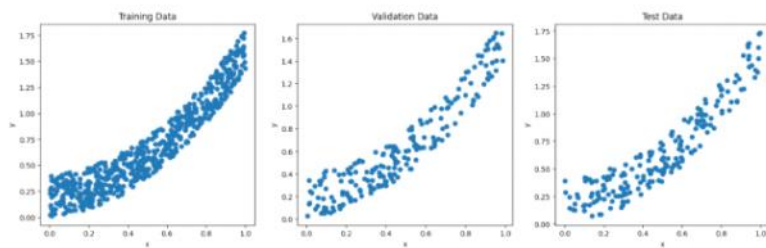
Cross Validation



Final test $R_{sqv} = 0.7$

Back and Forth
← {Purely For this purpose}

{Purely For this purpose} ⇒ only once



Cross validation heavily reduces number of training samples. only 60% becomes training

Solution \rightarrow {K-Fold Cross Validation}

Divide entire data into K chunks. Lets say 4 in this case.

| Train on | Unseen validate on | Metric |
|--------------------------------|--------------------|------------|
| 75% Chunk 1 Chunk 2 Chunk 3 | 25% Chunk 4 | $R_{sq} 4$ |
| 75% Chunk 1 Chunk 2 Chunk 4 | 25% Chunk 3 | $R_{sq} 3$ |
| 75% Chunk 1 Chunk 3 Chunk 4 | 25% Chunk 2 | $R_{sq} 2$ |
| 75% Chunk 2 Chunk 3 Chunk 4 | 25% Chunk 1 | $R_{sq} 1$ |

as test R_{sq} Metrics

$$\text{Final Metric} \rightarrow \text{Avg } R_{sq} = \left(\frac{R_{sq} 4 + R_{sq} 3 + R_{sq} 2 + R_{sq} 1}{4} \right)$$

{ Test Metric to report }

Purpose of Validation Set
to validate different hyperparameters } Satisfied

Report Model results on unseen data } Satisfied

\downarrow

In a way, we made
sure all data was unseen
at some point or another!