

ML - Bias Variance,

Regularization

Hyperparameter Tuning

Cross Validation

Revision

Polynomial Reg $\rightarrow n > 1$



Training
Performance

Select degree with highest R-sq.

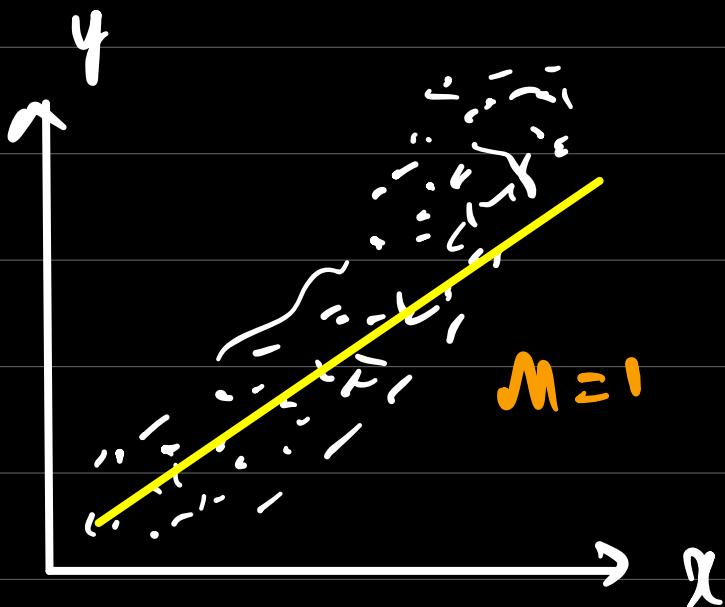


Best degree poly eq may not work the best
on Testing data. Why??

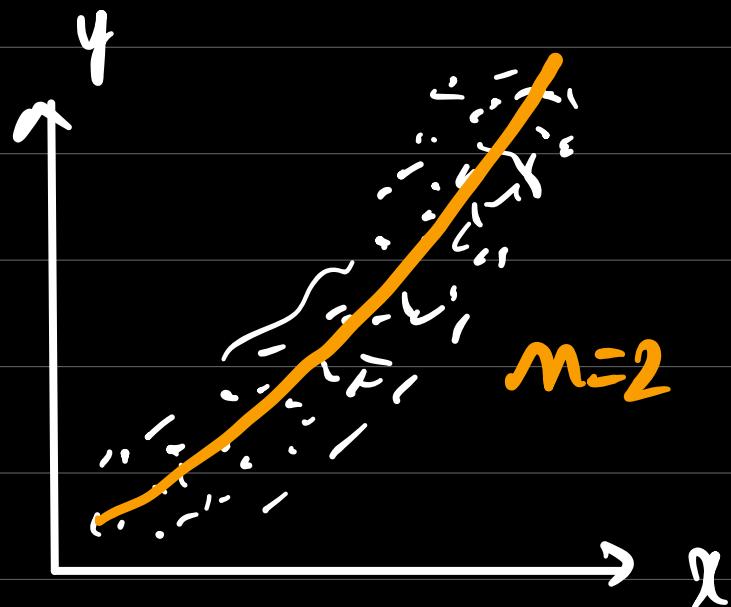
↓
Overfitting to train date (learnt patterns + noisy)

too
^

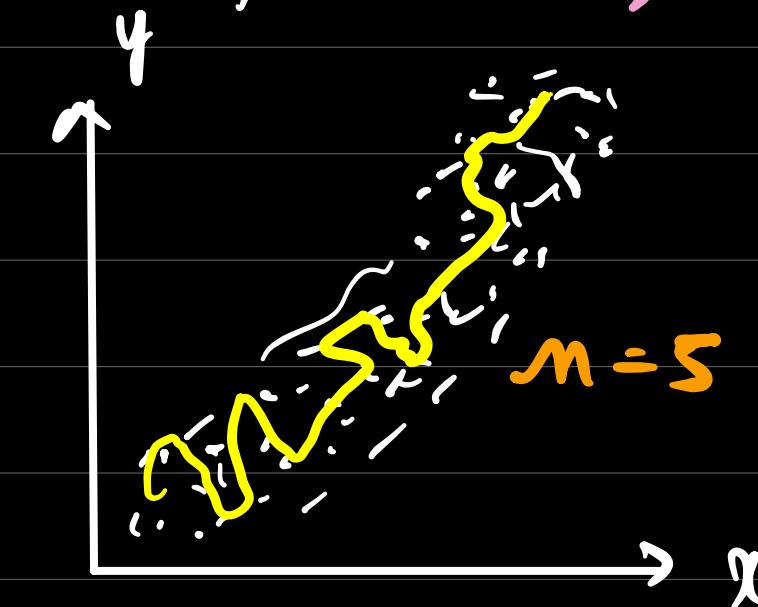
Underfitting - simple models (not even able to learn patterns)



