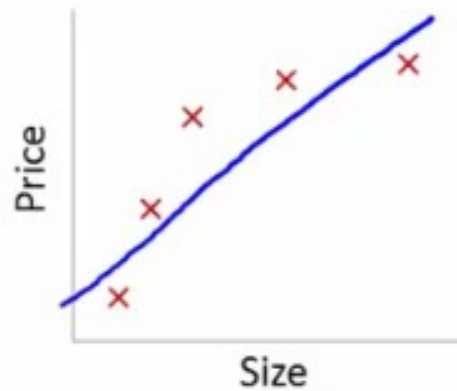


Fitting

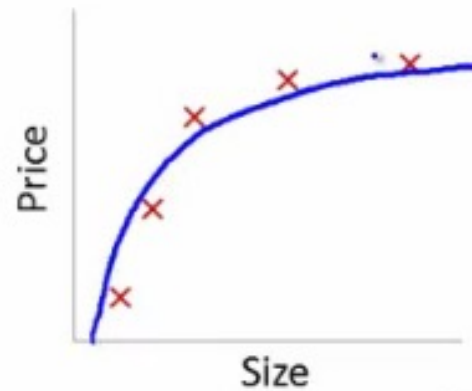
Just Fit

Fitting - regression



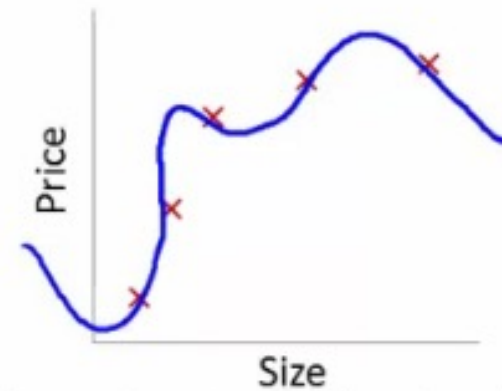
$$\theta_0 + \theta_1 x$$

High bias
(underfit)



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

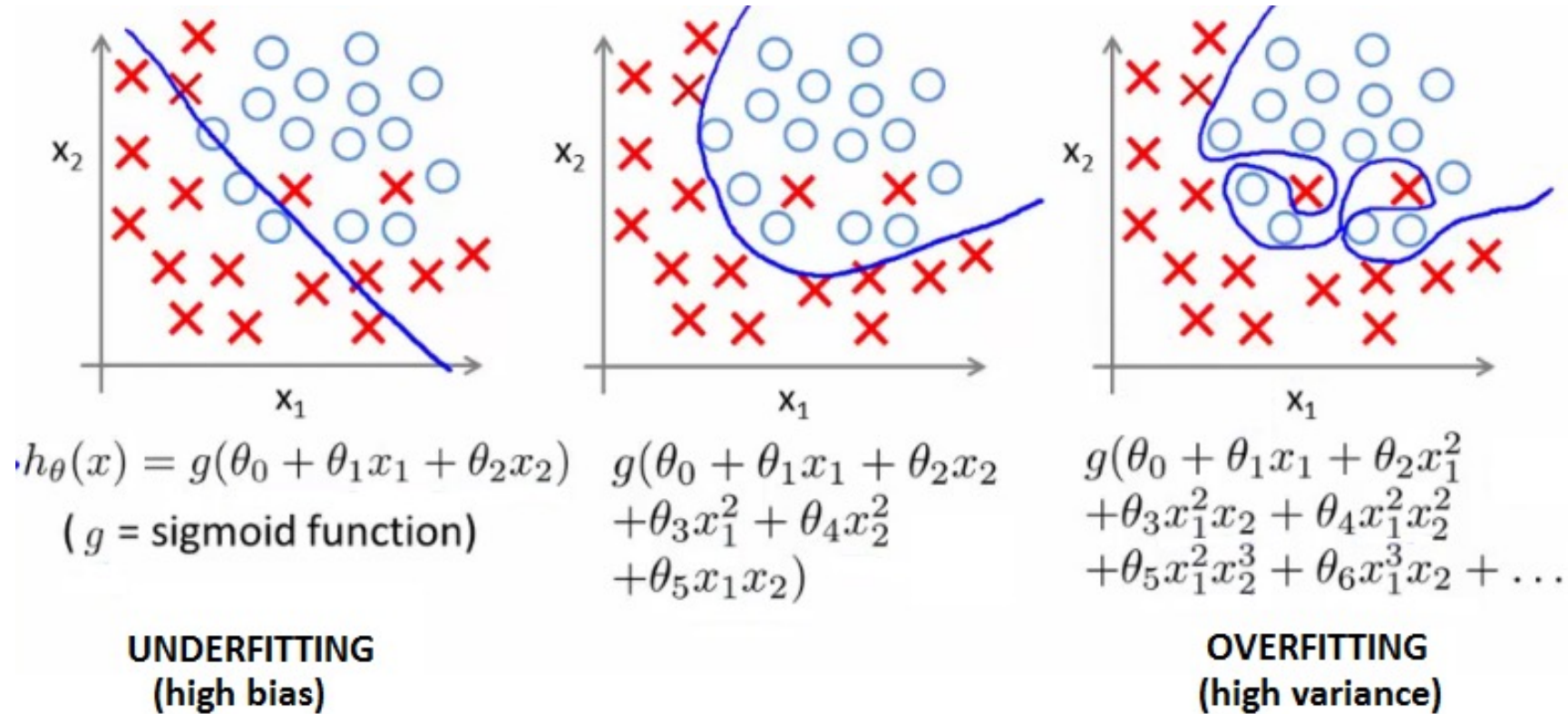
"Just right"



