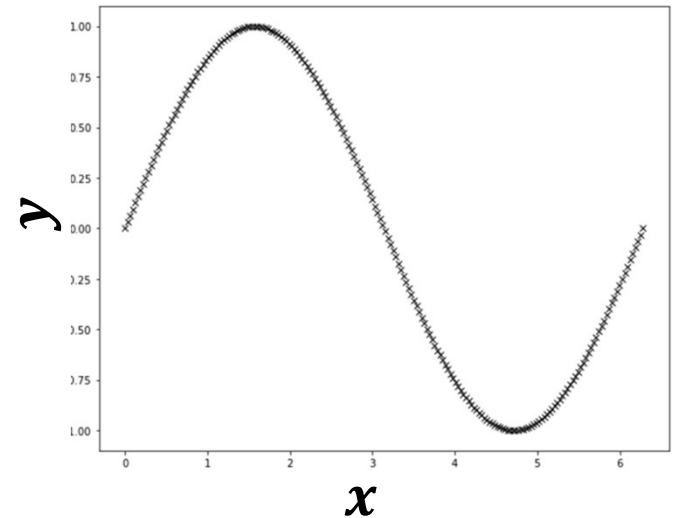


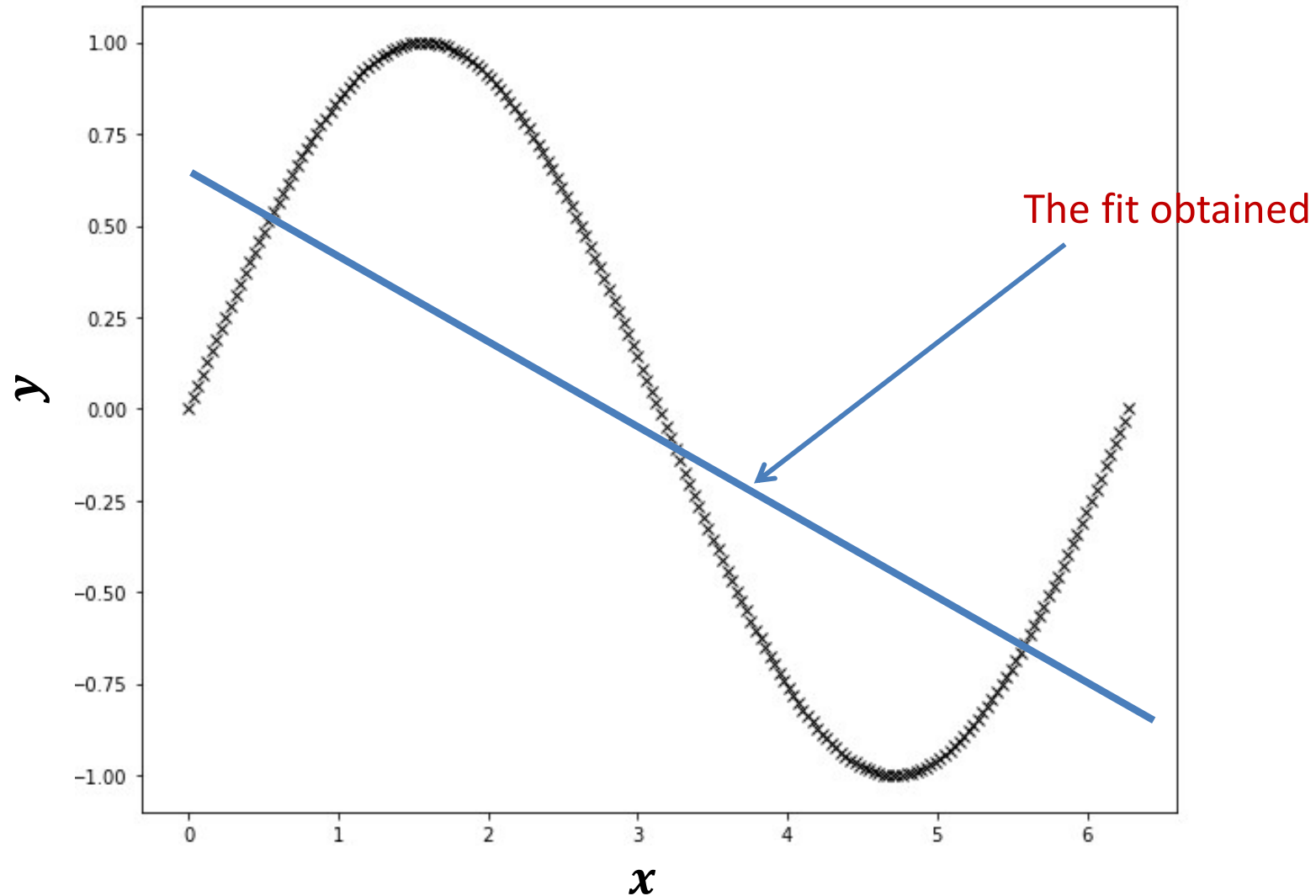
Locally Weighted Linear Regression

- Consider a use case where the relation between the feature vector (independent variables) and the target value (dependent variable) is not linear.
- $y = \sin(x)$
- Suppose we perform linear regression in such case.