Underfitting



JUST right fit



Overfitted

Bias Variance Trade-off

Bias: Model is constantly 'biased' towards doing wrong predictions. - HIGH BIAS

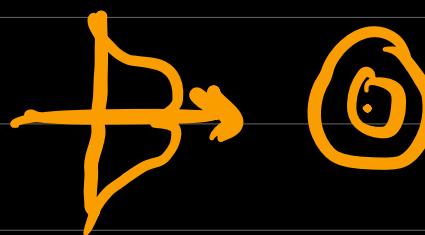
Variance: Just with a slight change in x

$$[1, 3, 4, 4, 2] \rightarrow [1.1, 3, 4.5, 2.1]$$

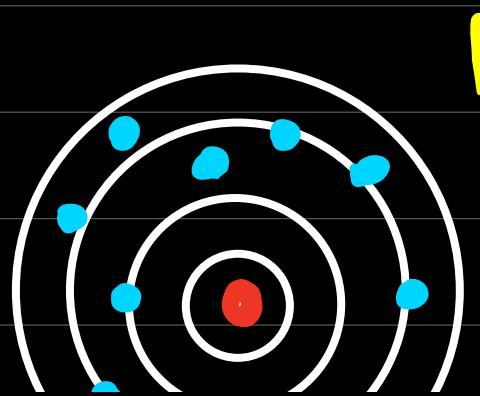
$$x_i \rightarrow y$$

$$x_i \rightarrow y \pm \Delta$$

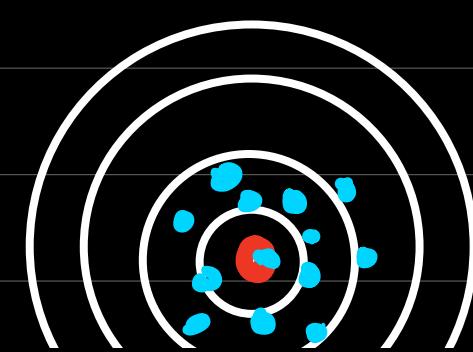
Predict is 'varying' a lot - HIGH VAR.

e.g.: You are learning archery 

- Prediction \hat{y}
- Ground truth y

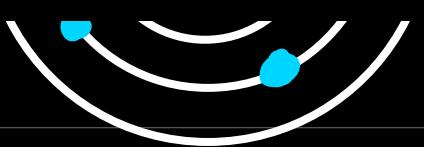


1



2

Overfitting



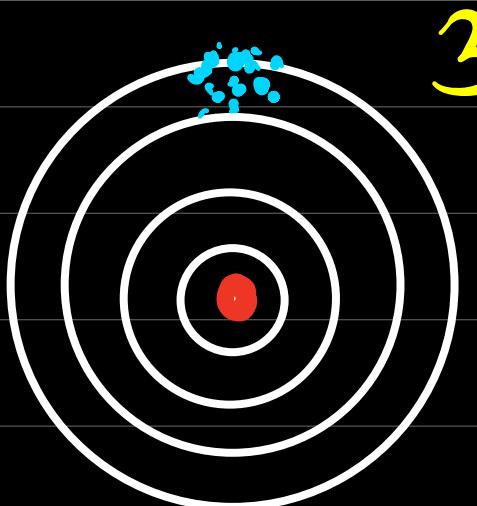
↑ BIAS

↑ VARIANCE



↓ BIAS

↑ VARIANCE



3

Underfitting

↑ BIAS
↓ VARIANCE



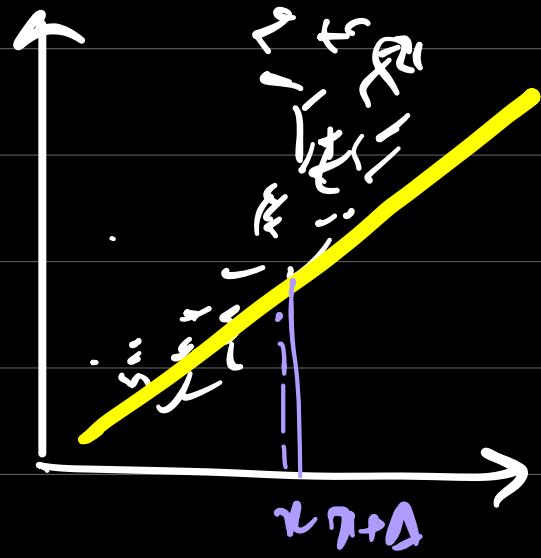
4

= Ideal

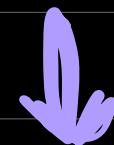
↓ BIAS
↓ VARIANCE

I ideally WANT ↓ Bias & Variance.

