

MLT: Week 5

Linear Regression: Least Squares and Kernel Regression

Sherry Thomas

Linear Regression Algorithm

Linear Regression Algorithm

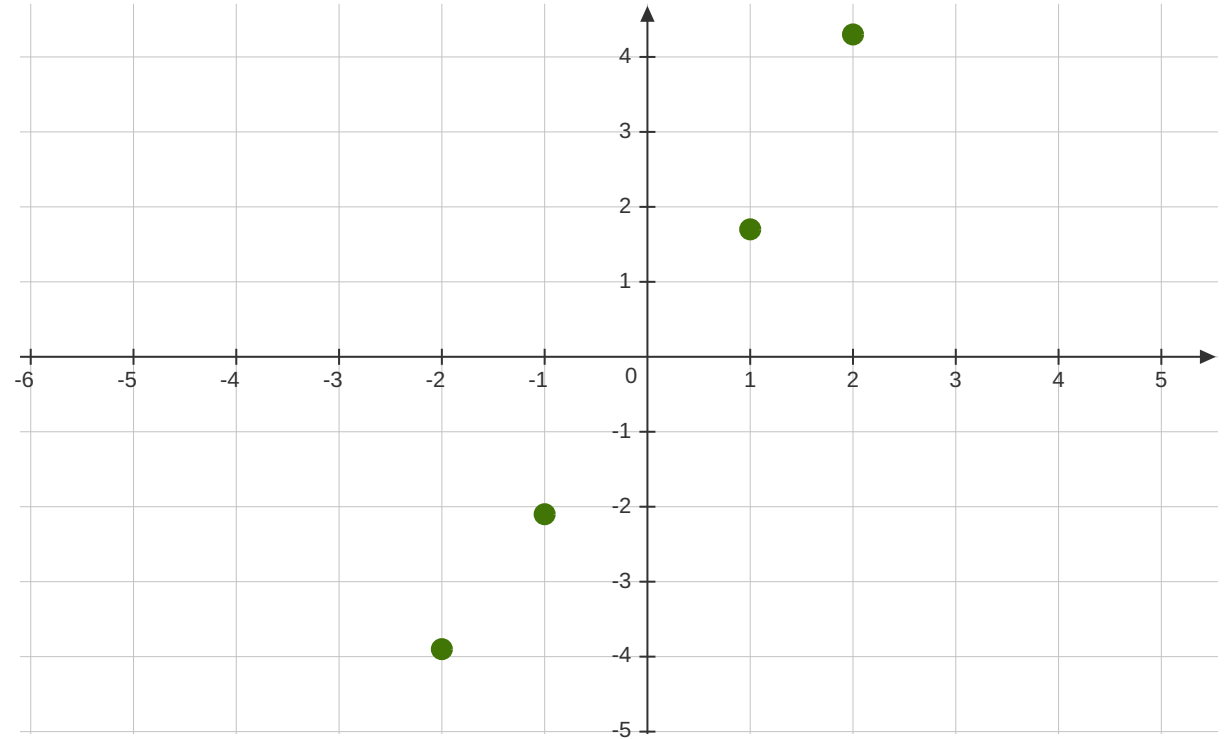
Given a dataset $\{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^d$,
let $\{y_1, \dots, y_n\}$ be the labels, where $y_i \in \mathbb{R}$.

$$X = \begin{bmatrix} -2 & -1 & 1 & 2 \end{bmatrix} \quad y = \begin{bmatrix} -3.9 \\ -2.1 \\ 1.7 \\ 4.3 \end{bmatrix}$$

Linear Regression Algorithm

Given a dataset $\{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^d$,
let $\{y_1, \dots, y_n\}$ be the labels, where $y_i \in \mathbb{R}$.

$$X = \begin{bmatrix} -2 & -1 & 1 & 2 \end{bmatrix} \quad y = \begin{bmatrix} -3.9 \\ -2.1 \\ 1.7 \\ 4.3 \end{bmatrix}$$



Linear Regression Loss Function

Linear Regression Loss Function

The most commonly used loss function for Regression is the Squared Sum Error. It is given by,

$$\sum_{i=1}^n ||\mathbf{w}^T \mathbf{x}_i - y_i||_2^2 \quad \forall i$$

Linear Regression Loss Function

The most commonly used loss function for Regression is the Squared Sum Error. It is given by,

$$\sum_{i=1}^n ||\mathbf{w}^T \mathbf{x}_i - y_i||_2^2 \quad \forall i$$

We use the line $y = x$ as a baseline to better understand the loss function.

