# Week 9: Linear Regression

## Table of Contents
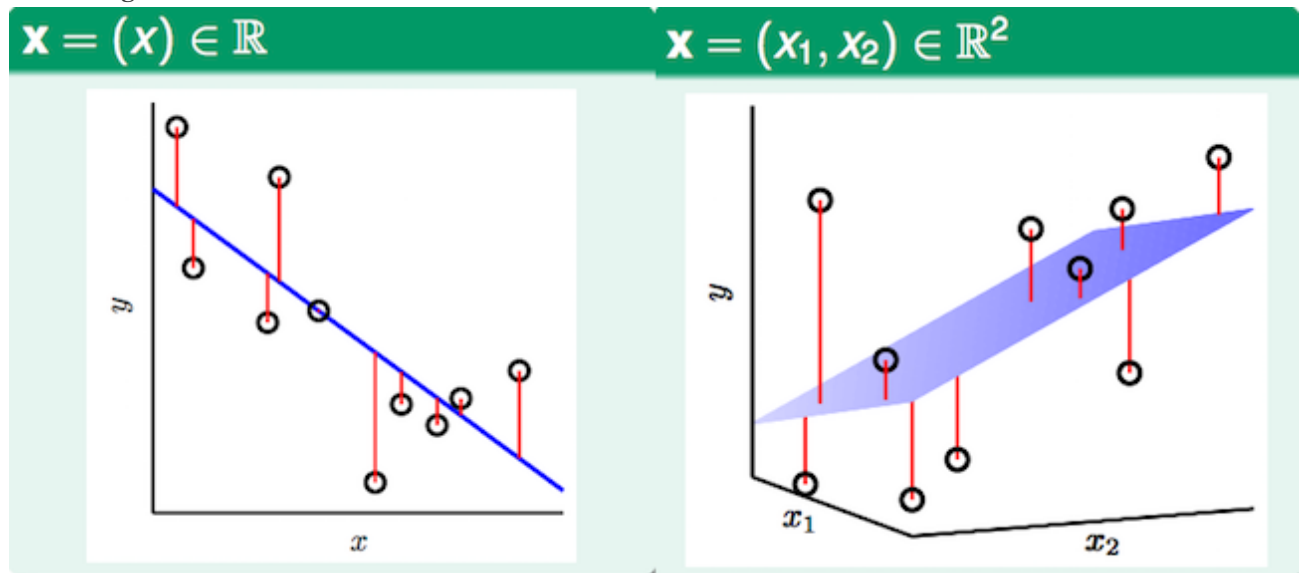
### Linear Regression Problem

1. For features $\mathbf{x} = (x_0, x_1, x_2, \ldots, x_d)$, approximate the target $y$ with a **weighted sum**

$$y \approx \sum_{i=0}^{d} w_i x_i$$

2. Linear regression hypothesis: $h(x) = \mathbf{w}^T \mathbf{x}$. Similar to perceptron, but taking the *real number* value instead of just the sign.

3. Linear regression illustrated



The goal of linear regression is to find the best-fitting *line/hyperplane* with small *residuals*

4. Linear regression error measure
   - Squared error $err(\hat{y} - y)^2$ is often used as the error measure for linear regression

$$E_{\text{in}}(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^{N} (\underbrace{h(\mathbf{x}_n)}_{\mathbf{w}^T \mathbf{x}_n} - y_n)^2 \qquad E_{\text{out}}(\mathbf{w}) = \underset{(\mathbf{x}, y) \sim P}{\mathcal{E}} (\mathbf{w}^T \mathbf{x} - y)^2$$

## Linear Regression Algorithm

1. Cost function of linear regression

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} (\hat{y}_n - y_n) = \frac{1}{N} \sum_{n=1}^{N} (w^T x_n - y_n)^2$$