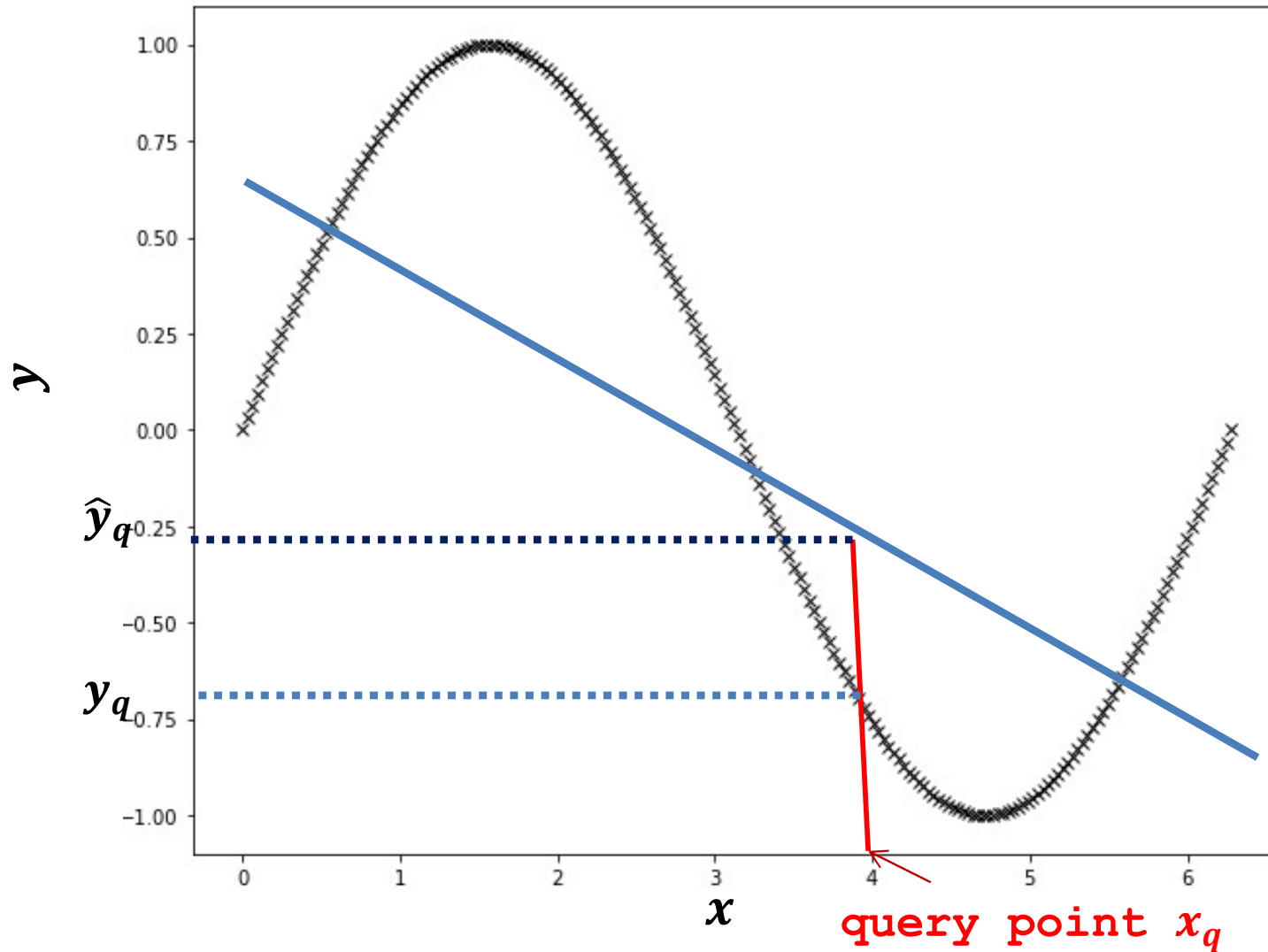


x	y
0	0
$\pi/12$	0.25882
$\pi/6$	0.5
$\pi/4$	0.70711
$5\pi/12$	0.86603
$\pi/3$	0.86603
$7\pi/12$	0.70711
$\pi/2$	0.5
$2\pi/3$	0.25882
$3\pi/4$	0
$5\pi/6$	-0.5
$11\pi/12$	-0.70711
π	-0.86603
$13\pi/12$	-0.86603
$7\pi/6$	-0.70711
$5\pi/4$	-0.5
$17\pi/12$	-0.25882
$3\pi/2$	0
$19\pi/12$	0.25882

