# Linear Regression

Shusen Wang

# Warm-up: Vector and Matrix

# Vector and Matrix

Vector ($n$-dim)

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

Matrix ($n \times d$)

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1d} \\ a_{21} & a_{22} & \cdots & a_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix}$$
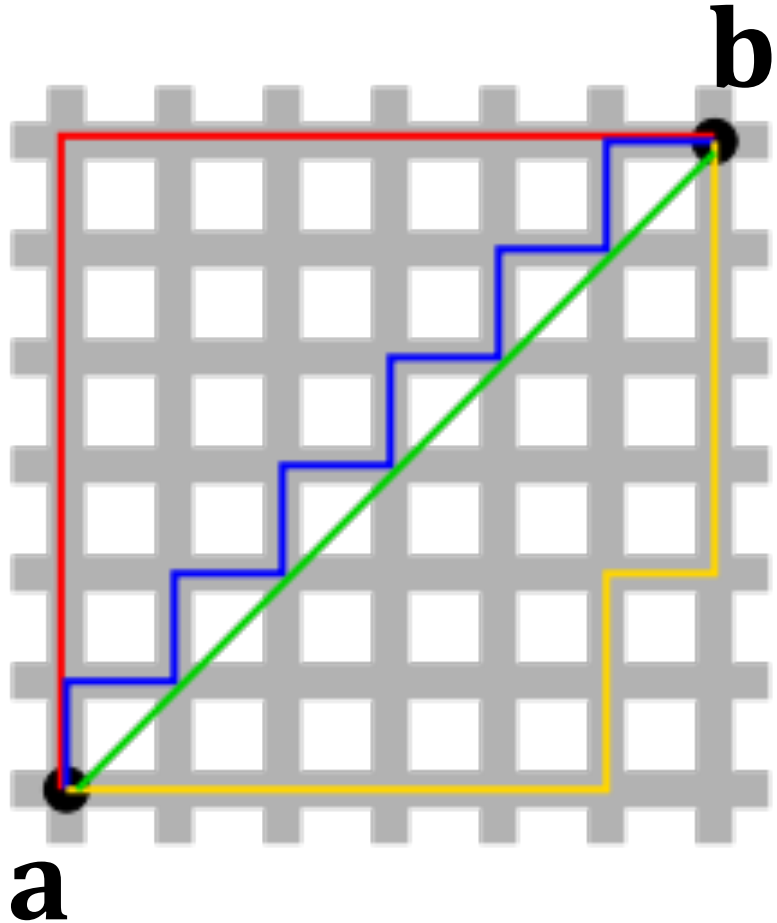
Row and columns

$$\mathbf{A} = \begin{bmatrix} \mathbf{a}_{:1} & \mathbf{a}_{:2} & \cdots & \mathbf{a}_{:d} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{1:} \\ \mathbf{a}_{2:} \\ \vdots \\ \mathbf{a}_{n:} \end{bmatrix}$$

# Vector Norms

- The $\ell_p$ norm: $\|\mathbf{x}\|_p := \left( \sum_i |x_i|^p \right)^{1/p}$.

- The $\ell_2$ norm: $\|\mathbf{x}\|_2 = \left( \sum_i x_i^2 \right)^{1/2}$ (the Euclidean norm).

- The $\ell_1$ norm $\|\mathbf{x}\|_1 = \sum_i |x_i|$.

- The $\ell_\infty$ norm is defined by $\|\mathbf{x}\|_\infty = \max_i |x_i|$.

# Vector Norms



- The $\ell_2$-distance (Euclidean distance): $\left|\left|\mathbf{a} - \mathbf{b}\right|\right|_2$ (green line)

- The $\ell_1$-distance (Manhattan distance): $\left|\left|\mathbf{a} - \mathbf{b}\right|\right|_1$ (red, blue, yellow lines)

# Transpose and Rank

**Transpose**:

$$\begin{bmatrix} 6 & 4 & 24 \\ 1 & -9 & 8 \end{bmatrix}^T = \begin{bmatrix} 6 & 1 \\ 4 & -9 \\ 24 & 8 \end{bmatrix}$$

**Square matrix**:  a matrix with the same number of rows and columns.

**Symmetric**:  a square matrix $\mathbf{A}$ is symmetric if $\mathbf{A}^T = \mathbf{A}$.

**Rank**:  the number of linearly independent rows (or columns).

**Full rank**:  a square matrix is full rank if the rank equals to #columns.

# Eigenvalue Decomposition

- Let $\mathbf{A}$ be any $n \times n$ symmetric matrix.