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

High variance
(overfit)

Fitting - classification



Overfitting and Underfitting

We want a model able to generalize from its experience (training example). Machine learning is concerned with perform accurately on new, unseen examples/tasks after having experienced a learning data set

Underfitting: the model does not fit the training data, because it is too simple and not able to learn. *It performs bad both on training and test data.*

Overfitting: the *model performs well on the training data but not on unseen data*, because it is too complex and not able to generalize. The model begins to "memorize" training data rather than "learning" to generalize from trend

Bias and Variance

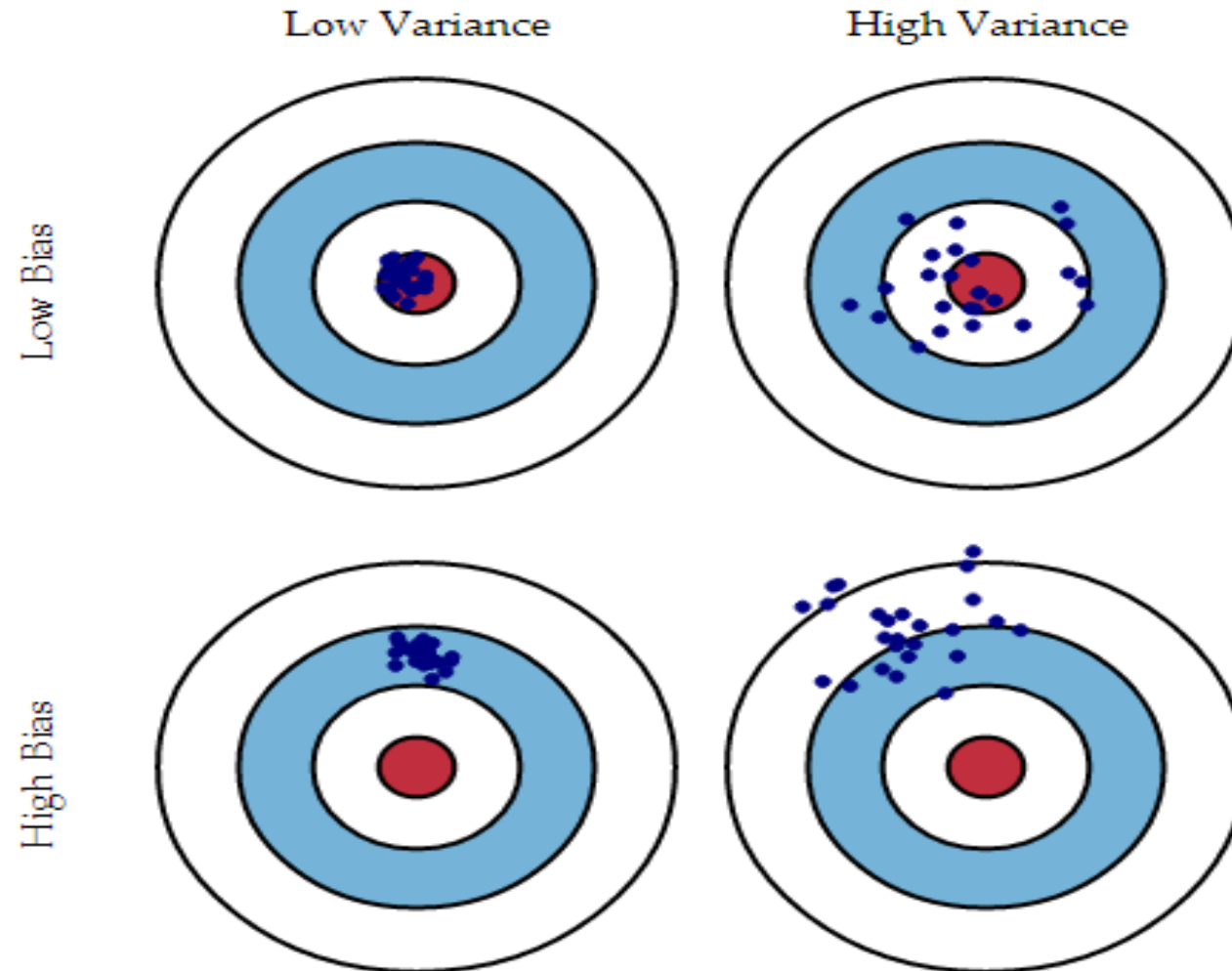
Bias is the systematic deviation of the estimator. It represents the mean of the error of the predicted values:

