



# Supervised Learning: Regression

# Objective



Objective

Define the set-up of  
Supervised Learning



Objective

Discuss basic  
regression models

# Supervised Learning

| **The set-up:** the given training data consist of  $\langle \text{sample}, \text{label} \rangle$  pairs, or  $(\mathbf{x}, y)$ ; the objective of learning is to figure out a way to predict label  $y$  for any new sample  $\mathbf{x}$ .

| Consider two types of problems:

- **Regression:**  $y$  continuous
- **Classification:**  $y$  is discrete, e.g., class labels.

# The Task of Regression

- | Given: A training set of  $n$  samples  $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$  where  $y^{(i)}$  is a continuous “label” (or target value) for  $\mathbf{x}^{(i)}$
- | To learn a model for predicting  $y$  for any new sample  $\mathbf{x}$ .
- | A simple model is **linear regression**: modeling the relation between  $y$  and  $\mathbf{x}$  via a linear function.

$$y \approx w_0 + w_1x_1 + \dots + w_dx_d = \mathbf{w}^t\mathbf{x}$$

(Note:  $\mathbf{x}$  is *augmented* by adding a dimension of constant 1 to the original sample.)

# Linear Regression

| We can introduce an error term to capture the residual

$$y = \mathbf{w}^t \mathbf{x} + e$$

| Applying this to all  $n$  samples, we have :

$$\mathbf{y} = \mathbf{X} \mathbf{w} + \mathbf{e}$$

$$\begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{pmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \begin{pmatrix} e^{(1)} \\ e^{(2)} \\ \vdots \\ e^{(n)} \end{pmatrix}$$

| *Learning* in this case is to figure out a good  $\mathbf{w}$ .

# Linear Regression (cont'd)

| Find an optimal  $\mathbf{w}$  by minimizing the squared error

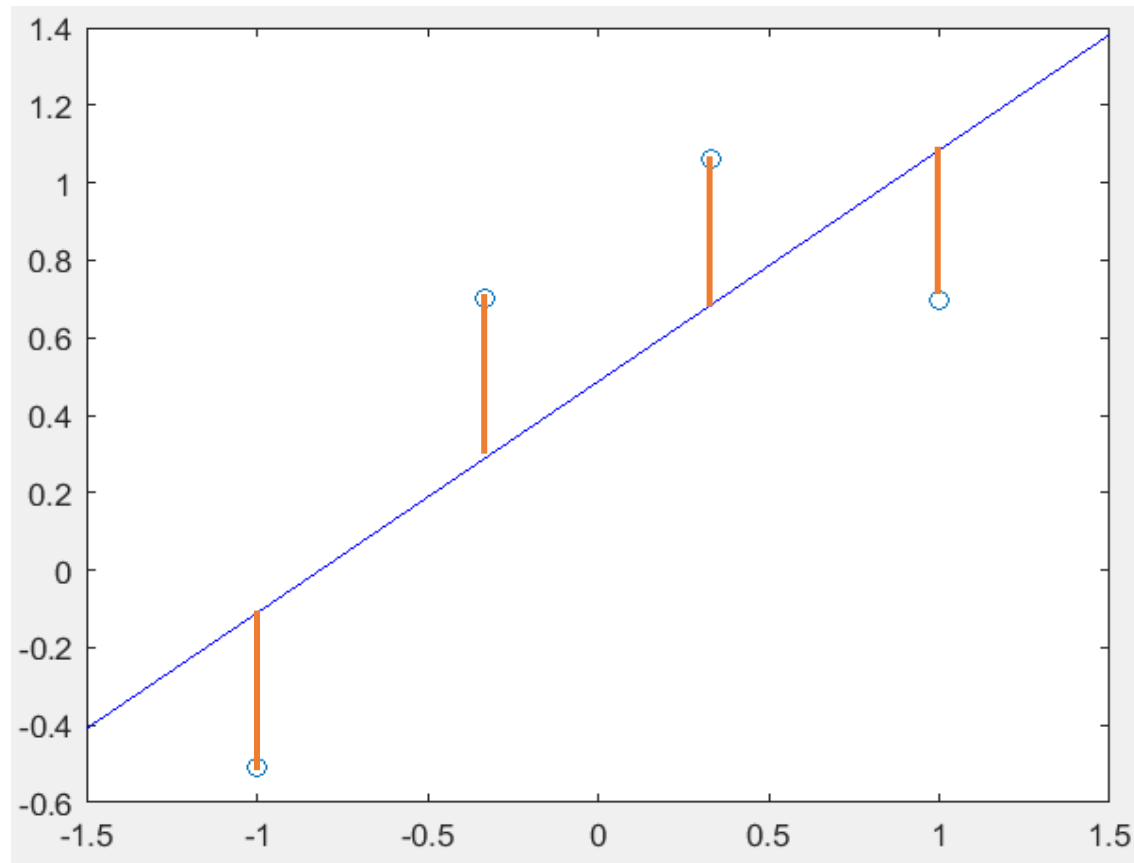
$$\|\mathbf{e}\|^2 = \|\mathbf{y} - \mathbf{X} \mathbf{w}\|^2$$

