

MULTI-LINEAR & POLYNOMIAL REGRESSION

LECTURE 1
SECTION 2
JUNE 1ST



INSTITUTE FOR APPLIED
COMPUTATIONAL SCIENCE
AT HARVARD UNIVERSITY

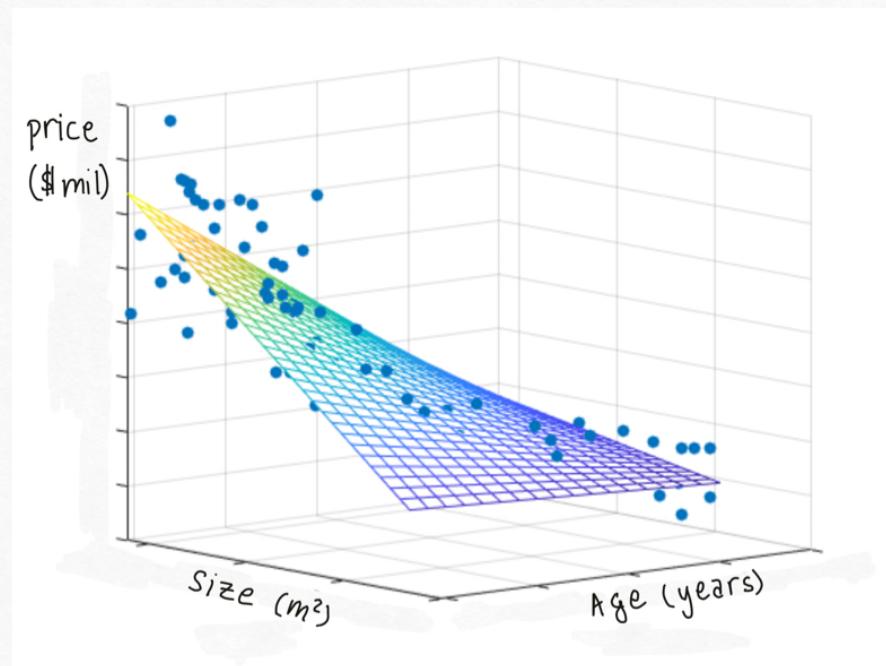


UNIVERSITY of
RWANDA

MULTILINEAR REGRESSION

MULTILINEAR REGRESSION:

What if the target y depends on multiple predictors?



A multi-linear model for y given $\mathbf{x} = [1, x_1, x_2, \dots, x_D]$ is:

$$y = w_D x_D + w_{D-1} x_{D-1} + \dots + w_1 x_1 + w_0$$

often we write this in vector notation:

$$\mathbf{y} = \mathbf{w}^T \mathbf{x}$$

where $\mathbf{w} = [w_0, w_1, \dots, w_D]$.

MODEL TRAINING: MULTI-LINEAR REGRESSION

Given a training dataset $\{(X^{(1)}, y^{(1)}), \dots, (X^{(N)}, y^{(N)})\}$ where

$$X_n = [X_1^{(n)}, X_2^{(n)}, \dots, X_D^{(n)}]$$

We find the optimal multi-linear model: $y = W^T X$

$$\text{Find } W^* = \underset{W}{\operatorname{argmin}} \mathcal{L}(W) = \underset{W}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N (y^{(n)} - W^T X^{(n)})^2$$

1. Compute the gradient: $\nabla_W \mathcal{L} = [\frac{\partial \mathcal{L}}{\partial w_0}, \frac{\partial \mathcal{L}}{\partial w_1}, \dots, \frac{\partial \mathcal{L}}{\partial w_D}]$

$$\begin{aligned}\nabla_W \mathcal{L} &= \frac{1}{N} \sum_{n=1}^N 2(y^{(n)} - W^T X^{(n)}) X^{(n)} \\ &= (y - W^T X) X\end{aligned}$$

where $y = [y_1, \dots, y_n]$ is the vector of target values.