$$\text{Bias}(h(X), f(X)) = (E[h(X)] - f(X))^2$$

Variance represents the mean of the variations of the predicted values with respect to the mean of the predicted values. It gives us an idea of the mean of the variations of the errors of the predicted values w.r.t. Y :

$$\text{Var}(h(X)) = E(h(X) - E[h(X)])^2$$

Bias and Variance - Intuition



Expected squared error and its relationship with expected value

- Expected value of a random variable X

$$E[X] = \sum_{i=1}^m p_i \cdot x_i$$

- Let X be the input and Y be the output in a problem we want to model

$$Y = f(X) + e$$

- We assume our hypothesis $h(X)$ to approximate $f(X)$. The expected squared error (a.k.a. **Mean Squared Error**) is then

$$Err(f(X)) = E[(Y - h(X))^2]$$

Bias and Variance trade-off

We assume to have an infinite number of *Dataset* D_i , that represents all the possible combinations of the training samples that could feed our model. They are all possible subsets of the corresponding random variable

We fix the model and we train it each time with a different D_i . We get an infinite number of hypothesis $h^{(D_i)}(x)$.

Each training sample is in the form $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$, where y is the sum of the ground truth and a gaussian error with zero mean and variance σ_e :

$$y^{(i)} = f(\mathbf{x}^{(i)}) + e^{(i)}$$

Each hypothesis $h^{(D_i)}(x)$ will be affected by an error in predicting values w.r.t. the ground truth.

We can model this error in the form of a Mean Squared Error (MSE):

$$MSE \left(h^{(D_i)}(x) \right) = E_x \left[(h^{(D_i)}(x) - f(x))^2 \right]$$

Bias and Variance trade-off (cont.d)

We want to estimate the expectation of the error w.r.t. all the possible D_i .

In details we want to compute the MSE of each hypothesis and then we compute the overall mean.

Performing this operation over all the infinite datasets is not feasible, but we can compute the expected value of the mean as the expectation of the MSE over all the datasets.

This expectation is named Generalization Error (GER):

$$\begin{aligned} GER &= E_D[MSE] \\ &= E_D \left[E_x \left[\left(h^{(D_i)}(x) - f(x) \right)^2 \right] \right] \\ &= E_x \left[E_D \left[\left(h^{(D)}(x^{(i)}) - f(x^{(i)}) \right)^2 \right] \right] \end{aligned}$$

The last step can be performed because of the linearity of expectation operator

Bias and Variance trade-off (cont.d)

Focus on the term $E_D \left[\left(h^{(D)}(x^{(i)}) - f(x^{(i)}) \right)^2 \right]$

We can define a new quantity, the best estimation of $f(x^{(i)})$, by feeding each hypothesis with the same training sample $x^{(i)}$, and then we compute the mean of all the values:

$$\bar{h}(x^{(i)}) = E_D \left[h^{(D)}(x^{(i)}) \right]$$

We can now perform a sum zero operation by adding and subtracting the latter from the initial term:

$$E_D \left[\left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) + \bar{h}(x^{(i)}) - f(x^{(i)}) \right)^2 \right]$$

We can solve the power and exploit the linearity of E:

$$E_D \left[\left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right)^2 \right] + E_D \left[\left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right)^2 \right] + E_D \left[2 \left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right) \left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right) \right]$$

Bias and Variance trade-off (cont.d)

$$E_D \left[\left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right)^2 \right] + E_D \left[\left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right)^2 \right] + E_D \left[2 \left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right) \left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right) \right]$$

Let us focus on the first term:

$$E_D \left[\left(h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right)^2 \right] = Var \left(h^{(D)}(x^{(i)}) \right)$$

It matches the variance definition.

