Consider a linear regression problem, where $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$, with squared loss. The hypothesis set is the set of *constant* functions, that is $\mathcal{H} = \{h_a : a \in \mathbb{R}\}$, where $h_a(\mathbf{x}) = a$. Let $S = ((\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m))$ denote the training set.

- Derive the hypothesis $h \in \mathcal{H}$ that minimizes the training error.
- Use the result above to explain why, for a given hypothesis $\hat{h}$ from the set of all linear models, the coefficient of determination $R^2 = 1 - \frac{\sum_{i=1}^{m}(\hat{h}(\mathbf{x}_i) - y_i)^2}{\sum_{i=1}^{m}(y_i - \bar{y})^2}$ where $\bar{y}$ is the average of the $y_i, i = 1, \ldots, m$ is a measure of how well $\hat{h}$ performs (on the training set).

$\mathcal{H} : \quad h_a \in \mathcal{H} \quad , \quad a \in \mathbb{R}$

$h_a(\vec{x}) = a \quad \forall \, \vec{x} \in \mathcal{X}$

## Solution

- Given $h_\partial \in \mathcal{H}$, the training error for such hypothesis:

$$L_S(h_\partial) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\partial(\vec{x}_i) - y_i \right)^2$$

since $h_\partial(\vec{x}) = \partial$
$\forall \vec{x} \in X$

$$= \frac{1}{m} \sum_{i=1}^{m} \left( \partial - y_i \right)^2$$

Now, finding $h_\partial \in \mathcal{H}$ that minimizes the training error corresponds to find $\partial$ that minimizes

$$L_S(h_\partial) = \frac{1}{m} \sum_{i=1}^{m} \left( \partial - y_i \right)^2 = (\quad)\partial^2 + (\quad)\partial + (\quad)$$

As a function of $\partial$



2

$\Rightarrow$ compute $\dfrac{d \, L_s(h_\partial)}{d\partial}$ and derive $\partial$ s.t. $\dfrac{d \, L_s(h_\partial)}{d\partial} = 0$

$$\frac{d \, L_s(h_\partial)}{d\partial} = \frac{d}{d\partial} \left( \frac{1}{m} \sum_{i=1}^{m} (\partial - y_i)^2 \right)$$

$$= \frac{1}{m} \sum_{i=1}^{m} \frac{d\left((\partial - y_i)^2\right)}{d\partial}$$

$$= \frac{1}{m} \sum_{i=1}^{m} 2(\partial - y_i)$$

$$\frac{d\left((\partial - y_i)^2\right)}{d\partial}$$

$$= 1 \cdot \frac{d\left((\partial - y_i)^2\right)}{d(\partial - y_i)}$$

$$= 1 \cdot 2(\partial - y_i)$$

$$\frac{2}{m} \sum_{i=1}^{m} (\partial - y_i) = 0$$

$$\Leftrightarrow \sum_{i=1}^{m} (\partial - y_i) = 0$$

3

$$\Leftrightarrow \left( \sum_{i=1}^{m} a \right) - \left( \sum_{i=1}^{m} y_i \right) = 0$$

$$\Leftrightarrow m \cdot a = \sum_{i=1}^{m} y_i \Big/ m = \bar{y}$$

$$\cdot) \quad R^2 = 1 - \left( \sum_{i=1}^{m} \left( \hat{h}(\vec{x}_i) - y_i \right)^2 \right) \Big/ \left( \sum_{i=1}^{m} \left( y_i - \bar{y} \right)^2 \right)$$

this is the error of $\hat{h}$ (error on the training set) relative to the error of the "best" naive predictor (which predicts a constant, without looking at $\vec{x}$)

# Polynomial models

Regression problem ($\bar{Y} = \mathbb{R}$), $\mathcal{X} = \mathbb{R}$.

How can we use as hypothesis set $\mathcal{H}$ the set of polynomials of degree $r$ with the machinery we have already developed?

polynomial of degree $r$ : $w_0 \cdot 1 + w_1 x + w_2 x^2 + w_3 x^3 + \ldots + w_{r-1} x^{r-1} + w_r x^r$

Given $x \in \mathbb{R}$, obtain vector: (feature expansion)

$$\vec{x}' = \begin{bmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^{r-1} \\ x^r \end{bmatrix} \implies$$

the hypothesis class of linear models for $\vec{x}'$ corresponds to polynomials of degree $r$ for $x$.

Given $\vec{x} \in \mathbb{R}^d$, $\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$, use the following expansion:

$$\vec{x}' = \begin{bmatrix} 1 \\ x_0 \\ x_0^2 \\ \vdots^r \\ x_0^r \\ x_1 \\ x_1^2 \\ x_1^r \\ \vdots \\ x_d \\ x_d^2 \\ \vdots^r \\ x_d^r \end{bmatrix} \implies$$ use linear models for $\vec{x}'$

Different expansion: $r = 2$, $\vec{x} \in \mathbb{R}^3$    $\vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$

$$\Rightarrow \text{obtain}$$

$$\vec{X'} = \begin{pmatrix} 1 \\ x_0 \\ x_1 \\ x_2 \\ x_0^2 \\ x_1^2 \\ x_2^2 \\ x_0 \, x_1 \\ x_1 \, x_2 \\ x_0 \, x_2 \end{pmatrix} \Rightarrow \text{build linear models on } \vec{X'}$$