- Eigenvalue decomposition: $\mathbf{A} = \sum_{i=1}^{n} \lambda_i \mathbf{v}_i \mathbf{v}_i^T$.

- Eigenvalues satisfy $|\lambda_1| \geq |\lambda_2| \geq \cdots \geq |\lambda_n|$.

- Eigenvectors satisfy $\mathbf{v}_i^T \mathbf{v}_j = 0$ for all $i \neq j$.

- $\mathbf{A}$ is full rank $\longleftrightarrow$ all the eigenvalues are nonzero.

# Vector and Matrix Derivatives

# Derivative of Scalar w.r.t. Scalar

Examples:

- $y = x^2; \quad \dfrac{dy}{dx} = 2x.$

- $y = e^x; \quad \dfrac{dy}{dx} = e^x.$

# Derivative of Vector w.r.t. Scalar

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a scalar $x \in \mathbb{R}$:

$$\frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} \frac{\partial y_1}{\partial x} \\ \frac{\partial y_2}{\partial x} \\ \vdots \\ \frac{\partial y_n}{\partial x} \end{bmatrix}$$

- Example:

$$\mathbf{y} = \begin{bmatrix} 3x^2 \\ x+1 \\ \log x \\ e^x \end{bmatrix}, \qquad \frac{\partial \mathbf{y}}{\partial x} = \begin{bmatrix} 6x \\ 1 \\ 1/x \\ e^x \end{bmatrix}$$

# Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial\, y}{\partial\, \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$

- Example 1:

$$y = \|\mathbf{x}\|_2^2 = \sum_{i=1}^{m} x_i^2, \qquad \frac{\partial\, y}{\partial\, \mathbf{x}} = 2\mathbf{x}.$$

# Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$

- Example 2:

$$y = \mathbf{x}^T \mathbf{z} = \sum_{i=1}^{m} x_i z_i, \qquad \frac{\partial y}{\partial \mathbf{x}} = \mathbf{z}.$$

# Derivative of Scalar w.r.t. Vector

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \\ \vdots \\ \frac{\partial y}{\partial x_m} \end{bmatrix}$$
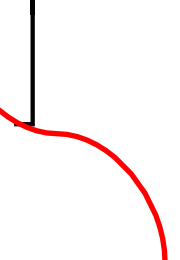
- Example 3:

$$y = \sum_{i=1}^{m} \log(1 + e^{-x_i}), \qquad \frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial \log(1+e^{-x_1})}{\partial x_1} \\ \vdots \\ \frac{\partial \log(1+e^{-x_m})}{\partial x_m} \end{bmatrix} = \begin{bmatrix} -\frac{1}{1+e^{x_1}} \\ \vdots \\ -\frac{1}{1+e^{x_m}} \end{bmatrix}$$

# Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix}$$

$m{\times}n$ matrix

The $(i, j)$-th entry is $\frac{\partial y_j}{\partial x_i}$

- Example 1:

$$\frac{\partial \mathbf{x}}{\partial \mathbf{x}} = \underbrace{\begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix}}_{m{\times}m}$$

# Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \qquad m{\times}n \text{ matrix}$$

- Example 2:

$$\mathbf{y} = \begin{bmatrix} a_1 x_1^2 \\ a_2 x_2^2 \\ \vdots \\ a_m x_m^2 \end{bmatrix} \in \mathbb{R}^m, \qquad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \underbrace{\begin{bmatrix} 2a_1 x_1 & 0 & \cdots & 0 \\ 0 & 2a_2 x_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 2a_m x_m \end{bmatrix}}_{m \times m}$$