**Goal**: Find **w** that minimizes cost funciton / in-sample error
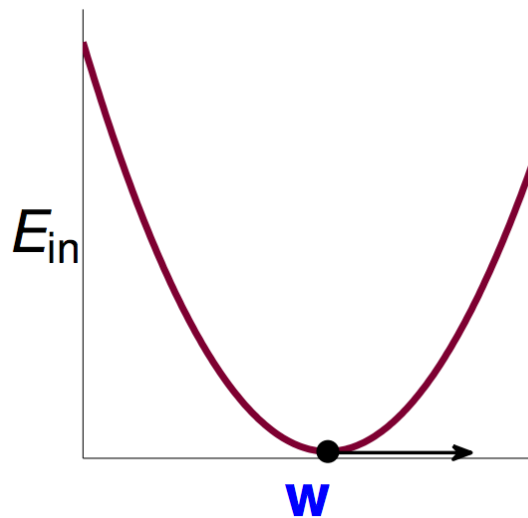
2. Matrix form of linear regression in-sample error

$$E_{in}(w) = \frac{1}{N} \sum_{n=1}^{N} (w^T x_n - y_n)^2 = \frac{1}{N} \sum_{n=1}^{N} (x_n^T w - y_n)^2$$

$$= \frac{1}{N} \left\| \begin{matrix} x_1^T w - y_1 \\ x_2^T w - y_2 \\ \cdots \\ x_N^T w - y_N \end{matrix} \right\|^2$$

$$= \frac{1}{N} \left\| \begin{bmatrix} -- x_1^T -- \\ -- x_2^T -- \\ \cdots \\ -- x_N^T -- \end{bmatrix} w - \begin{bmatrix} y_1 \\ y_2 \\ \cdots \\ y_3 \end{bmatrix} \right\|^2$$

$$= \frac{1}{N} \| \underbrace{X}_{N \times d+1} \underbrace{w}_{d+1 \times 1} - \underbrace{y}_{N \times 1} \|^2$$

## Optimizing In-Sample Error

1. $E_{in}(w)$ is a convex function

$$\underset{w}{min} \, E_{in}(w) = \frac{1}{N} \| Xw - y \|^2$$

- $X$ and $y$ come from the training dample $\mathcal{D}$, therefore $E_{in}$ is only a function of $w$
- $E_{in}(w)$ is **continuous, differentiable, and convex**, which are the necessary conditions for minimizing $E_{in}$ w.r.t. $w$

2. To minimize $E_{in}$, find $w$ that gives gradien of 0.

$$E_{in}(w) \equiv \begin{bmatrix} \dfrac{\partial E_{in}}{\partial w_0}(w) \\ \dfrac{\partial E_{in}}{\partial w_1}(w) \\ \dots \\ \dfrac{\partial E_{in}}{\partial w_d}(w) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

3. The gradient $\nabla E_{in}(w)$

$$E_{in}(w) = \frac{1}{N} \|Xw - y\|^2$$

$$= \frac{1}{N} (\underbrace{w^T X^T X w}_{A} - \underbrace{2w^T X^T y}_{b} + \underbrace{y^T y}_{c}) \quad \text{polynomial expansion}$$

Substitute and take partial derivative:

| Single $w$ | | Vector $w$ |
|---|---|---|
| $a, b, c$ are constants | | $A, \mathbf{b}$ are vectors, $c$ is a constant |
| | Generalizes to → | |
| $E_{in}(w) = \frac{1}{N}(aw^2 - 2bw + c)$ | | $E_{in}(w) = \frac{1}{N}(w^T A w - 2w^T \mathbf{b} + c)$ |
| $\nabla E_{in}(w) = \frac{1}{N}(2aw - 2b)$ | | $\nabla E_{in}(w) = \frac{1}{N}(2Aw - 2\mathbf{b})$ |