In the above line, \mathbf{w} is represented as $\mathbf{w} = [1]$. Using this in the loss function, we obtain,

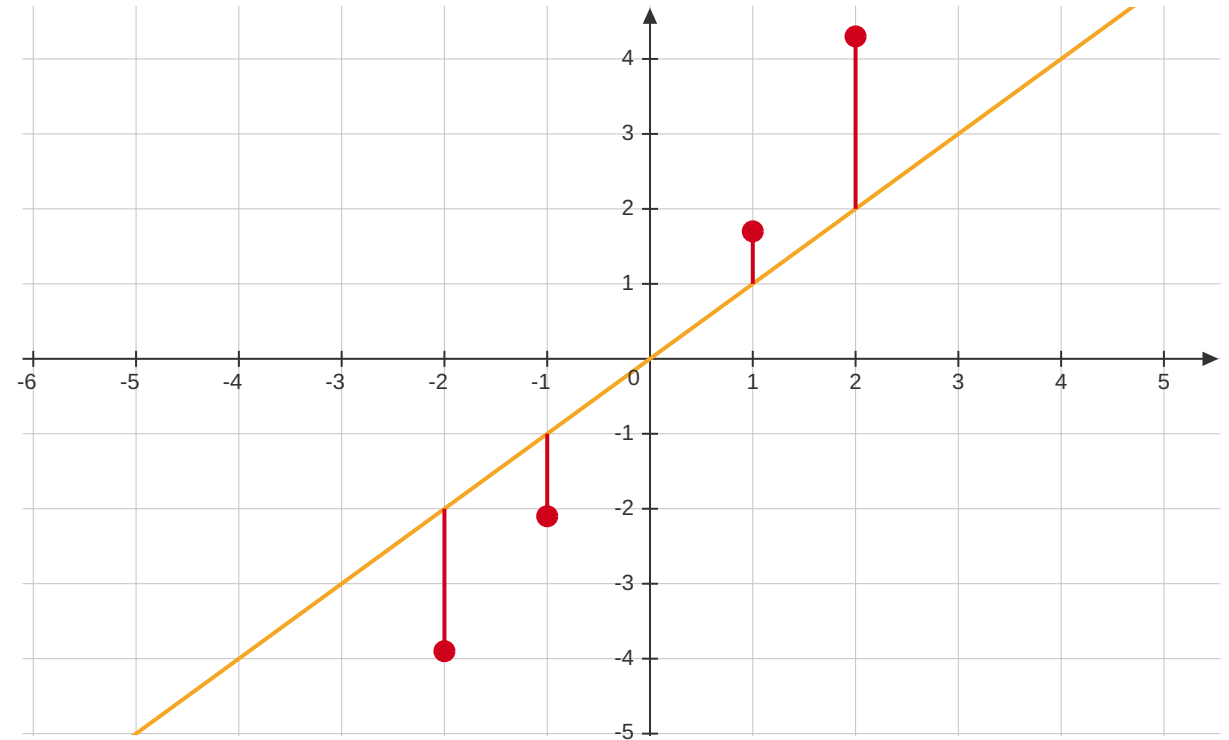
Linear Regression Loss Function

The most commonly used loss function for Regression is the Squared Sum Error. It is given by,

$$\sum_{i=1}^n ||w^T x_i - y_i||_2^2 \quad \forall i$$

We use the line $y = x$ as a baseline to better understand the loss function.