Locally Weighted Linear Regression



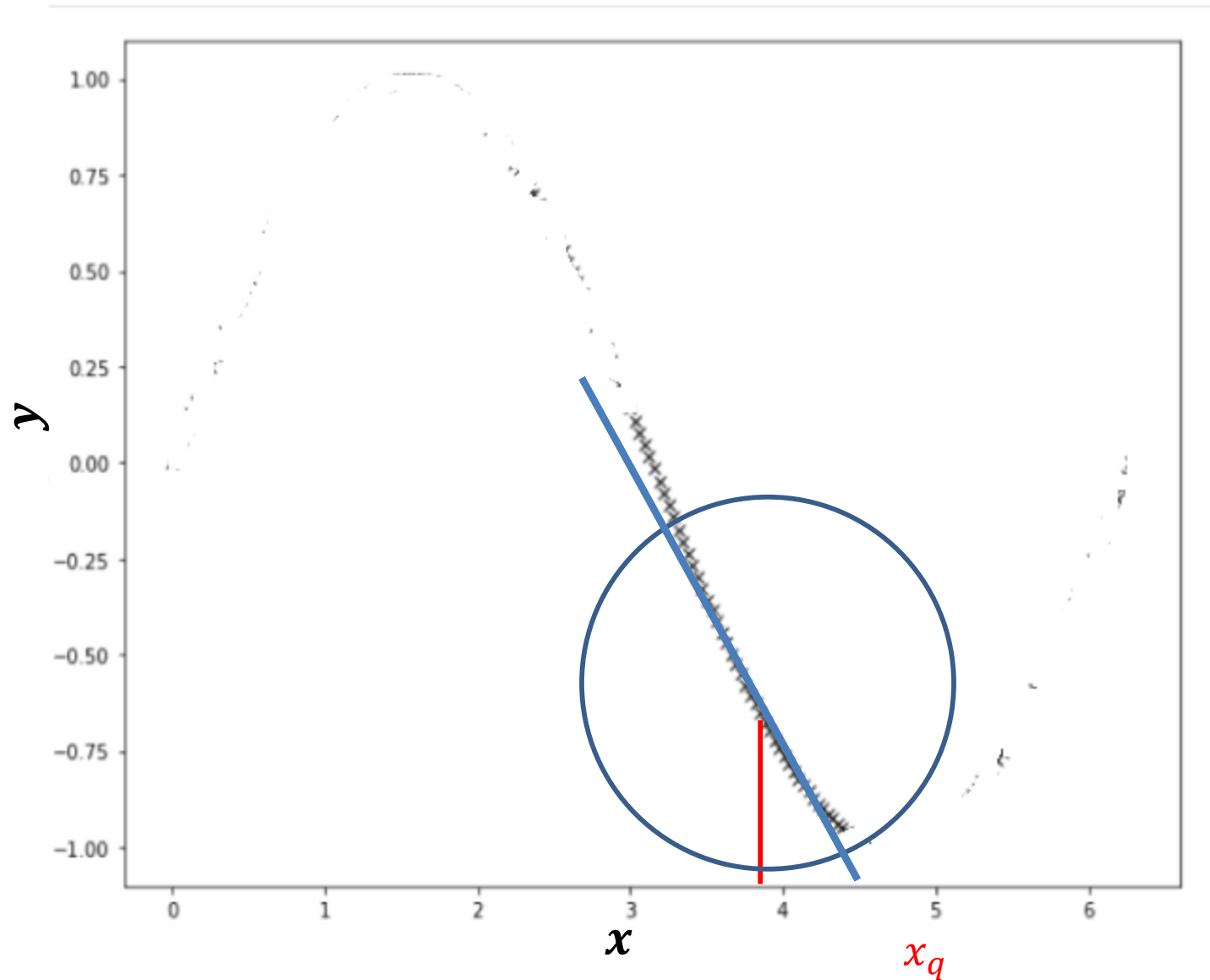
Locally Weighted Linear Regression



Locally Weighted Linear Regression

- Obviously, linear regression does not work.
- The principle of locally weighted regression is as follows.
- Suppose, we know the **query point x_q** where the value of y_q is to be calculated.
- The **samples in the vicinity of the query point** contribute more to linear regression, than the **samples farther from the query point**.

Locally Weighted Linear Regression



Locally Weighted Linear Regression

$$\textit{Find } \theta \textit{ to minimize } \frac{1}{2n} \sum_{i=1}^n w^i (\hat{y}^i - y^i)^2$$