Underfitting



Cant even learn
the patterns in
the data

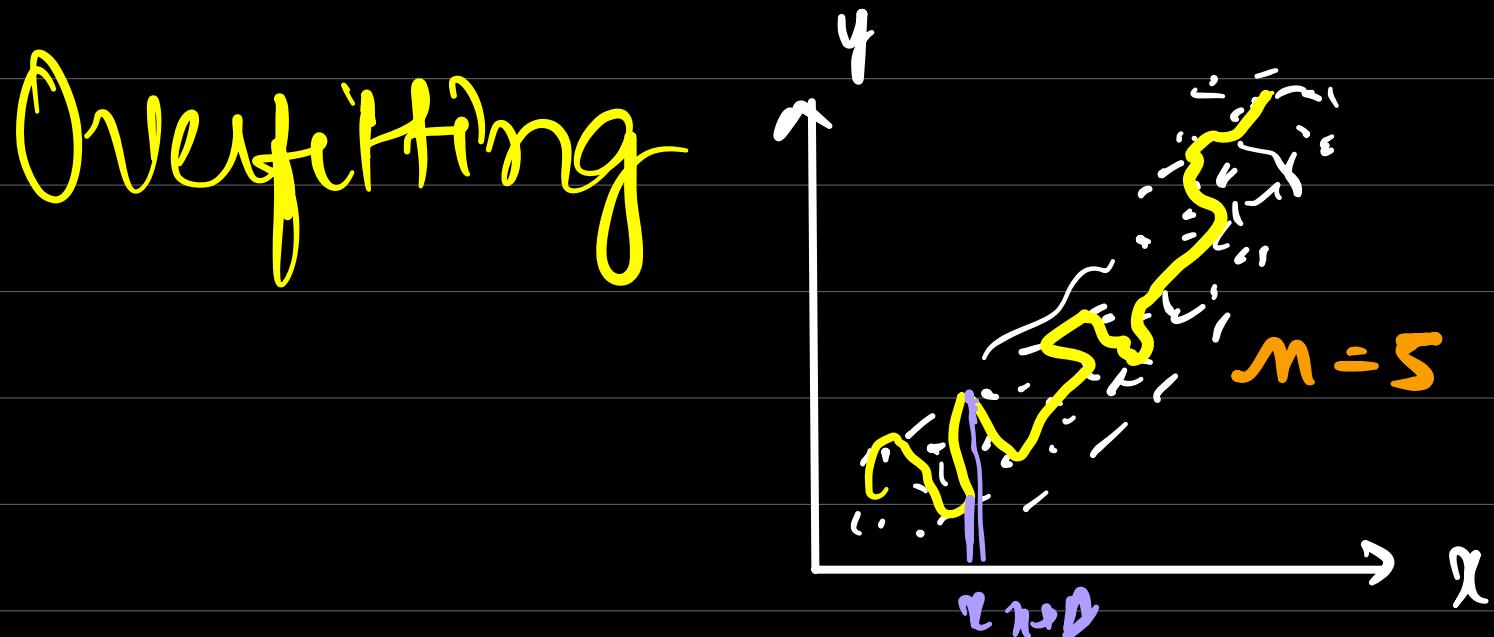


Bad predictions

Bias ↑↑

Variance ↓↓

Simple Models - Prob. of ↑ Bias.



Bias - Is model making bad pred.
always?

NO

Training Data - GOOD

Testing Data - May or not
be good.