In the above line, w is represented as $w = [1]$. Using this in the loss function, we obtain,



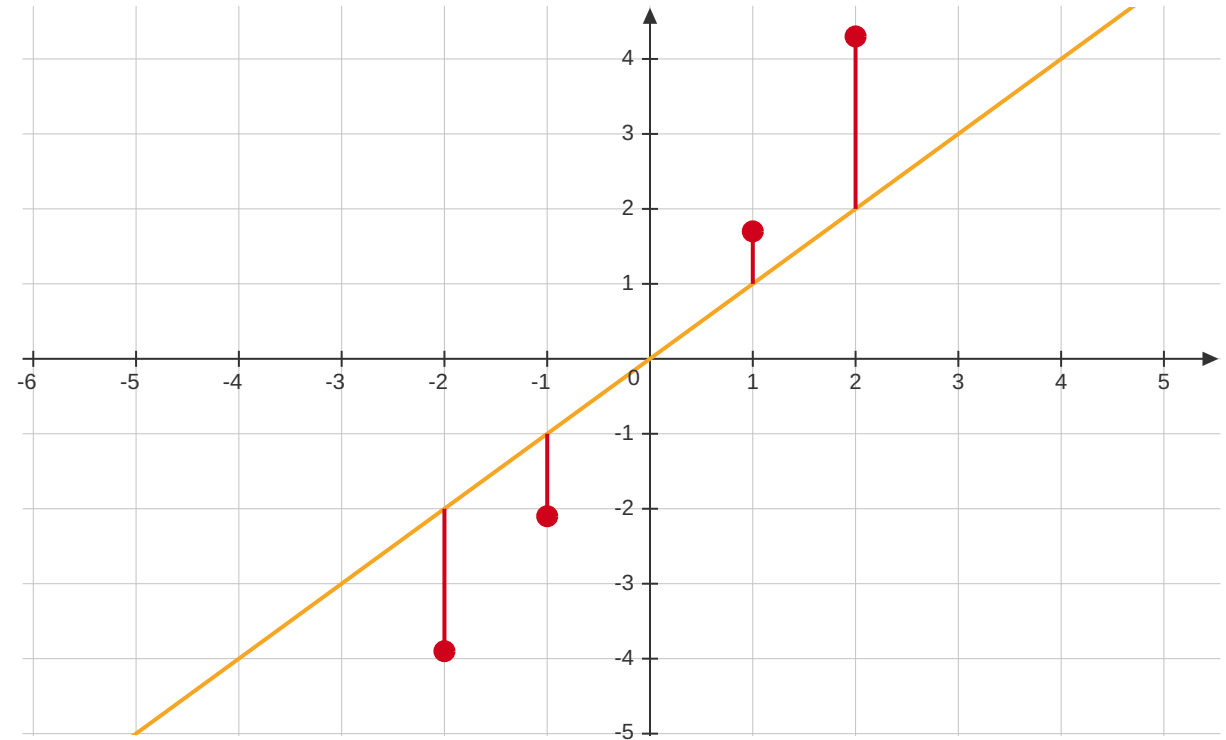
Linear Regression Loss Function

The most commonly used loss function for Regression is the Squared Sum Error. It is given by,

$$\sum_{i=1}^n ||\mathbf{w}^T \mathbf{x}_i - y_i||_2^2 \quad \forall i$$

We use the line $y = x$ as a baseline to better understand the loss function.

In the above line, \mathbf{w} is represented as $\mathbf{w} = [1]$. Using this in the loss function, we obtain,



$$||\mathbf{X}^T \mathbf{w} - \mathbf{y}||_2^2 = ||[1 \ 2 \ -1 \ -2]^T [1] - [1.7 \ 4.3 \ -2.1 \ -3.9]^T||_2^2 = ||[0.7 \ -2.3 \ 1.1 \ 1.9]^T||_2^2 = 3.25576^2 = 10.6$$

Optimizing the Error Function

Optimizing the Error Function

The minimization equation can be rewritten in the vectorized form as,

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|_2^2$$

Let this be a function of \mathbf{w} as follows:

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|_2^2$$

$$f(\mathbf{w}) = \frac{1}{2} (\mathbf{X}^T \mathbf{w} - \mathbf{y})^T (\mathbf{X}^T \mathbf{w} - \mathbf{y})$$

$$\therefore \nabla f(\mathbf{w}) = (\mathbf{X}\mathbf{X}^T)\mathbf{w} - (\mathbf{X}\mathbf{y})$$

Optimizing the Error Function

The minimization equation can be rewritten in the vectorized form as,

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|_2^2$$

Let this be a function of \mathbf{w} as follows:

