

"Polynomial Regression"

Polynomial Regression Overview:

given data $X \in \mathbb{R}^{n \times m}$ (with n samples and m features) and target $y \in \mathbb{R}^n$, polynomial regression try to make a model

like: $y \approx \phi(x) \cdot \theta$ where $\phi(x)$ is a transformation

of x to include polynomial features to a specified degree d . θ is the parameter vector to learn.

OLS with Single column x [$x \in \mathbb{R}^{n \times 1}$]

Step-1 Polynomial expansion for degree d .

$$\phi(x) = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ 1 & x_2 & x_2^2 & \dots & x_2^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix} = \phi(x) \in \mathbb{R}^{n \times (d+1)}$$

Step-2

hypothesis for a single row of x

$$\hat{y} = \sum_{j=0}^d \theta_j \cdot x_j$$

now for n samples MSE will be

$$MSE(\theta) = \frac{1}{n} \sum_{i=0}^n (\tilde{y}_i - \hat{\tilde{y}}_i)^2$$

$$MSE(\theta) = \frac{1}{n} \sum_{i=0}^n \left(\tilde{y}_i - \sum_{j=0}^d \theta_j x_i^j \right)^2$$

Step-3 Define total squared error (remove $\frac{1}{n}$)

$\frac{1}{n}$ is a constant factor and does not affect the location of minimum while differentiating.

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^d \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^d \end{bmatrix} [n \times (d+1)]$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} [1 \times (d+1)]$$

Predicted value for each x_i

$$\hat{y}_i = \phi_i \theta = \sum_{j=0}^d \theta_j \cdot x_i^j$$

so all prediction (n rows)

$$\hat{y} = \Phi \theta$$

new error (residuals)

$$e = y - \hat{y} = y - \Phi \theta$$

Step-4 sum of squared errors

$$J(\theta) = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \phi_i \theta)^2 = \|y - \Phi \theta\|^2$$

Step-5 expand $J(\theta)$

$$J(\theta) = \|y - \Phi \theta\|^2 = (y - \Phi \theta)^T \cdot (y - \Phi \theta) \quad [\|a\| = a^T \cdot a]$$

$$\therefore J(\theta) = (y^T - \theta^T \Phi^T) \cdot (y - \Phi \theta)$$

$$J(\theta) = y^T y - \theta^T \Phi^T y - y^T \Phi \theta + \theta^T \Phi^T \Phi \theta$$

$$J(\theta) = y^T y - 2\theta^T \Phi^T y + \theta^T \Phi^T \Phi \theta \quad \left[\Phi^T \Phi^T y \text{ and } y^T \Phi \Phi \text{ is symmetric matrix} \right]$$

Multivariate polynomial Regression using OLS:

Let, data input matrix: $X \in \mathbb{R}^{n \times d}$ where n is the number of rows and d is total columns in X .

Target: $y \in \mathbb{R}^{n \times 1} \rightarrow$ cause n rows has n target.

Desired polynomial degree p .

Note: we will expand the X matrix by adding $(p+1)$ polynomial term. Then we will apply linear regression on expanded X .

Step-1 polynomial Feature Transform

we want to build all polynomial terms of features up to degree p using combination.

For example: x has $d=2$ column and degree $p=2$

$$\text{For 1 sample, } \Phi(X) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

For n col and p degree number of combination = $\frac{(n+p)!}{n! p!}$

For n samples $\Phi(X) \in \mathbb{R}^{n \times m}$

where m is the polynomial terms, $m = \binom{d+p}{p}$

Step-2 Define the linear model using $\Phi(x)$

hypothesis, $\hat{y} = \Phi(x) \cdot \theta$ where:

$\Phi(x) \rightarrow$ transformed x

$\theta \in \mathbb{R}^{m \times 1}$ parameter vector

Step-3 Define the MSE cost function

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

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \Phi(x_i)^T \theta)^2$$

$$= \|y - \Phi(x) \cdot \theta\|^2$$

$$= (y - \Phi(x) \cdot \theta)^T \cdot (y - \Phi(x) \cdot \theta) \left[(a-b)^T \cdot (a-b) = a^T a - 2a^T b + b^T b \right]$$

$$= y^T y - 2y^T \Phi(x) \cdot \theta + \theta^T \Phi(x)^T \Phi(x) \theta$$

Final cost function $\rightarrow J(\theta) = y^T y - 2y^T (X_{poly}) \cdot \theta + \theta^T (X_{poly})^T X_{poly} \cdot \theta$