# Derivative of Vector w.r.t. Vector

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a vector $\mathbf{x} \in \mathbb{R}^m$:
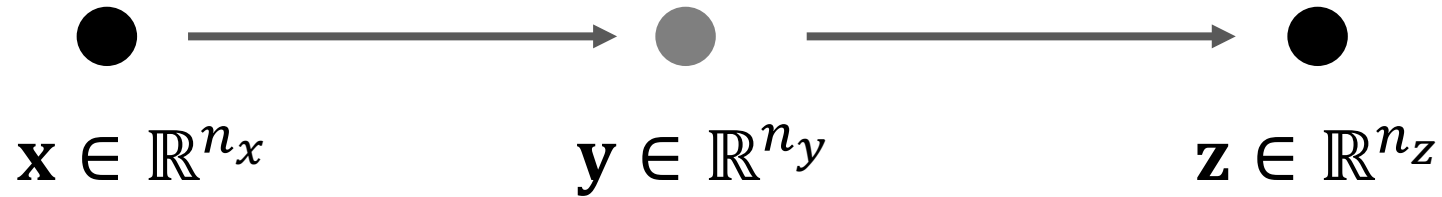
$$\frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_2}{\partial x_1} & \cdots & \frac{\partial y_n}{\partial x_1} \\ \frac{\partial y_1}{\partial x_2} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_2} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_1}{\partial x_m} & \frac{\partial y_2}{\partial x_m} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} \qquad m{\times}n \text{ matrix}$$

- Example 3:

$$\mathbf{A} \in \mathbb{R}^{n \times m}, \qquad \mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^n, \qquad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \mathbf{A}^T \in \mathbb{R}^{m \times n}$$

# Chain Rule

- Let $\mathbf{z} \in \mathbb{R}^{n_z}$ be a function of $\mathbf{y} \in \mathbb{R}^{n_y}$ and $\mathbf{y}$ be a function of $\mathbf{x} \in \mathbb{R}^{n_x}$.

$$\underbrace{\frac{d\mathbf{z}}{d\mathbf{x}}}_{n_x \times n_z} = \underbrace{\frac{d\mathbf{y}}{d\mathbf{x}}}_{n_x \times n_y} \underbrace{\frac{d\mathbf{z}}{d\mathbf{y}}}_{n_y \times n_z}$$

# Derivative of Scalar w.r.t. Matrix

- The derivative of a scalar $y \in \mathbb{R}$ w.r.t. a matrix $\mathbf{Z} \in \mathbb{R}^{p \times q}$:

  1. Vectorization: $\mathbf{x} = \text{vec}(\mathbf{Z}) \in \mathbb{R}^{pq \times 1}$.

  2. Compute $\dfrac{\partial y}{\partial \mathbf{x}} \in \mathbb{R}^{pq \times 1}$.

  3. Reshape the resulting $pq \times 1$ vector to $p \times q$ matrix.

# Derivative of Vector w.r.t. Matrix

- The derivative of a vector $\mathbf{y} \in \mathbb{R}^n$ w.r.t. a matrix $\mathbf{Z} \in \mathbb{R}^{p \times q}$ :

  1. Vectorization: $\mathbf{x} = \mathrm{vec}(\mathbf{Z}) \in \mathbb{R}^{pq \times 1}$.

  2. Compute $\dfrac{\partial\, \mathbf{y}}{\partial\, \mathbf{x}} \in \mathbb{R}^{pq \times n}$.

  3. Reshape the resulting $pq \times n$ matrix to $p \times q \times n$ tensor.

# Warm-up: Optimization

# Optimization: Basics

Optimization problem:     $\min_{\mathbf{w}} f(\mathbf{w})$;    s. t. $\mathbf{w} \in \mathcal{C}$ .

- $\mathbf{w} = [w_1, \cdots, w_d]$ :  optimization variables
- $f : \mathbb{R}^d \mapsto \mathbb{R}$  :  objective function
- $\mathcal{C}$ (a subset of $\mathbb{R}^d$) :  feasible set

# Optimization: Basics

Optimization problem: $\min_{\mathbf{w}} f(\mathbf{w})$; s. t. $\mathbf{w} \in \mathcal{C}$.

Constraint

- $\mathbf{w} = [w_1, \cdots, w_d]$ : optimization variables
- $f : \mathbb{R}^d \mapsto \mathbb{R}$ : objective function
- $\mathcal{C}$ (a subset of $\mathbb{R}^d$) : feasible set

# Optimization: Basics

Optimization problem: $\quad \min\limits_{\mathbf{w}} f(\mathbf{w}); \quad \text{s.t. } \mathbf{w} \in \mathcal{C}.$

Optimal solution: $\quad \mathbf{w}^\star = \underset{\mathbf{w} \in \mathcal{C}}{\operatorname{argmin}} f(\mathbf{w}).$

- $f(\mathbf{w}^\star) \leq f(\mathbf{w})$ for all the vectors $\mathbf{w}$ in the set $\mathcal{C}$.
- $\mathbf{w}^\star$ may not exist, e.g., $\mathcal{C}$ is the empty set.
- If $\mathbf{w}^\star$ exists, it may not be unique.