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}^T \mathbf{w} - \mathbf{y}\|_2^2$$

$$f(\mathbf{w}) = \frac{1}{2} (\mathbf{X}^T \mathbf{w} - \mathbf{y})^T (\mathbf{X}^T \mathbf{w} - \mathbf{y})$$

$$\therefore \nabla f(\mathbf{w}) = (\mathbf{X}\mathbf{X}^T)\mathbf{w} - (\mathbf{X}\mathbf{y})$$

Setting the above equation to zero, we get

$$(\mathbf{X}\mathbf{X}^T)\mathbf{w} - (\mathbf{X}\mathbf{y}) = 0$$

$$(\mathbf{X}\mathbf{X}^T)\mathbf{w}^* = \mathbf{X}\mathbf{y}$$

$$\therefore \mathbf{w}^* = (\mathbf{X}\mathbf{X}^T)^+ \mathbf{X}\mathbf{y}$$

where $(\mathbf{X}\mathbf{X}^T)^+$ represents the pseudo-inverse of $\mathbf{X}\mathbf{X}^T$.

Using Least Squares Regression

Using Least Squares Regression

To account for the bias term, we augment our feature matrix with an additional row of ones.

Hence, our current dataset is,

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -2 & -1 & 1 & 2 \end{bmatrix}$$

Using least squares regression to obtain the weight vector, we get,

$$\begin{aligned} w &= (XX^T)^+ Xy \\ &= \left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \right)^+ \begin{bmatrix} 1 & 1 & 1 & 1 \\ -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} -3.9 \\ -2.1 \\ 1.7 \\ 4.3 \end{bmatrix} \\ w &= \left(\begin{bmatrix} 4 & 0 \\ 0 & 10 \end{bmatrix} \right)^+ \begin{bmatrix} 0 \\ 20.2 \end{bmatrix} = \begin{bmatrix} 0.25 & 0 \\ 0 & 0.1 \end{bmatrix} \begin{bmatrix} 0 \\ 20.2 \end{bmatrix} = \begin{bmatrix} 0 \\ 2.02 \end{bmatrix} \end{aligned}$$

Using Gradient Descent

Using Gradient Descent

Given that w^* is the solution of an unconstrained optimization problem, we can solve it using gradient descent, where:

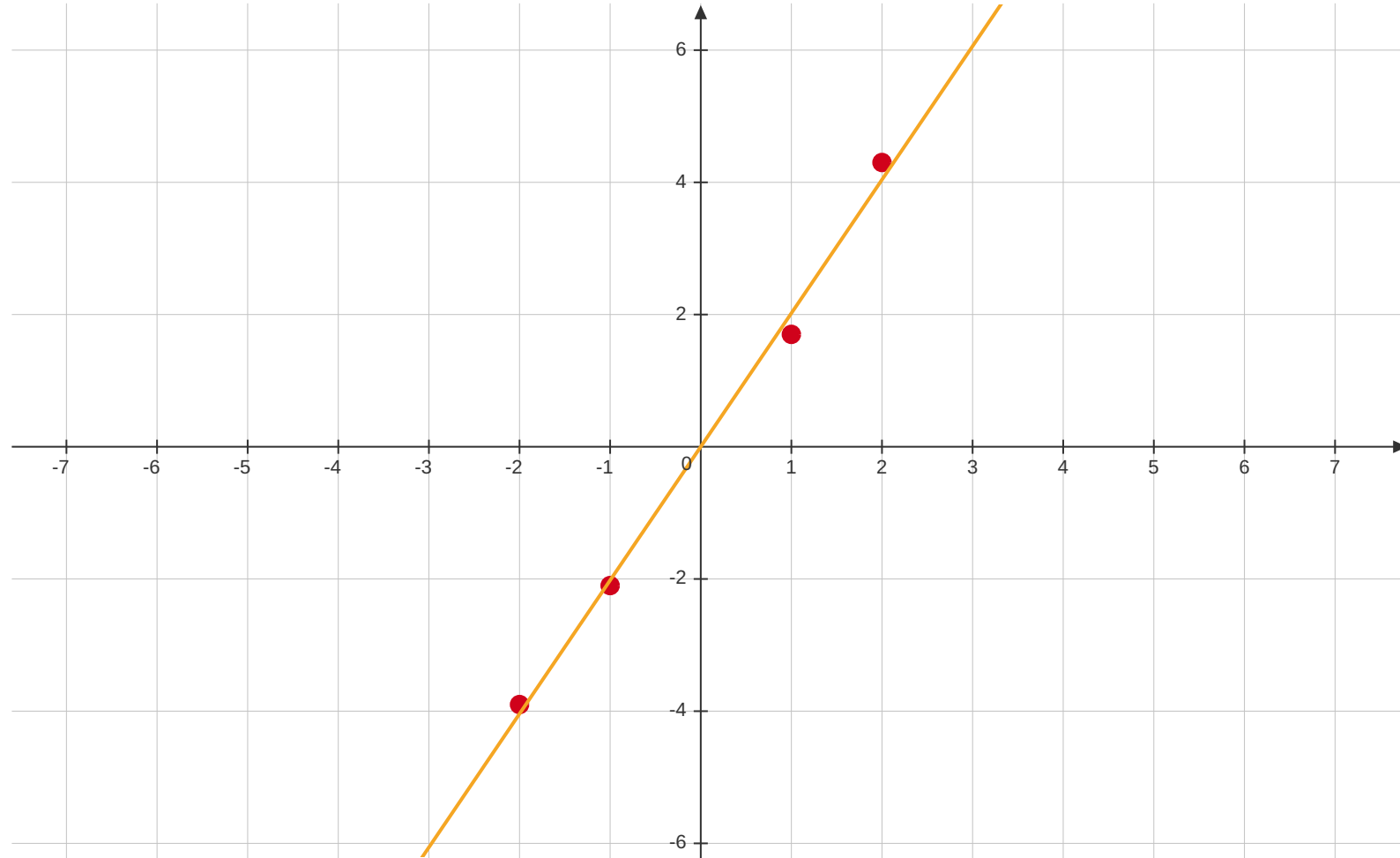
$$\begin{aligned}w^{t+1} &= w^t - \eta^t \nabla f(w^t) \\ \therefore w^{t+1} &= w^t - \eta^t [(XX^T)w^t - (Xy)]\end{aligned}$$

Here, η is a scalar used to control the step size of the descent, and t is the current iteration.

Using gradient descent to obtain the weight vector, we have:

$$\begin{aligned}w^1 &= w^0 - \eta^0 [(XX^T)w^0 - (Xy)] \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \times \left[\left(\begin{bmatrix} 1 & 1 & 1 & 1 \\ -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \right) \begin{bmatrix} 0 \\ 0 \end{bmatrix} - \begin{bmatrix} 1 & 1 & 1 & 1 \\ -2 & -1 & 1 & 2 \end{bmatrix} \begin{bmatrix} -3.9 \\ -2.1 \\ 1.7 \\ 4.3 \end{bmatrix} \right] \\ &= \begin{bmatrix} 0 \\ 0 \end{bmatrix} - 0.1 \times \left[- \begin{bmatrix} 0 \\ 20.2 \end{bmatrix} \right] \\ w^1 &= \begin{bmatrix} 0 \\ 2.02 \end{bmatrix}\end{aligned}$$

