- Data contains information, and we assume it is generated from some unknown underlying distribution (this is why we need ML).

- Machine Learning = learning patterns of this unknown distribution from data.

- Information = surprise:
    - If a data point was fully expected → low surprise → low information
    - If it reveals something new → high surprise → high information

- Seeing a data point and learning something new from it = gaining information.

- This idea of "surprise = information" is the foundation for information theory concepts like entropy (average surprise).

- The same idea applies to any data type: images, text, audio, etc.

- Many ML terms later (entropy, loss, likelihood) are built on this basic idea

Real world data has : Pattern and Noise, and pattern is not direclty inferrable becaue there is random noise . Our goal is to capture pattern.

In real life we don't know how data is genreated, for understanding purpose will assume we know true process, so that we can see what model is supposed to learn, for example: lets take sin function and add noise to it, (synthetic data)

Link: code and explaination for sine prediction

In reality, data comes from an unknown and potentially infinite population/distribution , but we only observe a finite dataset.
From this limited data, we must generalize to unseen data, not simply memorize the given samples.
This is inherently difficult because finite samples may not fully represent the true distribution .
Sampling data also involves uncertainty, and probability theory provides a mathematical framework to quantify this uncertainty.
Once we obtain a probabilistic representation of the data (i.e., uncertainty is modeled), we still need to make predictions or decisions.
Decision theory tells us how to make optimal predictions under uncertainty.

**Curve fitting** means:
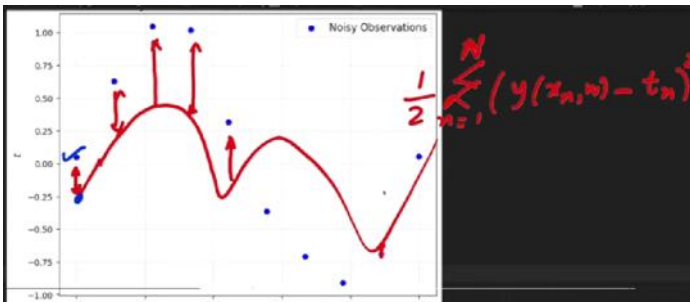Finding a function that best fits (approximates) the given data points.



there is only one original feature xxx; higher-dimensional feature space is created by applying polynomial basis functions x2,x3,...x^2, x^3 while y remains the target variable.
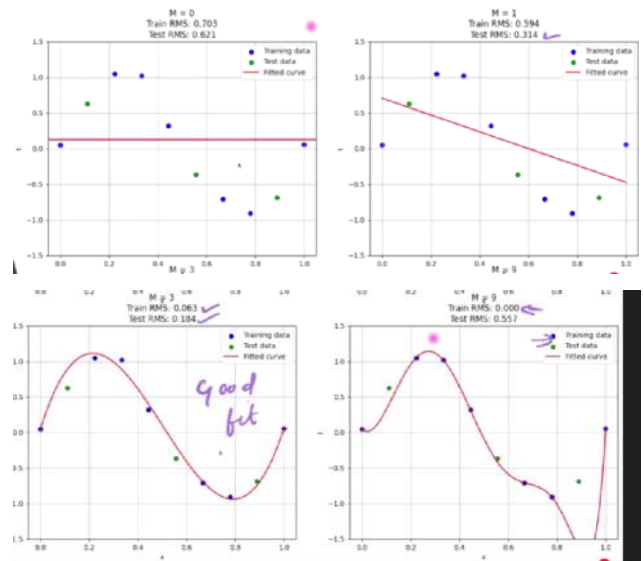


How do we know our curve is good?
we get W* , set of W for which error function is minized.
it should not memorize (less error during training, but sucks in testing), it should generalize



code



We can see, if we increase the number of parameters, its overfitting, so shall we reduce the number of paramters but it would limit th learning cabaplity

but isnt higher degree polynomial superset of of lower degree polynomial?
  ➢ It might fit the noise too.
  ➢ Goal is not find best fit in training , but to generalize
(the reason its having very high coefficient values, so if we can maintain or penalize it to have lower coeffiient values we are good with high parameters)







This is called regularization, M = no. of degree is hyperparameter, fixed during training

Regularization says: "Don't use large weights unless really necessary."

λ controls how much we punish complexity.

$$E(w) = \frac{1}{2}\sum_{n=1}^{N}\left(y(x_n, w) - t_n\right)^2 + \nearrow \frac{\lambda}{2}\|w\|^2$$

$\lambda, \quad M \quad \textcircled{w}$

$\uparrow \quad \uparrow$

$O$

Fix