$$\textit{where } w^i = \exp\left(-\frac{(x^i - x_q)^2}{2\tau^2}\right)$$

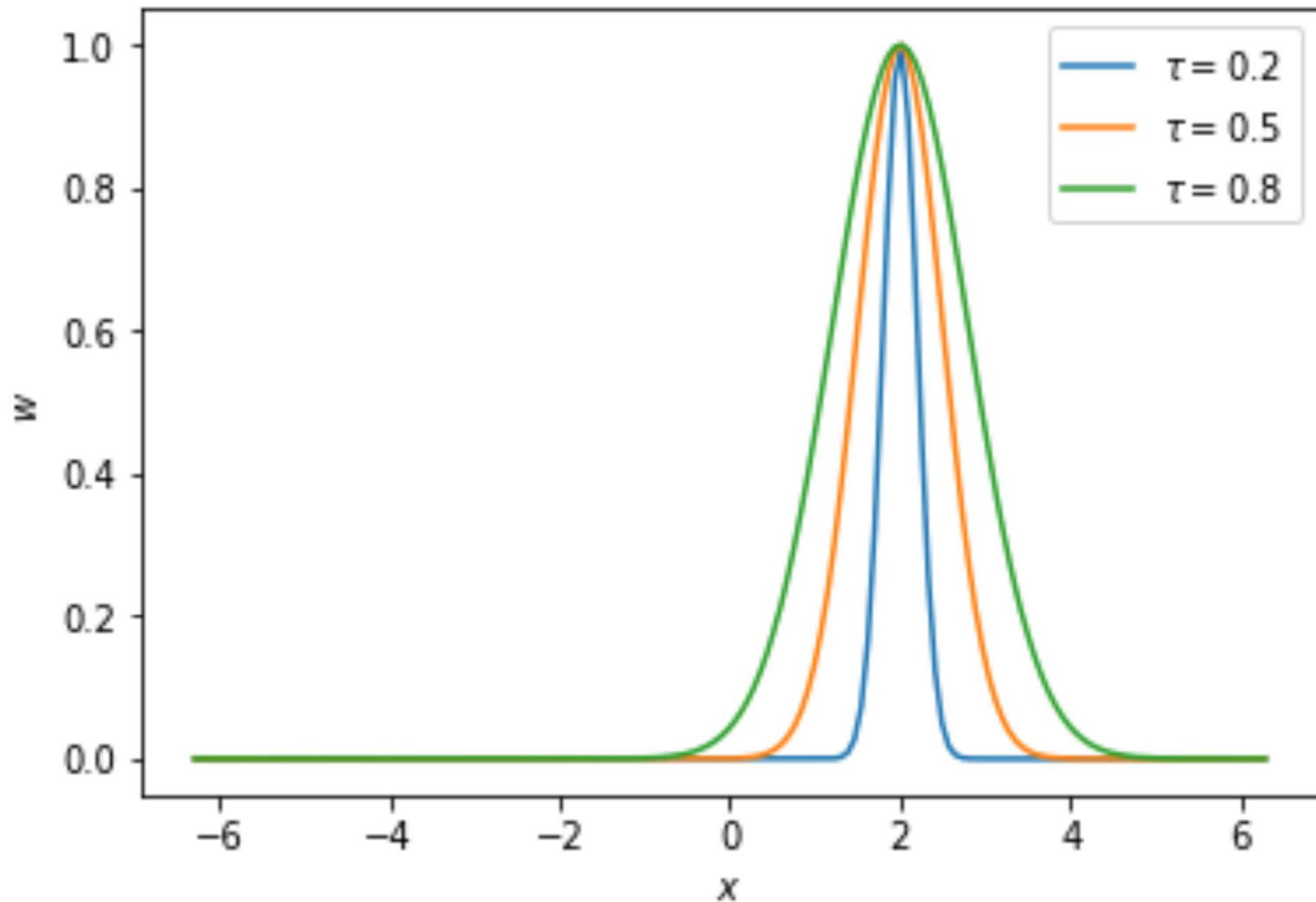
w^i is the weight given to the error of regression contributed by sample i for updating the model parameters.

The farther is x^i from the query point x_q , smaller is the weight.

Locally Weighted Linear Regression

- τ controls how quickly the weight of a training example falls off with distance of $x(i)$ from the query point x_q ;
- τ is called the bandwidth parameter,

Locally Weighted Linear Regression



Locally Weighted Linear Regression

$$\text{Find } \theta \text{ to minimize } J(\theta) = \frac{1}{2n} \sum_{i=1}^n w^i (h_{\theta}(x^i) - y^i)^2$$

$$\text{where } w^i = \exp\left(-\frac{(x^i - x_q)^2}{2\tau^2}\right)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{n} \sum_{i=1}^n w^i (h_{\theta}(x^i) - y^i) x_j^i = \frac{1}{n} \sum_{i=1}^n w^i E^i x_j^i$$