The second term:

$$E_D \left[\left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right)^2 \right] = (\bar{h}(x^{(i)}) - f(x^{(i)}))^2$$

Matches the squared bias

Bias and Variance trade-off (cont.d)

The third term:

$$\begin{aligned} & 2 \left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right) E_D \left[h^{(D)}(x^{(i)}) - \bar{h}(x^{(i)}) \right] \\ &= 2 \left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right) \left(E_D \left[h^{(D)}(x^{(i)}) \right] - E_D \left[\bar{h}(x^{(i)}) \right] \right) \\ &= 2 \left(\bar{h}(x^{(i)}) - f(x^{(i)}) \right) \left(\bar{h}(x^{(i)}) - \bar{h}(x^{(i)}) \right) \\ &= 0 \end{aligned}$$

We can move back to the MSE:

$$MSE(h^{(D)}(x^{(i)})) = Bias^2 \left(h^{(D)}(x^{(i)}) \right) + Var \left(h^{(D)}(x^{(i)}) \right)$$

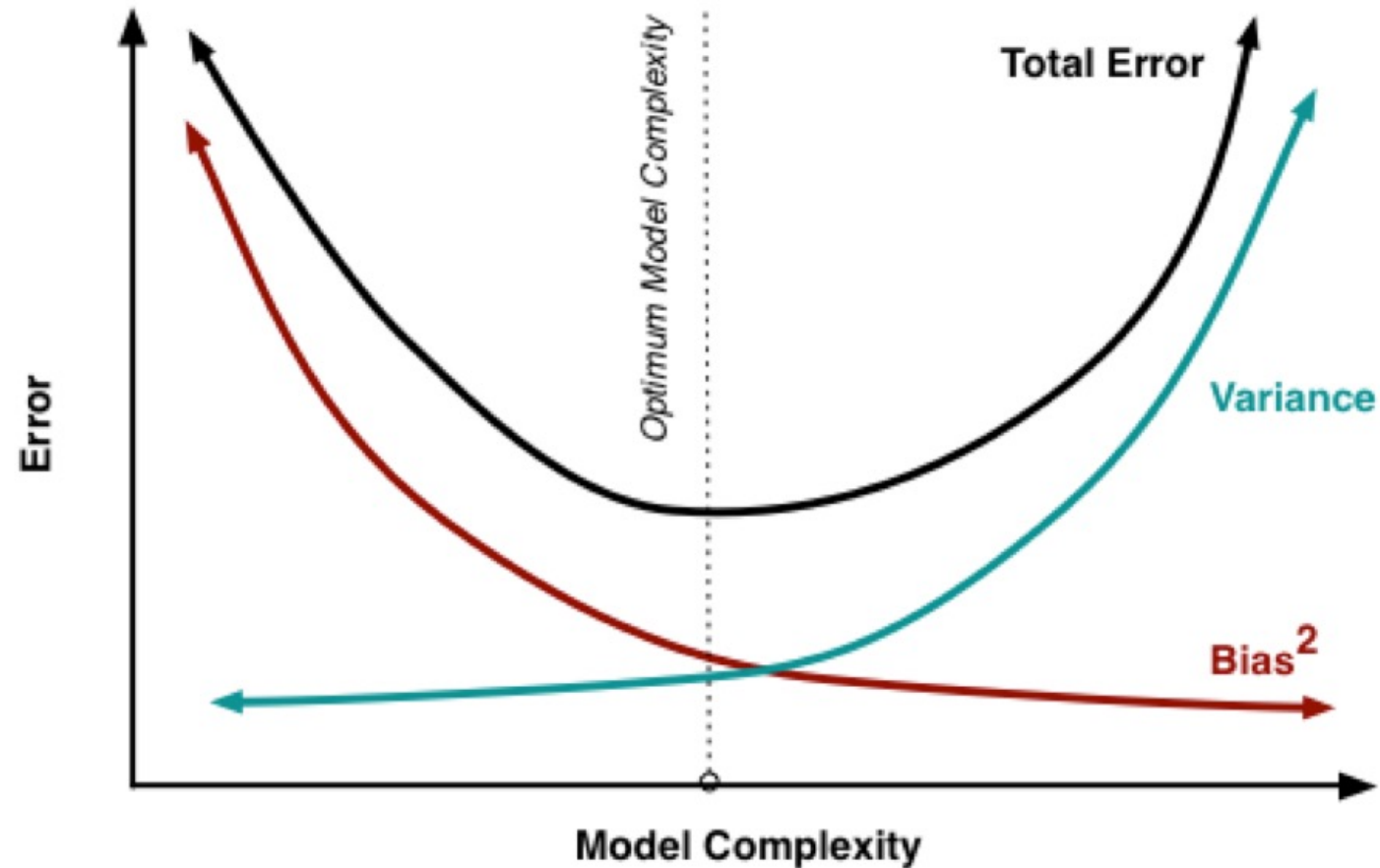
Bias and Variance trade-off (cont.d)

In order to reduce the Generalization Error, we should reduce the bias or the variance:

- Bias represents how much the best estimation is able to get close to the ground truth. The high bias for a specific model denotes a too simple model that is not able to predict, whatever Dataset you are training on.
- Variance represents how much a single hypothesis can be different from the best estimation. The high variance means that the same model, trained with different datasets, generates very different hypothesis. The hypothesis are very complex, they are prone to overfit and they are not able to generalize.