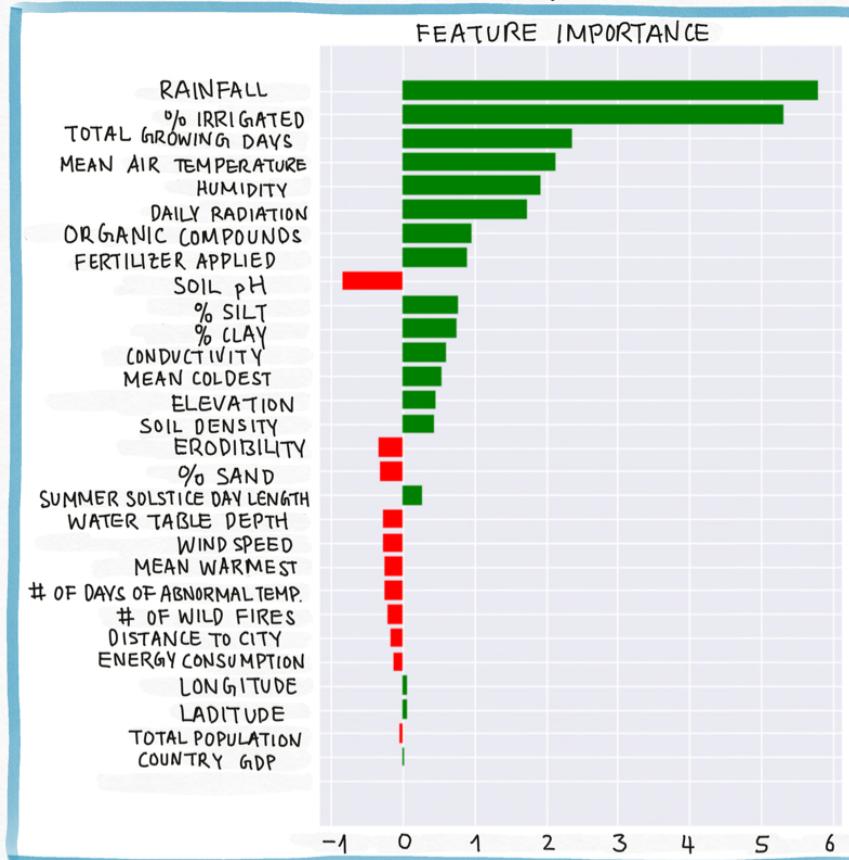
where $X = \begin{bmatrix} 1 & X_1^{(1)} & X_2^{(1)} & \cdots & X_D^{(1)} \\ 1 & X_1^{(2)} & X_2^{(2)} & \cdots & X_D^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & X_1^{(N)} & X_2^{(N)} & \cdots & X_D^{(N)} \end{bmatrix}$ is the matrix of predictor values

2. Solve for stationary points: $\nabla_W \mathcal{L} = 0$

$$W^* = (X^T X)^{-1} X^T y$$

INTERPRETING MULTI-LINEAR REGRESSION:

For linear models, it's important to interpret the model parameters.



When there are a large number of predictors: X_1, \dots, X_d , there will be a large number of model parameters: w_1, \dots, w_d .

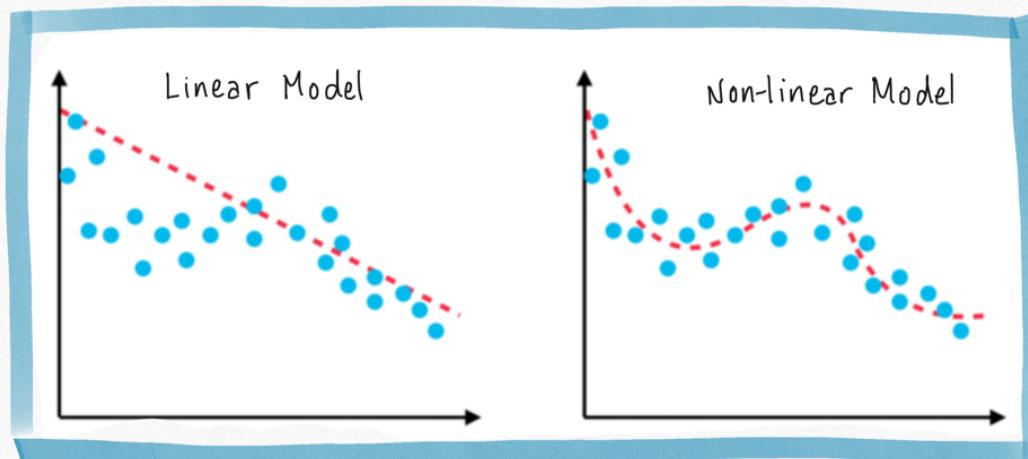
Looking at the values of w is impractical, so we visualize these values a **feature importance** graph.

The feature importance graph shows which predictor has the most impact on the model's prediction.

POLYNOMIAL REGRESSION

FITTING NONLINEAR DATA:

Multi-linear models can fit large datasets with many predictors. But the relationship between predictor and target isn't always linear.



We want a model:

$$y = f_w(x)$$

where f is a non-linear function, and w is a vector of the parameters of f .

POLYNOMIAL REGRESSION:

The simplest type of non-linear function is polynomial:

$$y = w_k x^k + w_{k-1} x^{k-1} + \dots + w_2 x^2 + w_1 x^1 + w_0$$

But this looks a lot like multi-linear regression where the predictors are powers of x !