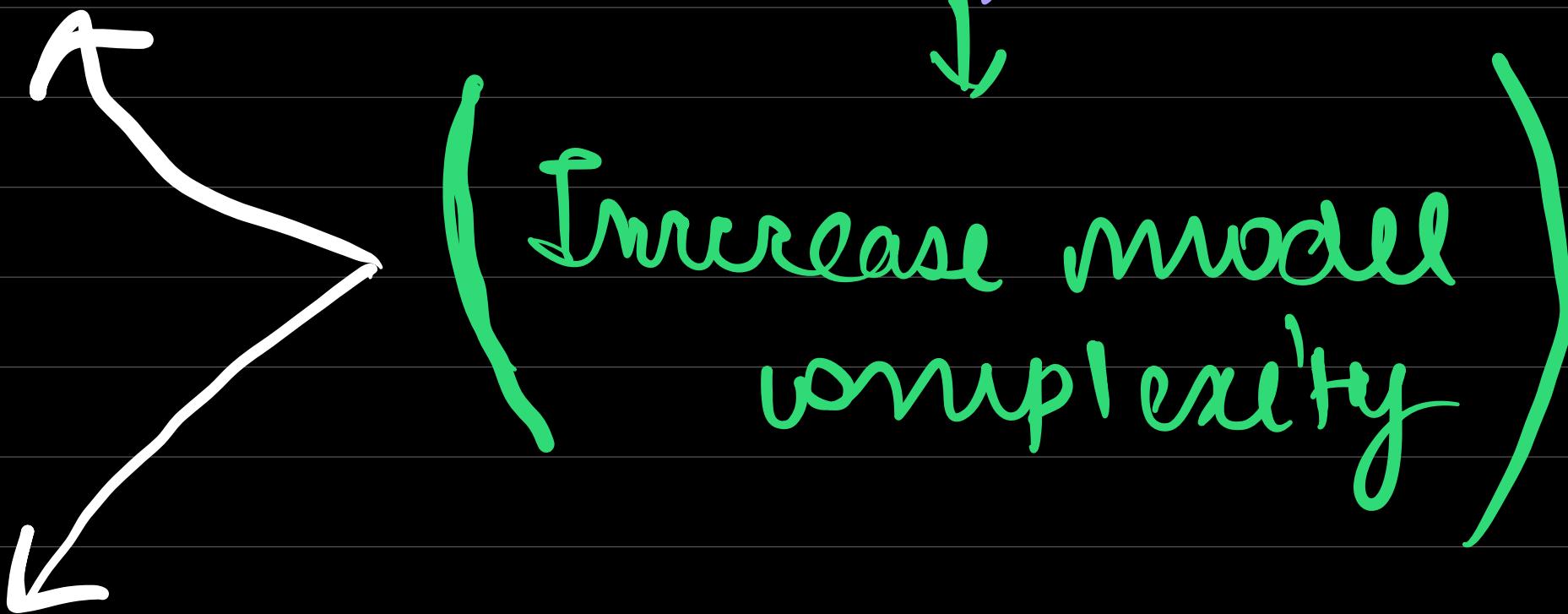
↓ Bias, ↑ Variance.

Practical tips for detecting U or O ?

Performance	Underfitting	Oversetting
Train	BAD	GOOD
Test	BAD	BAD or NOT AS GOOD
 Still better than an underfitted model		