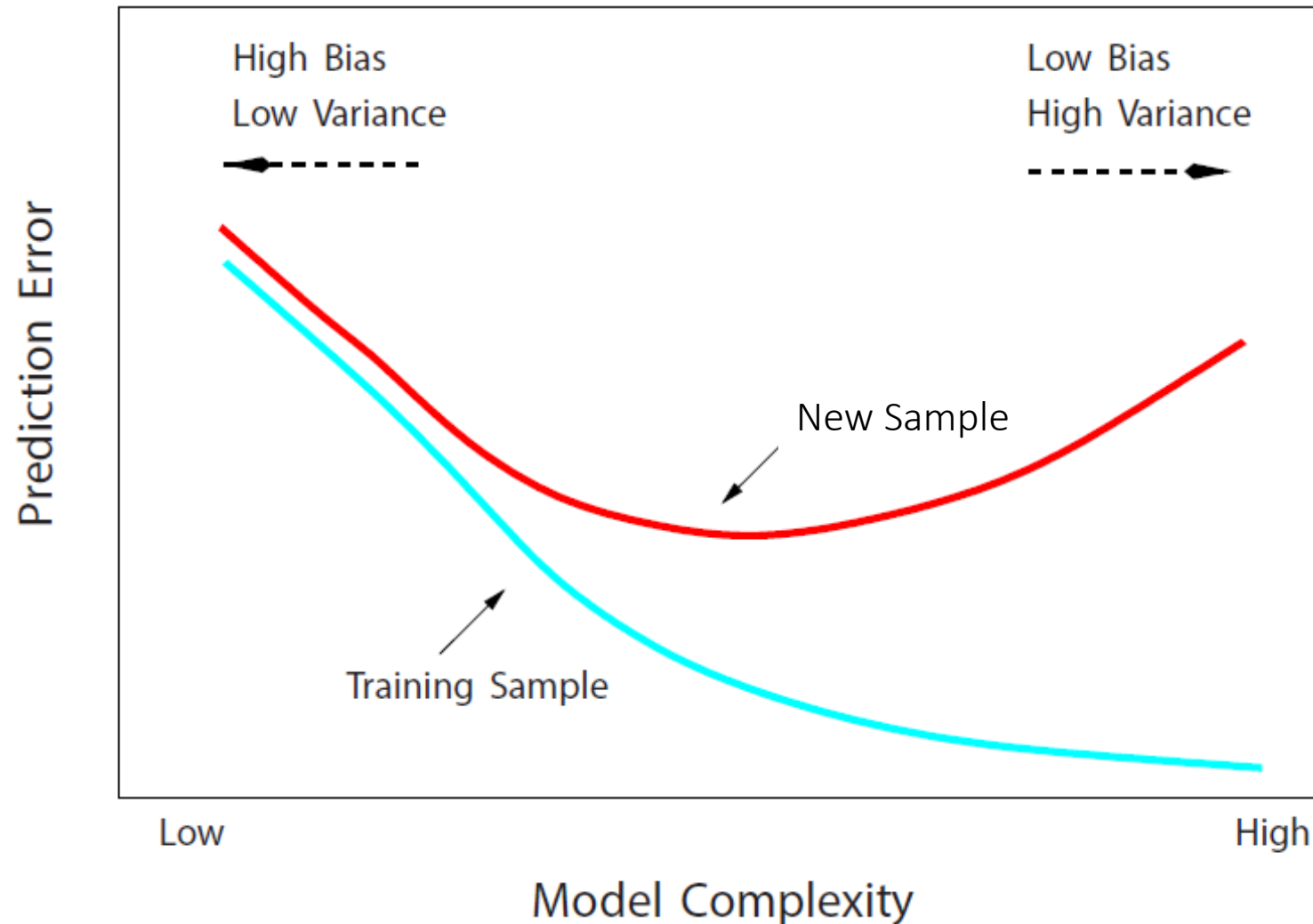
The optimum would be to maintain low bias and variance, but, a model given, the bias and the variance behave at the opposite w.r.t. the complexity of the model.