Regression Line using Gradient Descent



Kernel Regression Algorithm

Kernel Regression Algorithm

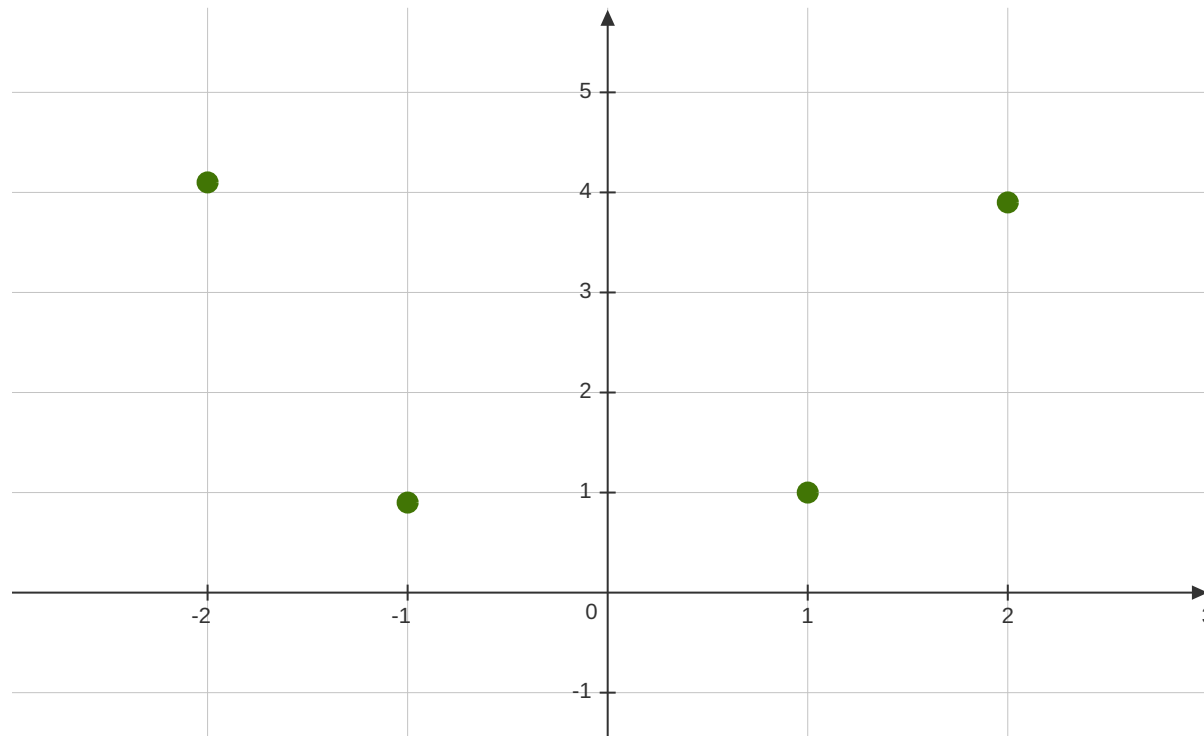
Given a dataset $\{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^d$,
let $\{y_1, \dots, y_n\}$ be the labels, where $y_i \in \mathbb{R}$.

$$X = \begin{bmatrix} 1 & 2 & -1 & -2 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 3.9 \\ 0.9 \\ 4.1 \end{bmatrix}$$

Kernel Regression Algorithm

Given a dataset $\{x_1, \dots, x_n\}$ where $x_i \in \mathbb{R}^d$,
let $\{y_1, \dots, y_n\}$ be the labels, where $y_i \in \mathbb{R}$.

$$X = \begin{bmatrix} 1 & 2 & -1 & -2 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 3.9 \\ 0.9 \\ 4.1 \end{bmatrix}$$



Using Kernel Regression

Using Kernel Regression

Let $w^* = X\alpha^*$ for some $\alpha^* \in \mathbb{R}^n$.

$$X\alpha^* = w^*$$

$$\therefore X\alpha^* = (XX^T)^+Xy$$

$$(XX^T)X\alpha^* = (XX^T)(XX^T)^+Xy$$

$$(XX^T)X\alpha^* = Xy$$

$$X^T(XX^T)X\alpha^* = X^TXy$$

$$(X^TX)^2\alpha^* = X^TXy$$

$$K^2\alpha^* = Ky$$

$$\therefore \alpha^* = K^{-1}y$$

where $K \in \mathbb{R}^{n \times n}$ and K can be also obtained using a kernel function like the Polynomial Kernel or RBF Kernel

Let's use the polynomial kernel of degree of two. By applying the kernel function to the dataset, we obtain,

Let's use the polynomial kernel of degree of two. By applying the kernel function to the dataset, we obtain,

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^2$$

$$k(\mathbf{x}_1, \mathbf{x}_1) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \right)^2 = (2 + 1)^2 = 9$$

$$k(\mathbf{x}_1, \mathbf{x}_2) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \right)^2 = (6 + 1)^2 = 16$$

$$\therefore \mathbf{K} = \begin{bmatrix} 9 & 16 & 1 & 16 \\ 16 & 36 & 0 & 36 \\ 1 & 0 & 9 & 0 \\ 16 & 36 & 0 & 36 \end{bmatrix}$$

Let's use the polynomial kernel of degree of two. By applying the kernel function to the dataset, we obtain,