Training

Bad Performance → Underfitting



Good Performance → look at the
testing performance

Testing

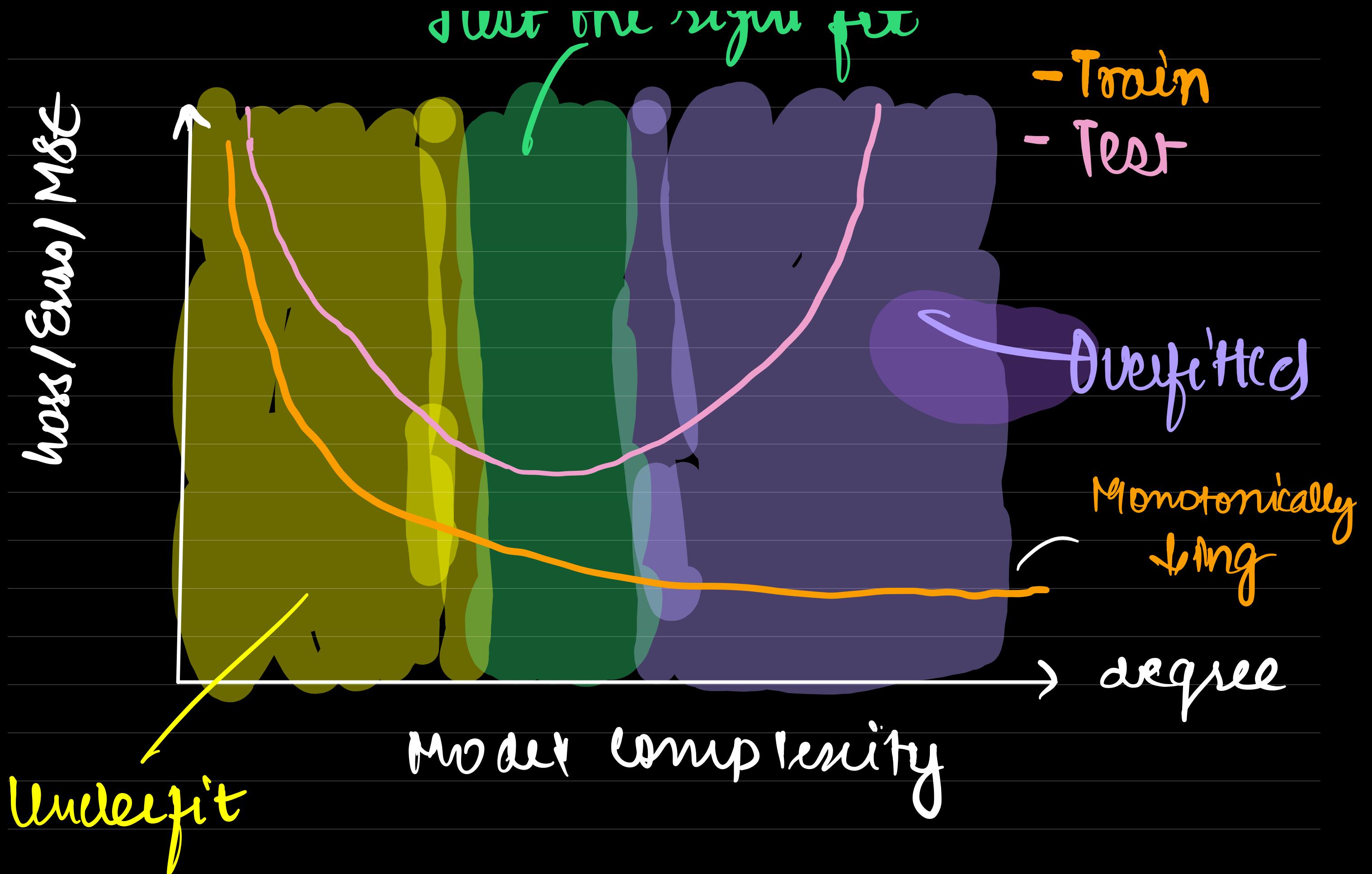
Bad Performance → Overfitting



Reduce model
complexity

Good Performance - DANCE
Deploy

Final U. solution



Q Train - 90% Test - 90%

↓ ↓
~ Good Good - No overfitting

Just the right fit

Q Test Performance > Train Performance.

① Data Leakage (Reusing test data)

② Sampling data from rare region.
qst.

- Test data as a nice / ideal / beautiful / no
strange subset - NO OUTLIERS

NO WEIRDNESS

Regularisation.

degree = 1, 2, 3, 4, 5, 6, 7, 8, 10, 20, 30, 40

$$y = 0.0002x^{20} + 1.8x^{18} + 20x + 40$$

less rmp.

more rmp.

Maint havee knowned at this

$$y = 20x + 40 + \dots$$

GOAL: To select a higher degree or
complex

YET

somewhat prevent overfitting

AUTOMATICALLY

Assume, chose d=30 model

$$y = w_0 + w_1 x_1 + w_2 x^2 + w_3 x^3 + \dots + w_{29} x^{29} + w_{30} x^{30}$$

Goal: Prevent overfitting

Reduce w's of less imp features.
magnitude.

$$\text{LOSS} = \text{MSE}$$

UNREG



Training = Minimizing loss.

= Minimizing errors. $y - \hat{y}$

Training
data.

\therefore Chosen a high degree polynomial.

\xrightarrow{d} Reduce errors \rightarrow Overfitting

solution??

↓ dealing overfitting

loss

= MSE +

LINREG

↑ ↓ overfitting

$$LR = 1 + 1.5f_1 + 0.2f_2 + 0 \cdot f_3$$

least imp feature - f_3

reduce this

$$MSE = \underbrace{\sum_{i=1}^m e_i}_{m} = \underbrace{\sum_{i=1}^m (y_i - \hat{y}_i)^2}_{m}$$