Bias and Variance w.r.t. complexity



The best trade-off is choosing the model complexity that show the minimum sum between bias and variance

Training set and test set w.r.t. complexity



Fitting

Regularization

Addressing overfitting:

Options:

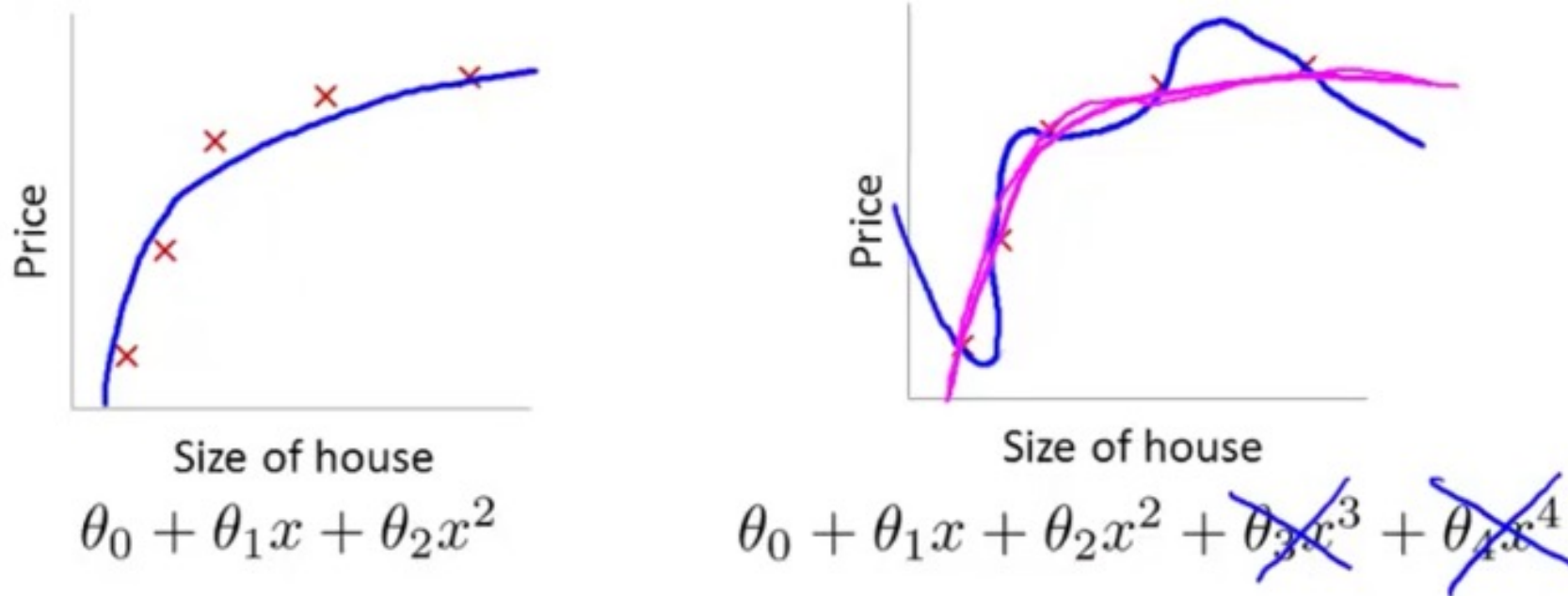
1. Reduce number of features

- Manually select which features to keep
- Model selection algorithm

2. Regularization

- Keep all the features, but reduce magnitude/values of the parameters
- Works well when we have a lot of features, each of which contributes a bit to predicting y

Regularization Intuition



We want to reduce the contribution of the high degree terms of the polynomial to make the curve smoother. If we added an heavy term to the high degree parameters we are imposing the minimization function to make them really small:

$$\min_{\theta} J(\theta) = \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + 10000\theta_3 + 10000\theta_4$$

$$\theta_3 \approx 0 \quad \theta_4 \approx 0$$

Regularization

If we have small values for the parameters leads to have Simple hypothesis that is less prone to overfitting

$$\begin{aligned}
 J(\theta) &= \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 \right) \\
 &\downarrow \\
 J(\theta) &= \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \|\theta\|_2^2 \right) \quad \text{Squared L2-Norm} \\
 &\downarrow \\
 J(\theta) &= \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right) \quad \theta_0 \text{ not included}
 \end{aligned}$$

Lambda is the regularization parameter that controls the influence of regularization w.r.t. the hypothesis:

The bigger is, the smaller will be the theta and smoother will be the curve

Regularization for regression

In regularized linear regression, we choose to minimize

$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right)$$

What if lambda is set to an extremely large value?

- Algorithm works fine; setting to be very large can't hurt it.
- Algorithm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit well even training data).
- Gradient descent will fail to converge.

Regularization for regression

In regularized linear regression, we choose to minimize

$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right)$$

What if lambda is set to an extremely large value?

- Algorithm works fine; setting to be very large can't hurt it.
- Algorithm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit even training data well).
- Gradient descent will fail to converge.