| The solution can be found to be:

$$\hat{\mathbf{w}} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{y}$$

| In practice, some iterative approaches may be used (e.g., gradient descent search).

# A simple example



# Generalizing Linear Regression

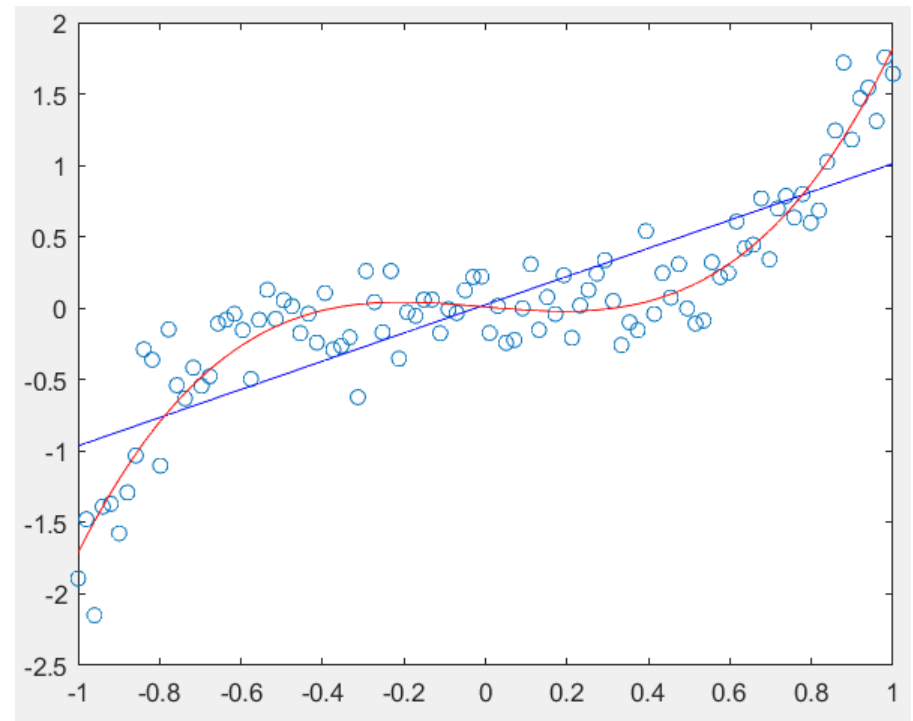
Introducing some **basis functions**  $\phi_j(\mathbf{x})$ :

$$y = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_{M-1}\phi_{M-1}(\mathbf{x})$$

Compare:

➤ Blue: Linear Regression

➤ Red: With  $\phi_j(x) = x^j$





# Regularized Least Squares

| E.g., use a new error function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

- $\lambda$  is the regularization coefficient
- $E_D(\mathbf{w})$  is the data-dependent error
- $E_W(\mathbf{w})$  is the *regularization term*, e.g.,  $E_W(\mathbf{w}) = \|\mathbf{w}\|^q$

| Help to alleviate overfitting.