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^2$$

We find α using the following equation,

$$k(\mathbf{x}_1, \mathbf{x}_1) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \right)^2 = (2 + 1)^2 = 9$$

$$k(\mathbf{x}_1, \mathbf{x}_2) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \right)^2 = (6 + 1)^2 = 16$$

$$\therefore \mathbf{K} = \begin{bmatrix} 9 & 16 & 1 & 16 \\ 16 & 36 & 0 & 36 \\ 1 & 0 & 9 & 0 \\ 16 & 36 & 0 & 36 \end{bmatrix}$$

Let's use the polynomial kernel of degree of two. By applying the kernel function to the dataset, we obtain,

$$k(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^T \mathbf{x}_j + 1)^2$$

$$k(\mathbf{x}_1, \mathbf{x}_1) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} + 1 \right)^2 = (2 + 1)^2 = 9$$

$$k(\mathbf{x}_1, \mathbf{x}_2) = \left(\begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} + 1 \right)^2 = (6 + 1)^2 = 16$$

$$\therefore K = \begin{bmatrix} 9 & 16 & 1 & 16 \\ 16 & 36 & 0 & 36 \\ 1 & 0 & 9 & 0 \\ 16 & 36 & 0 & 36 \end{bmatrix}$$

We find α using the following equation,

$$\alpha = K^{-1}y$$

$$\alpha = \begin{bmatrix} 9 & 16 & 1 & 16 \\ 16 & 36 & 0 & 36 \\ 1 & 0 & 9 & 0 \\ 16 & 36 & 0 & 36 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 3.9 \\ 0.9 \\ 4.1 \end{bmatrix}$$
$$\alpha = \begin{bmatrix} -0.1813 \\ 0.1707 \\ -0.1798 \\ 0.1737 \end{bmatrix}$$

Finding the Predictions using α

Finding the Predictions using α

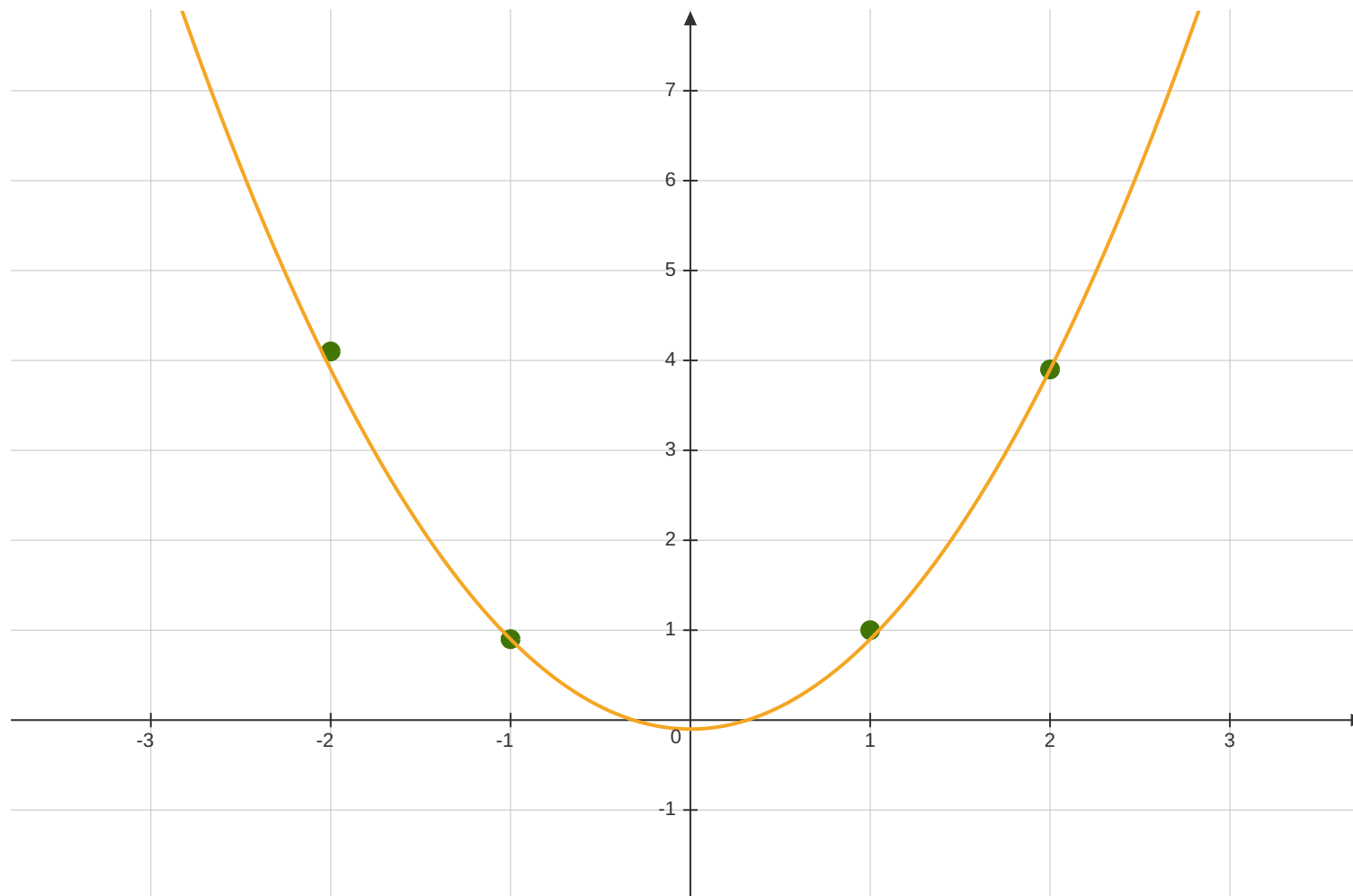
We use the following equation to get y_{pred}