Fitting

Regularized linear regression

Regularization for linear regression

$$J(\theta) = \frac{1}{2m} \left(\sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right)$$

We want to derive the cost function to modify the update rule in the gradient descent:

Repeat until convergence {

$$\begin{cases} \theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_0^{(i)} \\ \theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j = 1, 2, \dots, n \end{cases}$$

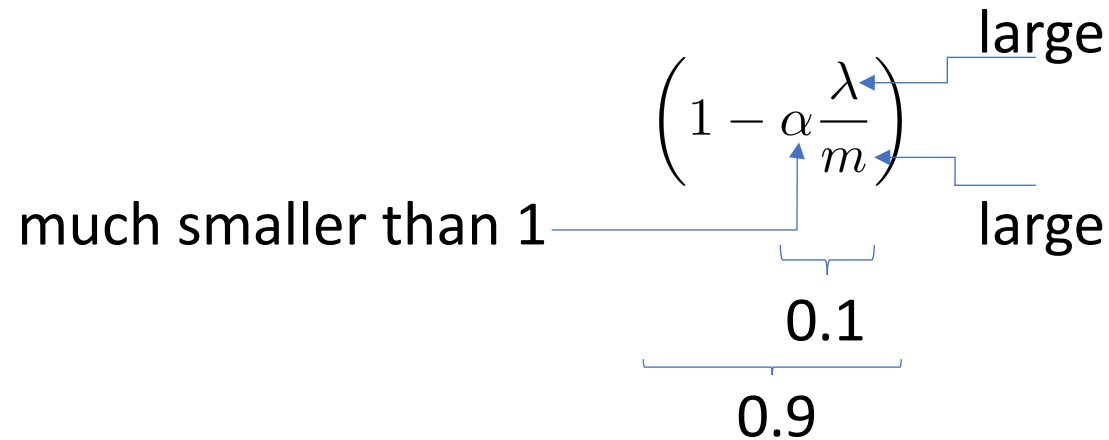
}

Regularization for linear regression

The update function can be re-written as

$$\theta_j = \theta_j \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} \quad j = 1, 2, \dots, n$$

It is worth to note that the second part is same as before the regularization but we update the theta moving from a ratio of the old theta smaller than 1:



Fitting

Regularized logistic regression

Regularization for logistic regression

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(\mathbf{x}^{(i)})) \right) + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

We want to derive the cost function to modify the update rule in the gradient descent:

Repeat until convergence {

$$\begin{cases} \theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_0^{(i)} \\ \theta_j = \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad j = 1, 2, \dots, n \end{cases}$$

}

Fitting

Regularization with L1 norm

Regularization with L1-norm

It is possible to regularize the cost function using L1 norm:

$$\ell^1 = \sum_{j=1}^n |\theta_j|$$

But L1-norm is no more derivable and standard gradient descent cannot be used.

Why should we consider (in some very specific cases) to use it?

Let us focus on the minimization of the norms:

$$\begin{array}{ll} \text{Norm } \ell^2 : & \min \sum_{i=1}^m (\dots)^2 \\ \text{s.t.} & \|\theta\|_2^2 < t \end{array}$$

$$\begin{array}{ll} \text{Norm } \ell^1 : & \min \sum_{i=1}^m (\dots) \\ \text{s.t.} & \|\theta\|_1 < t \end{array}$$

Limiting our analysis to two parameters

$$\begin{array}{ll} \text{Norm } \ell^2 : & \min \sum_{i=1}^m (\dots)^2 \\ \text{s.t.} & \theta_1^2 + \theta_2^2 < t \end{array}$$

$$\begin{array}{ll} \text{Norm } \ell^1 : & \min \sum_{i=1}^m (\dots) \\ \text{s.t.} & |\theta_1| + |\theta_2| < t \end{array}$$