Example: Selling price v.s. size

x (m^2)	y (\$ mil)
50	0.5
60	1.6
70	5.8
80	8.3



x	x^2	x^3	y
50	2500	12500	0.5
60	3600	21600	1.6
70	4900	34300	5.8
80	6400	51200	8.3

We transform the training data by adding powers of the single predictor x . Then fit a multi-linear model on the transformed data:

$$y = w_3 x^3 + w_2 x^2 + w_1 x + w_0$$

MODEL TRAINING: POLYNOMIAL REGRESSION

Given a dataset $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})\}$, we find the optimal polynomial model: $y = w_k x^k + w_{k-1} x^{k-1} + \dots + w_1 x + w_0$.

1. We transform the data by adding new predictors: $\tilde{X} = [1, \tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_k]$ where $\tilde{x}_k = x^k$.
2. Find the parameters $w^* = [w_0^*, w_1^*, \dots, w_k^*]$ by multi-linear regression

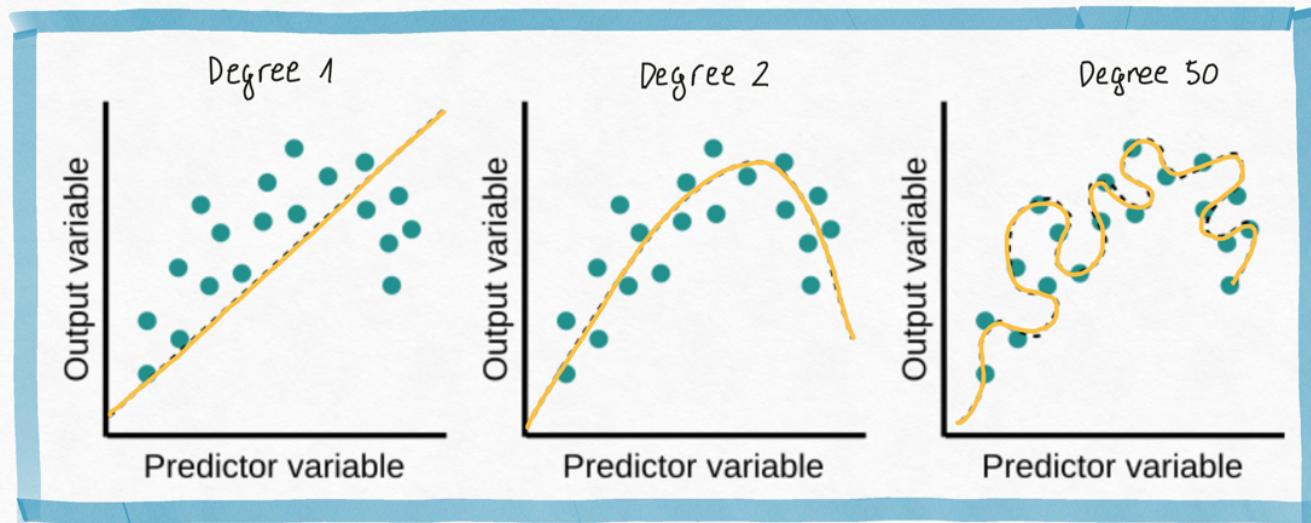
$$\begin{aligned} w^* &= \underset{w}{\operatorname{argmin}} \mathcal{L}(w) = \underset{w}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N (y^{(n)} - w^T \tilde{X}^{(n)})^2 \\ &= (\tilde{X}^T \tilde{X})^{-1} \tilde{X}^T y \end{aligned}$$

where $y = [y^{(1)}, \dots, y^{(N)}]$ is the vector of response values,

where $\tilde{X} = \begin{bmatrix} 1 & \tilde{x}_1^{(1)} & \tilde{x}_2^{(1)} & \dots & \tilde{x}_k^{(1)} \\ 1 & \tilde{x}_1^{(2)} & \tilde{x}_2^{(2)} & \dots & \tilde{x}_k^{(2)} \\ \vdots & & & & \\ 1 & \tilde{x}_1^{(N)} & \tilde{x}_2^{(N)} & \dots & \tilde{x}_k^{(N)} \end{bmatrix}$ is the matrix of predictor values in the transformed dataset.

MODEL SELECTION: CHOOSING THE DEGREE

Fitting a polynomial model requires choosing a degree.