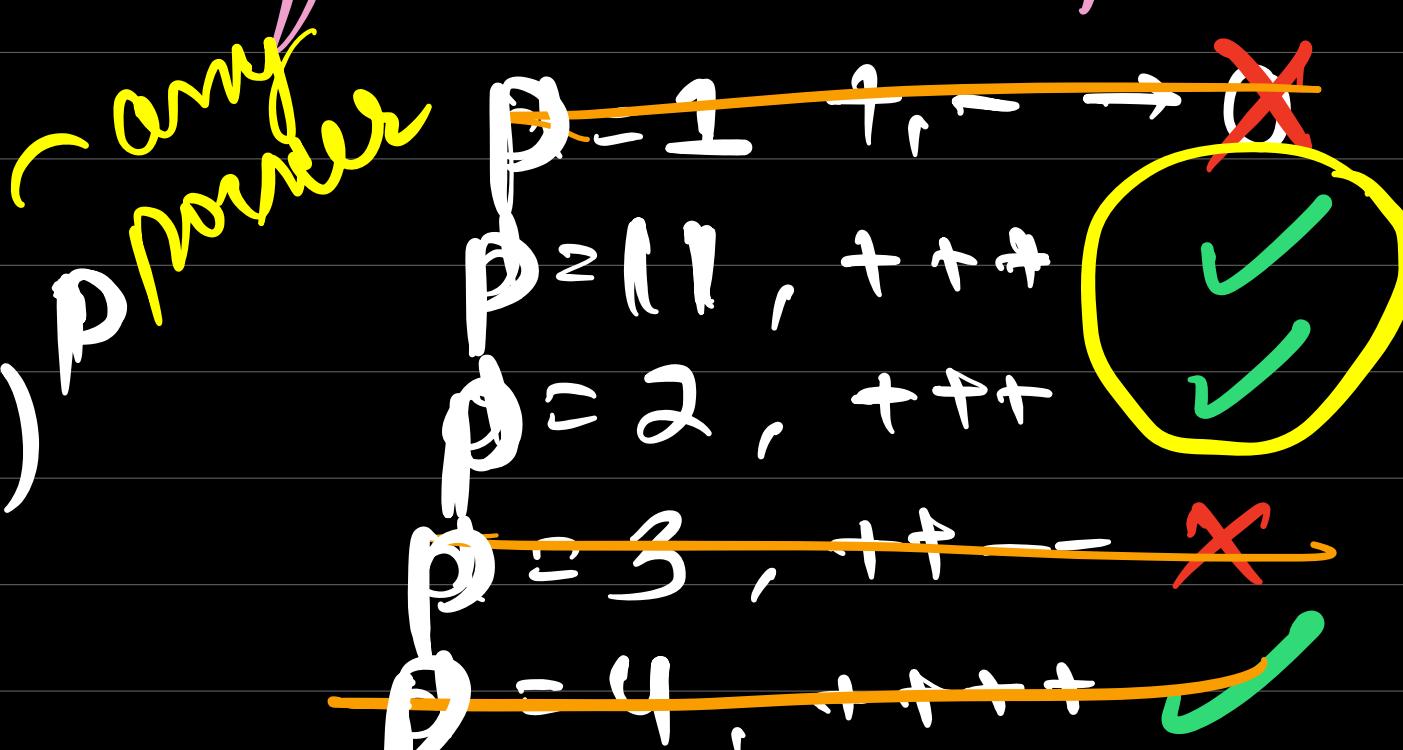
Goal 1: $\hat{y}_i \longleftrightarrow y_i$

Goal 2: Reduce the w's for less rup

features.

2, 4, 6, 8

$$\sum_{j=1}^d (w_j)^p$$



Reduce irrelevant w's so that $\sum w_j$ will
be zero

Error minimisation - three options

① MAE

$$\sum_{i=1}^m |e_i|$$

$|w_j|$

$$\left(\sum_{j=1}^d |w_j| \right)$$

② MSE

$$\sum_{i=1}^m |e_i|^2$$

$$\left(\sum_{j=1}^d (w_j)^2 \right)$$

To minimize.

option 1

option 2.

loss
update

$\text{MSE} + \text{Regularization}$
loss.

Minimize

Reduce
error

May ↑ overfitting

Reduce irrelevant
w's

↓ Overfitting

TUG OF WAR.

Find equilibrium.

$$L_1 - \sum_{j=1}^d |w_j| - \text{hasso}$$

Regularisation

$$L_2 - \sum_{j=1}^d (w_j)^2 - \text{Ridge.}$$

Regularisation

Q Should we give equal weightage
to MSE and Regressions?