Substitute again and simplify. The following applies to both 1D and multi-dimentional cases:

$$\nabla E_{in}(w) = \nabla \frac{1}{N}(w^T X^T X w - 2w^T X^T y + y^T y)$$

$$= \frac{2}{N}(X^T X w - X^T y)$$

4. Optimal linear regression weights

    - Task: Find $w_{LIN}$ such that $\frac{2}{N}(X^TXw - X^Ty) = \nabla E_{in}(w) = 0$
    - When $X^TX$ is **invertible**
        - **Unique** solution
        - $X^TX$ is often invertible
            - $X$ has dimension $N \times (d+1) \to X^TX$ has dimension $(d+1) \times (d+1)$
                - $N$ being number of training samples, $d$ being the number of variables in the model, or degrees of freedom, or VC dimension
            - $N >> d$ most of the time, so there's likely enough 0s in the resulting matrix for $X^TX$ to be invertible

$$w_{LIN} = \underbrace{(X^TX)^{-1}X^T}_{\text{pseudo-inverse } X^\dagger}y$$

    - When $X^TX$ is **singular**
        - **many** optimal solution, one of them being

$$w_{LIN} = X^\dagger y \quad \text{with different definition for } X^\dagger$$

    - Practical suggestion
        - Use **well-implemented, existing** routine to obtain $X^\dagger$ directly, instead of calculating $(X^TX)^-1X^T$ on a case-by-case basis
        - Helps with cases where $X^TX$ is *almost* singular, as such edge cases are already taken care of by built-in routine

5. The linear regression algorithm, in d-dimensions

    - From $\mathcal{D}$, construct input matrix $X$ and output vector $\mathbf{y}$ as:

$$X = \underbrace{\begin{bmatrix} -\,-\,x_1^T\,-\,- \\ -\,-x_2^T\,-\,- \\ \ldots \\ -\,-x_N^T\,-\,- \end{bmatrix}}_{N\times(d+1)} \qquad \mathbf{y} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \ldots \\ y_3 \end{bmatrix}}_{N\times 1}$$

    - Calculate *pseudo-inverse* $\underbrace{X^\dagger}_{(d+1)\times N}$
    - Return $\underbrace{\mathbf{w}_{LIN}}_{(d+1)\times 1} = X^\dagger y$

## Generalization of Lienar Regression

1. Guarantee of linear regression analytic solution: **Average** in-sample error $\overline{E_{in}}$ (across all training samples) is **smaller** than the noise level contained in training data, and decreses as sample size $N$ grows:

$$\overline{E_{in}} = \mathcal{E}_{D\sim P_N}\{E_{in}(w_{LIN} \text{ w.r.t } \mathcal{D})\} = \text{ noise level} \cdot (1 - \frac{d+1}{N})$$
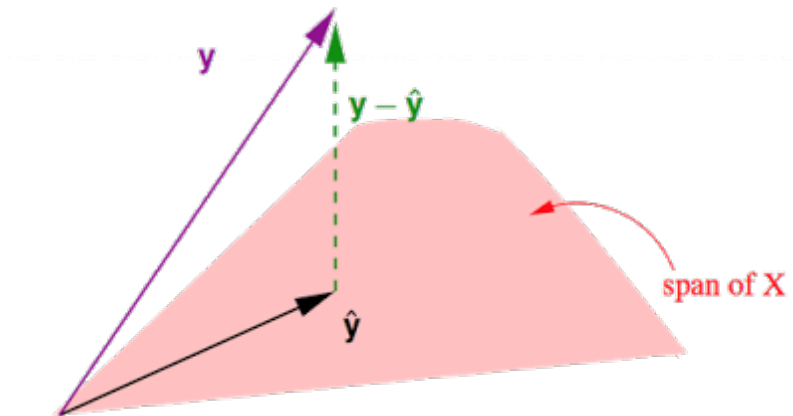
**The Hat Matrix**

1.  Summary: Given optimal $w_{LIN}$ , in-sample error of linear regression can be represented as:

$$E_{in}(w_{LIN}) = \frac{1}{N} \|y - \hat{y}\|^2$$

$$= \frac{1}{N} \|y - \underbrace{XX^\dagger y}_{w_{LIN}}\|^2$$

$$= \frac{1}{N} \|(\underbrace{I}_{identity} - XX^\dagger)y\|^2$$

- $XX^\dagger$ is known as **hat matrix** $H = X(X^TX)^{-1}X^T$

2.  Geometric view of hat matrix



In n-dimensional $\mathbb{R}^N$ :