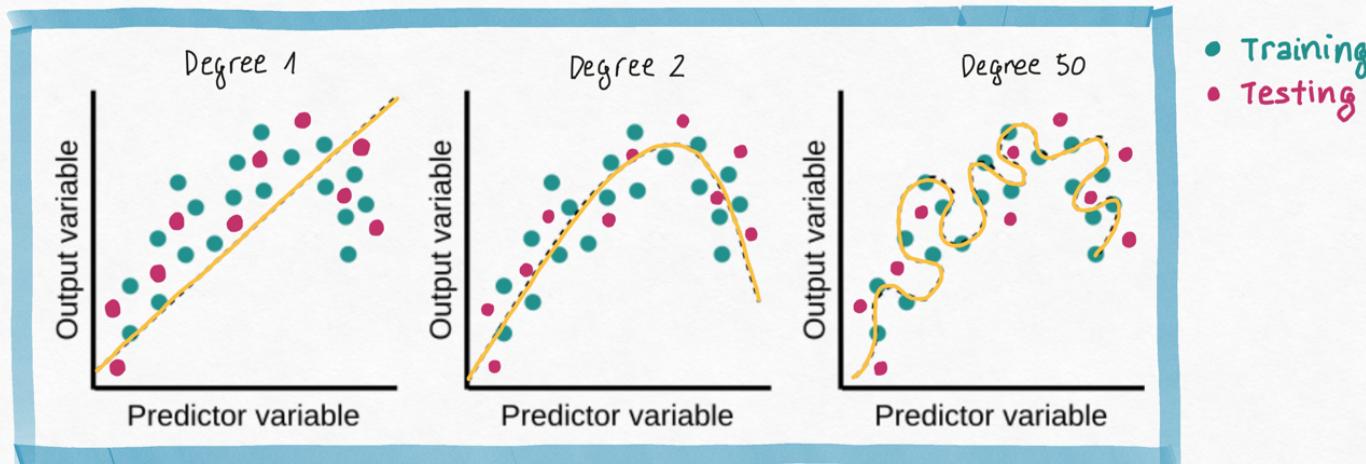
Underfitting: when the degree is too low, the model cannot fit the trend.

We want a model that fits the trend and ignores the noise?

Overfitting: when the degree is too high, the model fits all the noisy data points

OVERTFITTING OR UNDERFITTING?

If the dataset contains more than one predictor then we cannot easily visualize the polynomial $y = f(x)$.



Train MSE: 2.1

Test MSE: 1.9

When we underfit,
the train and test
errors are similar
and both large.

Train MSE: 0.5

Test MSE: 0.65

When we fit the
trend correctly
both train and
test error "should"
be small.

Train MSE: 0

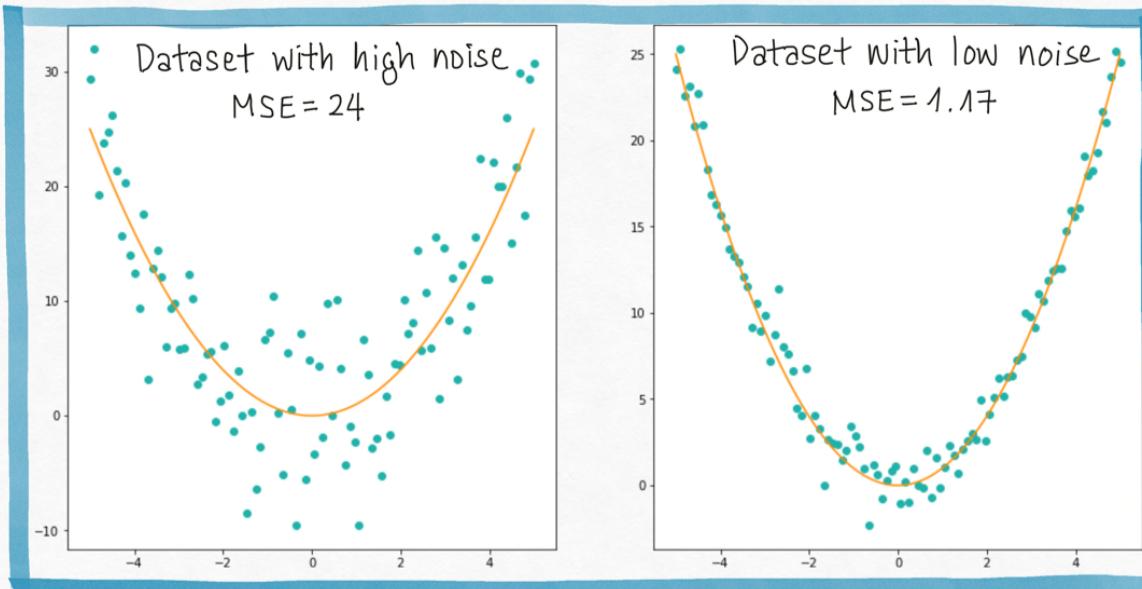
Test MSE: 0.9

When we overfit,
the train error is
very small, and test
error is much larger

UNDERFITTING OR HIGH NOISE?

We said that an ideal model would have similar train and test MSE, and both would be small; an underfitting model would have similar train and test MSE but both would be large.

Is this always true?



MSE is not a reliable way to detect underfitting, we need another method that accounts for noise in data.