LWR-Batch Gradient Descent

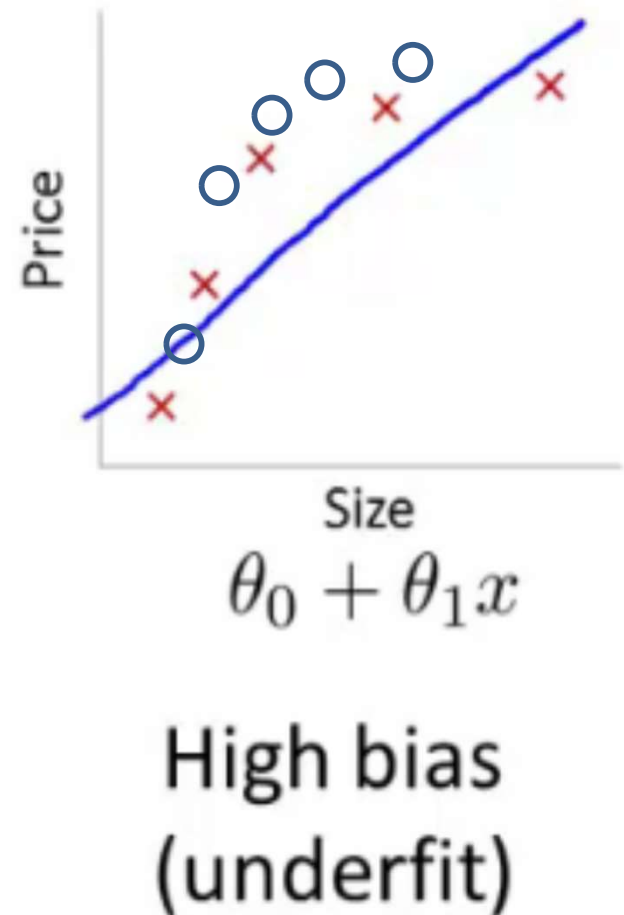
- $g = W * E * X = \begin{bmatrix} W^1 \\ W^2 \\ \vdots \\ W^n \end{bmatrix} \begin{bmatrix} E^1 \\ E^2 \\ \vdots \\ E^n \end{bmatrix} * \begin{bmatrix} x_0^1 & x_1^1 & x_2^1 \\ x_0^2 & x_1^2 & x_2^2 \\ \vdots & \vdots & \vdots \\ x_0^n & x_1^n & x_2^n \end{bmatrix}$
- $grad = \frac{1}{n} np.sum(g, axis = 0)$
- $[\theta_0, \theta_1, \theta_2] = [\theta_0, \theta_1, \theta_2] - \alpha * gra$

Bias and Variance of a Model

- Consider the problem of predicting y from $x \in R$.
- In regression we have several options (linear and non-linear) of fitting $y = h_{\theta}(x)$ to a dataset.
- The following questions are important to judge the choice.
- How well does the **fit/model** perform?
- How is its training error often referred to as **bias**?
- How is its testing error over different datasets often referred to as **variance**?