- $X$ matrix can be viewed as a hyperplane (red area)
- Geometrically, the smallest possible residual $y - \hat{y}$ should be **perpendicular** to the $X$ hyperplane
- $H$ creates $\hat{y}$, the projection of $y$ onto $X$ hyperplane
  - For smallest residual, let

$$Hy = \hat{y}$$
$$y - Hy = y - \hat{y}$$
$$(I - H)y = y - \hat{y}$$

In other words, $I - H$ creates a *perpendicular* projection of $y$ onto $X$ hyperplane

3.  Properties of $H$

- Symmetric

$$H^T = (X(X^TX)^{-1}X^T)^T$$
$$= X((X^TX)^{-1})^TX^T$$
$$= X(X^TX)^{-1}X^T$$
$$= H$$

- Idempotent

$$H^2 = (X(X^T X)^{-1} X^T)(X(X^T X)^{-1} X^T)$$
$$= X \underbrace{(X^T X)^{-1} (X^T X)}_{I} (X^T X)^{-1} X^T$$
$$= X(X^T X)^{-1} X^T$$
$$= H$$

- ○ Positive semi-definite
    - ▪ All eigenvalues are non-negative
    - ▪ $\lambda = eigenvalues, b = eigenvectors$

$$Hb = \lambda b$$
$$H^2 b = \lambda Hb = \lambda(\lambda b)$$
$$\because H = H^2$$
$$H^2 b = Hb = \lambda b$$
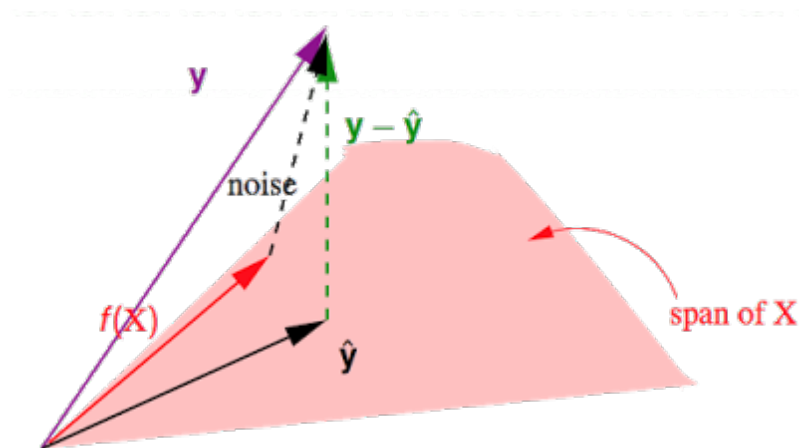$$\therefore \lambda^2 b = \lambda b$$
$$\lambda(\lambda - 1)b = 0$$
$$\lambda = 0 \text{ or } \lambda = 1$$

- ○ Trace of $H$
    - ▪ *Trace* of a matrix is the sum of all its eigenvalues

$$trace(I - H) = N - (d + 1)$$

4. Hat matrix when $y$ contains noise



- ○ Assume training input $y$ comes from some ideal target function $f(X) \in span + noise$

$$y = f(X) + noise$$
$$(I - H)noise = y - \hat{y}$$

- ○ Substituting into definition of $E_{in}$

$$E_{in}(w_{LIN}) = \frac{1}{N} \|y - \hat{y}\|^2$$

$$= \frac{1}{N} \|(I - H)noise\|^2$$

$$= \frac{1}{N} trace(I - H)\|noise\|^2$$

$$= \frac{1}{N}(N - (d+1))\|noise\|^2$$

- Averaging across all possible training samples of size N from the population results in the analytical guarantee of linear regression

$$\overline{E_{in}} = \text{noise level} \cdot (1 - \frac{d+1}{N})$$

$$\overline{E_{out}} = \text{noise level} \cdot (1 + \frac{d+1}{N})$$

5. The learning curve

- Both in-sample and out-of-sample errors converge to noise leve $\sigma^2$
- Generalization error $E_{out} - E_{in}$ is bounded
  - With respect to the same ideal target function + noise
  - The bounded difference can be expressed as function of VC dimension $d$ and sample size $N$, $\frac{2(d+1)}{N}$