Analogies

$$w_{\text{new}} = w_{\text{old}} + \alpha (\text{-Gradient})$$

Regularisation
Rate

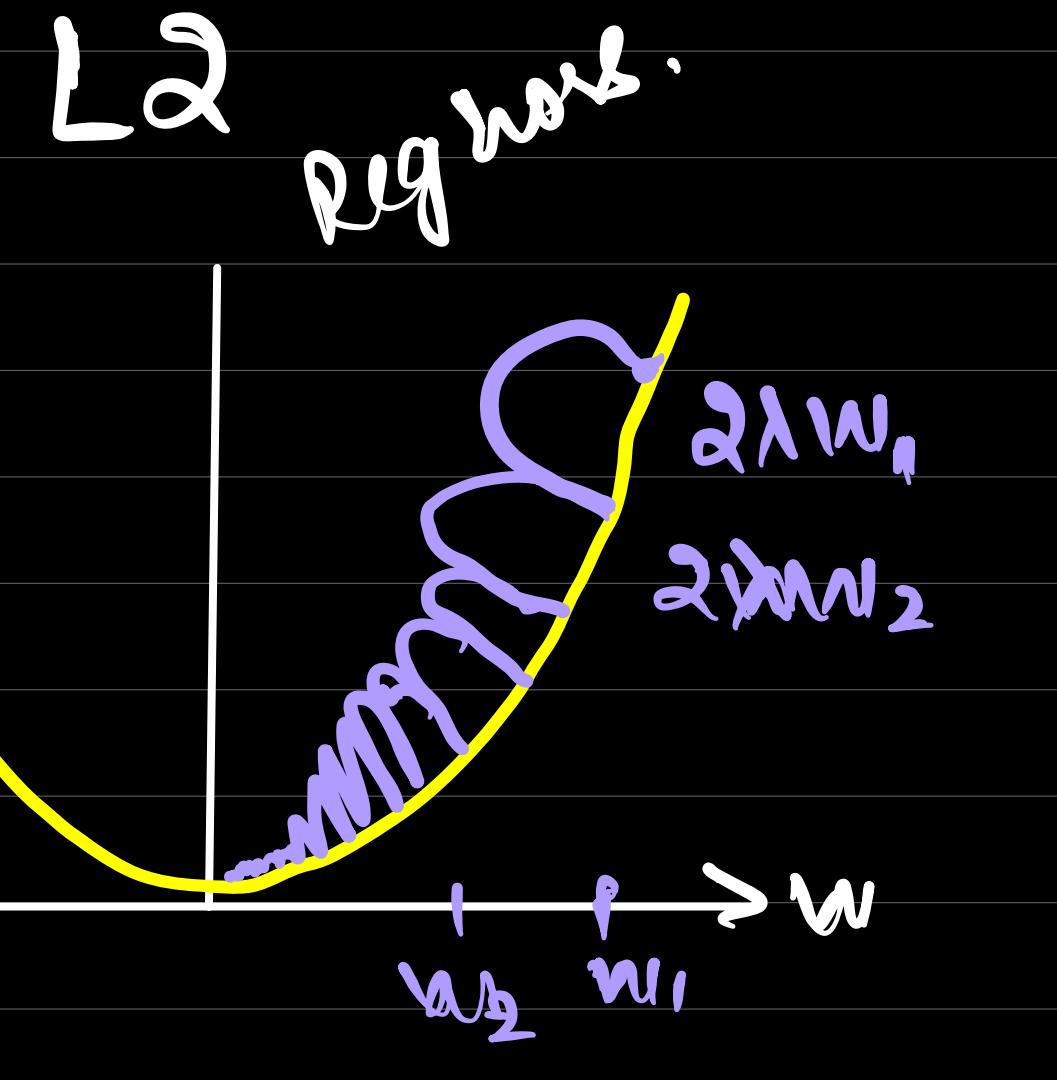
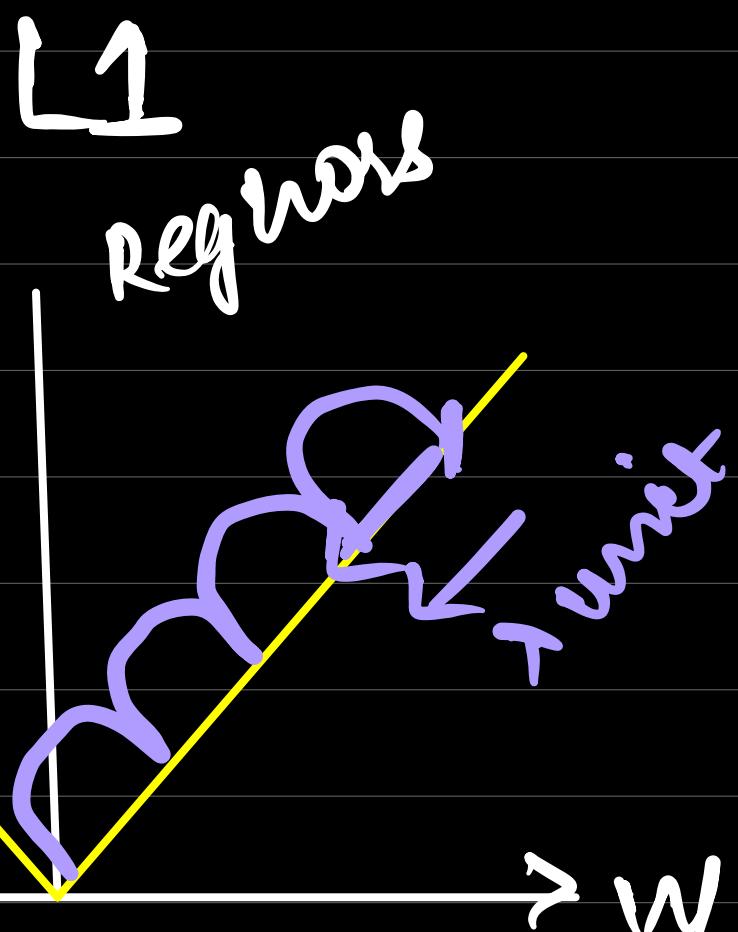
wf