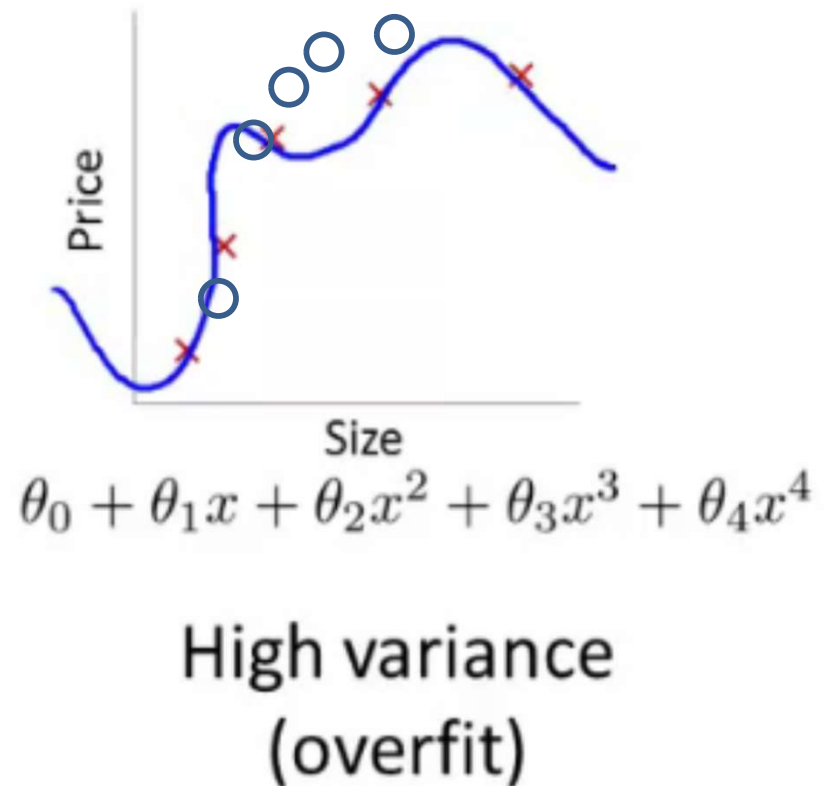
Bias and Variance of a Model

- Consider the problem of predicting y from $x \in R$.
- The figure shows the result of fitting $y = h_{\theta}(x) = \theta_0 + \theta_1 x_1$ to a dataset.
- The fit is not very good.
- The model has high bias and suffers from underfit.



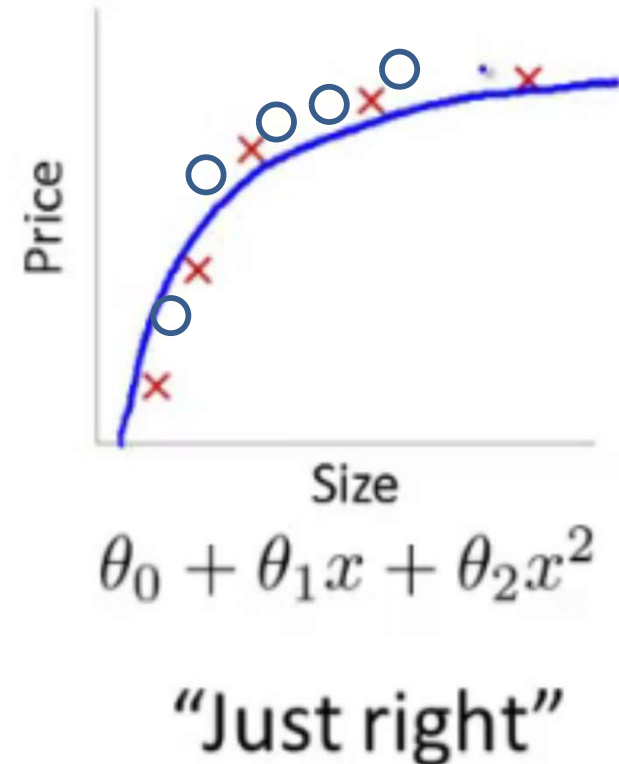
Bias and Variance of a Model

- The figure is the result of fitting a 4-th order polynomial.
- The Fitted curve passes through the data perfectly.
- The model has **less bias**.
- But, the variance i.e the test error may be high.
- This is often referred to as overfitting.



Bias and Variance of a Model

- Instead, if we perform non-linear regression, and fit $y = h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2$, then we obtain a slightly better fit to the data.
- This model has low bias and low variance and proves to be a right fit.



Over fitting

- Over fitting refers to a model that fits the training data too well.
- Let $error_t(h)$ and $error_D(h)$ be the errors of hypothesis h over training samples and the entire distribution of samples respectively.
- Hypothesis $h \in H$ is said to **overfit** training data *iff* there is an alternate hypothesis $h' \in H$ such that
 - $error_t(h) < error_t(h')$ (h performs well on training data but h' doesn't)
 - $error_D(h) > error_D(h')$ (h does not perform well on test data but h' does)