Step-4 Take gradient with respect θ to minimize $J(\theta)$

$$\nabla_{\theta} (y^T y) = 0$$

$$\nabla_{\theta} (-2y^T X_{poly} \theta) = -2y^T X_{poly}$$

$$\nabla_{\theta} (\theta^T X_{poly}^T X_{poly} \theta) = 2X_{poly}^T X_{poly} \theta$$

Total gradient

$$\nabla_{\theta} J(\theta) = -2X_{\text{poly}} y^T + 2X_{\text{poly}}^T X_{\text{poly}} \theta$$

To minimize $J(\theta)$ set gradient to 0

$$-2X_{\text{poly}} y^T + 2X_{\text{poly}}^T X_{\text{poly}} \theta = 0$$

$$\Rightarrow X_{\text{poly}}^T X_{\text{poly}} \theta = X_{\text{poly}} y^T$$

$$\Rightarrow \theta = \left(X_{\text{poly}}^T X_{\text{poly}} \right)^{-1} X_{\text{poly}} y^T$$

"Polynomial Regression using Gradient Descent"

single feature input:

Step-1 Transform feature x

Let, $x \in \mathbb{R}^{n \times 1} \rightarrow$ a single column input.

and polynomial degree = 2

$$X_{\text{poly}} = \phi(x) = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix} \in \mathbb{R}^{n \times 3}$$

Take the hypothesis: $\hat{y} = \phi(x) \theta$ where $\hat{y} \in \mathbb{R}^{n \times 1}$, $\theta \in \mathbb{R}^{3 \times 1}$

Step-2 Define the cost function (MSE)

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

$$= \frac{1}{2n} \sum_{i=1}^n (X_{poly}^i \theta - y_i)^2$$

$$J(\theta) = \frac{1}{2n} \|X_{poly} \theta - y\|^2$$

$\phi(x)$ and X_{poly}
is same thing

Step-3 Gradient of the cost function

$$\begin{aligned} \nabla_{\theta}(J(\theta)) &= \frac{1}{2n} \cdot 2 \phi(x)^T (\phi(x) \cdot \theta - y) \\ &= \frac{1}{n} X_{poly}^T (X_{poly} \theta - y) \end{aligned} \quad \left\{ \begin{array}{l} \nabla_{\theta} \|A\theta - y\|^2 \\ = 2A^T (A\theta - y) \end{array} \right.$$

Step-4 update rule

$$\theta := \theta - \alpha \cdot \nabla_{\theta}(J(\theta))$$

$$\theta := \theta - \alpha \left[\frac{1}{n} X_{poly}^T (X_{poly} \theta - y) \right]$$

Note $X_{poly} \cdot \theta - y \rightarrow$ is the vector of prediction error.
shape $(n \times 1)$

Multi-feature Input (using gradient descent).

let $X \in \mathbb{R}^{n \times d} \rightarrow$ the original data input

after polynomial transformation

$X_{\text{poly}} = \phi(X) \in \mathbb{R}^{n \times m}$ where $n \rightarrow$ number of rows
 $m \rightarrow$ includes all the combinations up to degree p

$$m = \frac{(c+p)!}{c! p!}$$

\rightarrow where $c \rightarrow$ number of columns in X before transformation

$p \rightarrow$ number of degree.

Define the cost function

MSE cost function

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

$$= \frac{1}{2n} \sum_{i=1}^n (\phi(x) \cdot \theta - y)^2 \quad [\hat{y} = \phi(x) \cdot \theta]$$

$$= \frac{1}{2n} \|\phi(x) \cdot \theta - y\|^2$$

calculate the gradient:

$$\nabla_{\theta} (J(\theta)) = \nabla_{\theta} \left(\frac{1}{2n} \|\phi(x) \cdot \theta - y\|^2 \right)$$

$$= \frac{1}{n} [\phi(x)]^T (\phi(x) \cdot \theta - y)$$

update theta

$$\theta := \theta - \alpha \cdot \nabla_{\theta} (J(\theta))$$

learning rate