$$L_{\text{new}} = \text{MSE} + \lambda \text{ Regularisation L}$$

Remember

Weight
Update

$$\text{Calculated } \frac{\partial L}{\partial w} = \frac{\partial \text{MSE}}{\partial w} + \lambda \frac{\partial \text{Reg L}}{\partial w}$$



$$w > 0 \quad \frac{\partial L}{\partial w} = 1$$

$$w < 0 \quad \frac{\partial L}{\partial w} = -1$$

$$w = 0 \quad \frac{\partial L}{\partial w} = 0$$

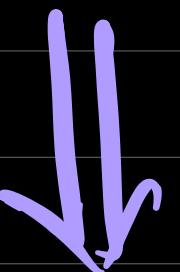
$$\frac{\partial L}{\partial w} = -2w$$

In case of L2, steps towards an update will become smaller & smaller.

⇒ L1 Regularisation leads to a sparse solution i.e., it has tendency to make less relevant

feature to go off by making
 $w_j = 0$

\Rightarrow it has property of making mag.
of w's overall small but doesn't
make anyone of them as zero.



feature selection - to remove all

use imp feature.

Use L1 Reg

Care about only and don't want
performance remove even less
imp features

Use L2 Reg

Use both L1 and L2

$$\text{Loss} = \text{MSE} + \lambda_1 \sum_{j=1}^d |w_j| + \lambda_2 \sum_{j=1}^d (w_j)^2$$

Elastic Net Regularisation.

Hyper-parameter Tuning

Parameters - $w_0, w_1, w_2, w_3, \dots, w_d$

Things which model learns
during training

Hyper-parameters - learning rate (η, α)
Regularization Rate
Degree of nonlinearity

- Try many values set all set by DB training the model.
NOT learnt from training, but from diff experiments

The process of selecting the right value of
Hyperparameter - Hyperparameter Tuning

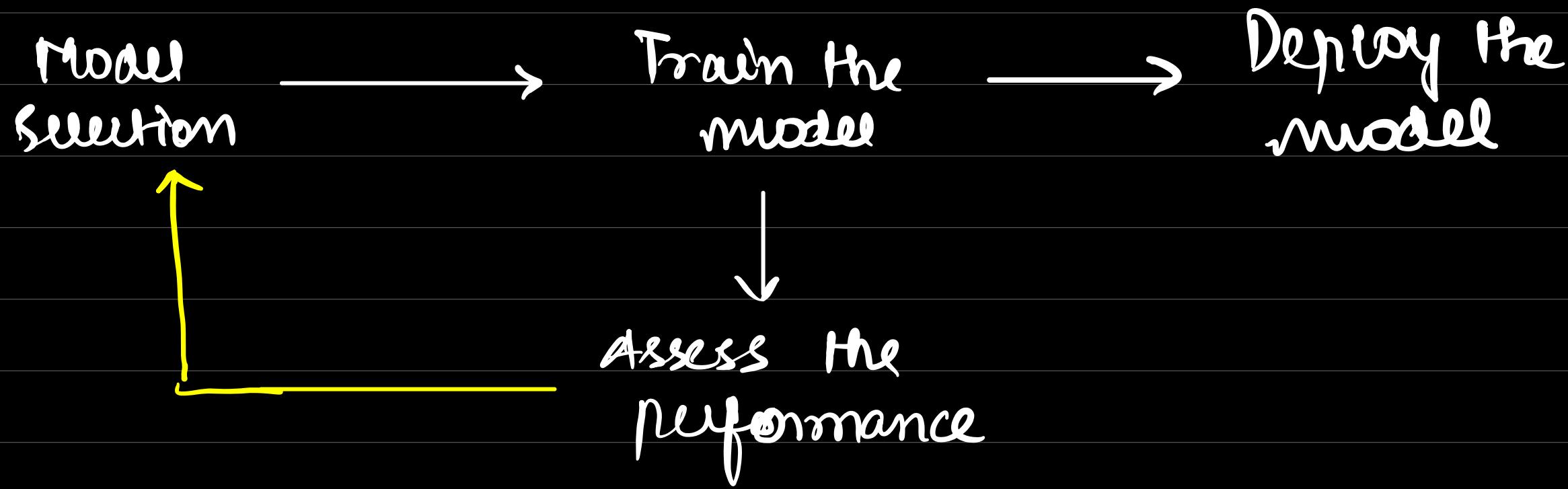
$$d = [1, 2, 3, 4, 5, 10, 20, 30]$$

for deg in d :

Train a model for degree = (deg)

Check the train & test performance

(at which ever ↑ test perf. - use that deg.)



Similar job for 'd, t, f'