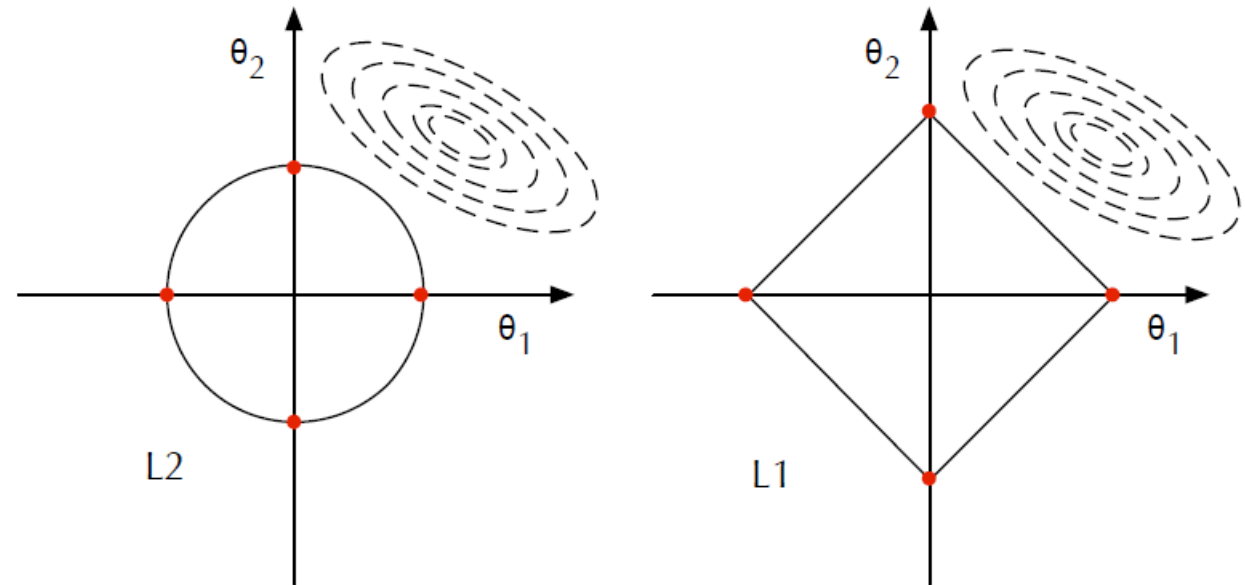
Regularization with L1-norm

$\sum (\dots)^2$ denotes a parabolic curve equation;
 $\theta_1^2 + \theta_2^2 < t$ denotes the internal area of a circle;
 $|\theta_1| + |\theta_2| < t$ denotes the internal area of a rhombus.

During the minimization we want to minimize that parabolic curve and lie inside the circle or the rhombus.

It can be proved that it is easier, in the L1 case, to reach the intercept with the axis and hence to set that parameter to zero.

This lead to a simple but effective feature selection.

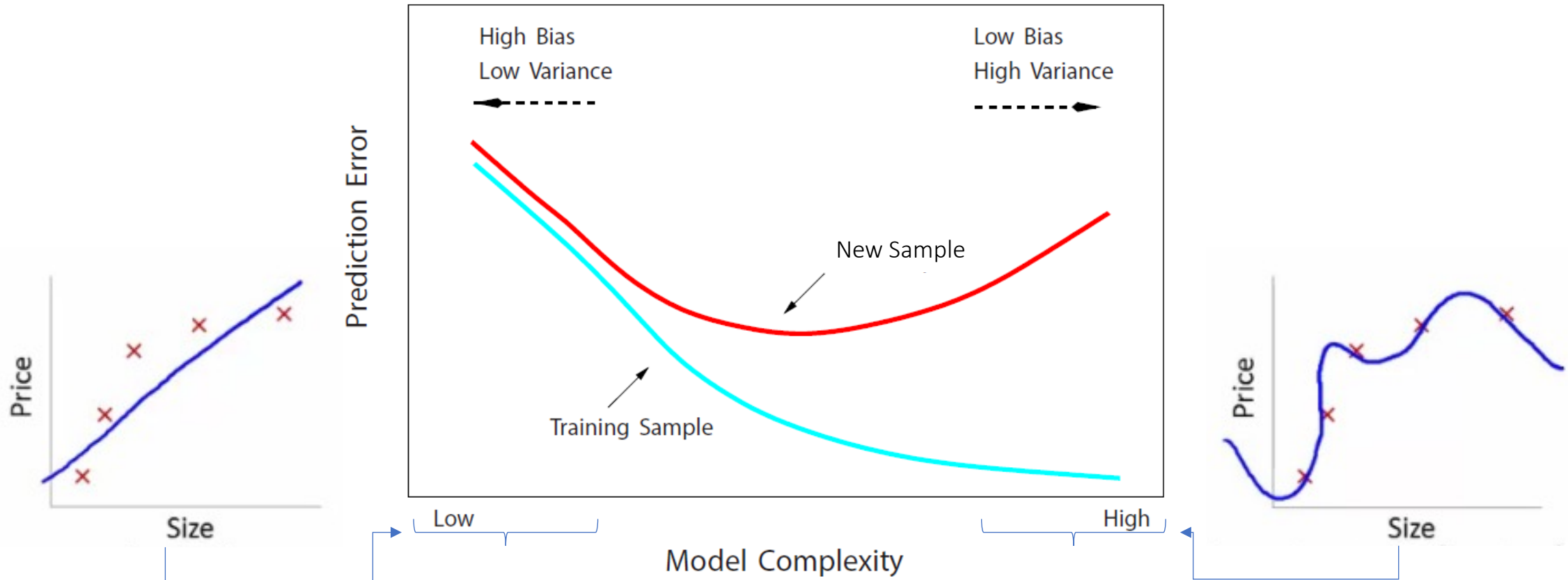


Fitting

Regularization vs Bias/Variance

Training set and test set w.r.t. complexity

Let us recap the Bias variance behavior w.r.t. the complexity of the model



Training set and new samples w.r.t. regularization