$$\mathbf{w}^* \phi(\mathbf{X}) = \sum_{i=1}^n \alpha_i^* k(\mathbf{x}_i, \mathbf{x}_i)$$

where α_i^* gives the importance of the i^{th} datapoint towards \mathbf{w}^* .

$$\begin{aligned} y_{pred} &= \mathbf{K}^T \alpha \\ &= \begin{bmatrix} 9 & 16 & 1 & 16 \\ 16 & 36 & 0 & 36 \\ 1 & 0 & 9 & 0 \\ 16 & 36 & 0 & 36 \end{bmatrix} \begin{bmatrix} -0.1813 \\ 0.1707 \\ -0.1798 \\ 0.1737 \end{bmatrix} \\ y_{pred} &= \begin{bmatrix} 0.92 \\ 3.94 \\ 0.98 \\ 4.06 \end{bmatrix} \end{aligned}$$

Regression Line using Kernel Regression



Question 1:

Q. Gaussian kernel regression with parameter $\sigma^2 = 1/2$ was applied to the following dataset with two features:

$$X = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \quad y = [2.1 \ 1 \ 2 \ 1.2]^T$$

The weight vector can be written as $w = \phi(X) \alpha$ where $\phi(X)$ is the transformation mapping corresponding to the kernel. The vector α is given by $[2.1 \ -2.1 \ 3 \ 0]^T$ which is obtained as $K^{-1}y$ where K is the kernel matrix. What will be the prediction for point $[1 \ 1]^T$?

A. The kernel function is given by,

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right) = \exp\left(\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\left(\frac{1}{2}\right)}\right) = \exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2)$$

The prediction is given by,

$$\begin{aligned} \mathbf{K}^T \alpha &= \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_{test}) & k(\mathbf{x}_2, \mathbf{x}_{test}) & k(\mathbf{x}_3, \mathbf{x}_{test}) & k(\mathbf{x}_4, \mathbf{x}_{test}) \end{bmatrix} \begin{bmatrix} 2.1 \\ -2.1 \\ 3 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} k\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) & k\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) & k\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) & k\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix}\right) \end{bmatrix} \begin{bmatrix} 2.1 \\ -2.1 \\ 3 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} e^{-1} & e^{-1} & e^0 & e^{-2} \end{bmatrix} \begin{bmatrix} 2.1 \\ -2.1 \\ 3 \\ 0 \end{bmatrix} \\ &= 2.1e^{-1} - 2.1e^{-1} + 3e^0 + 0e^{-2} \\ \mathbf{K}^T \alpha &= 3 \end{aligned}$$