# Least Squares Regression

# Linear Regression

**Input:** vectors $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \cdots, y_n \in \mathbb{R}$

**Output:** a vector $\color{red}{\mathbf{w}} \in \mathbb{R}^d$ and scalar $\color{red}{b} \in \mathbb{R}$ such that $\mathbf{x}_i^T \color{red}{\mathbf{w}} + \color{red}{b} \approx y_i$.

1-dim ($d = 1$) example:

Solution:
$$y_i \approx \color{red}{0.15}\, x_i + \color{red}{5.0}$$

# Linear Regression

**Input:** vectors $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{R}^d$ and labels $y_1, \cdots, y_n \in \mathbb{R}$

**Output:** a vector $\mathbf{w} \in \mathbb{R}^d$ and scalar $b \in \mathbb{R}$ such that $\mathbf{x}_i^T \mathbf{w} + b \approx y_i$.

**Question** (regard training):
how to compute $\mathbf{w}$ and $b$?
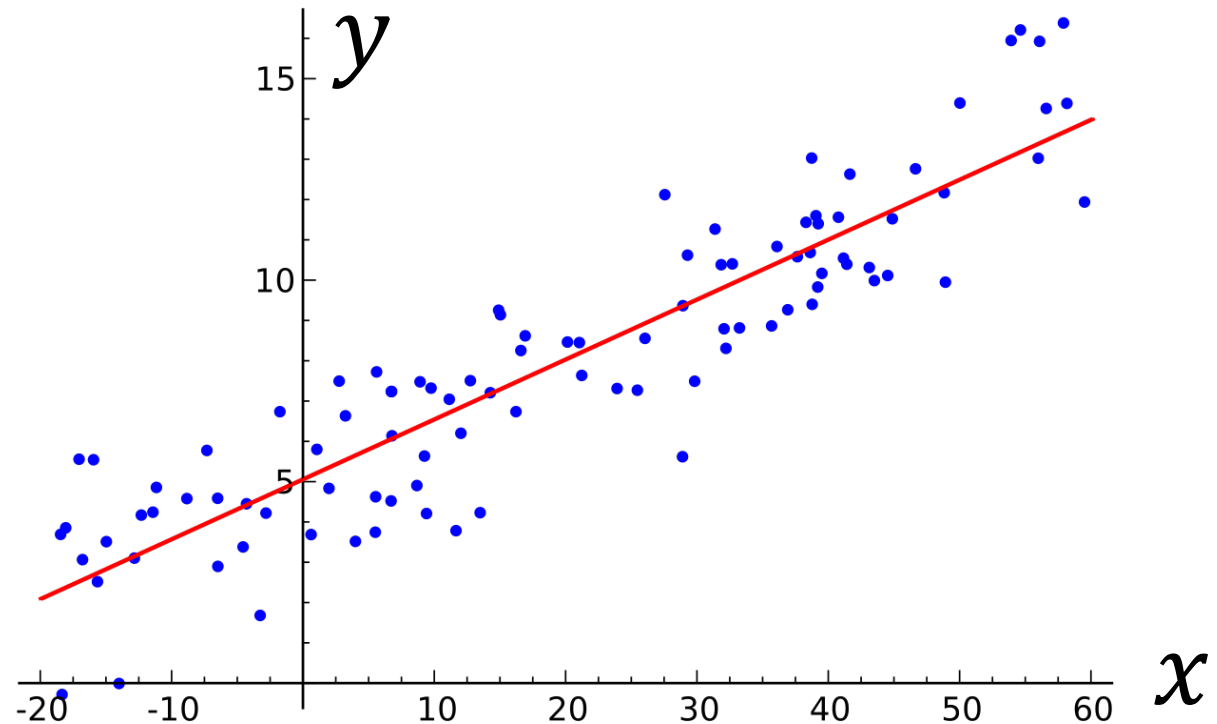
# Least Squares Regression

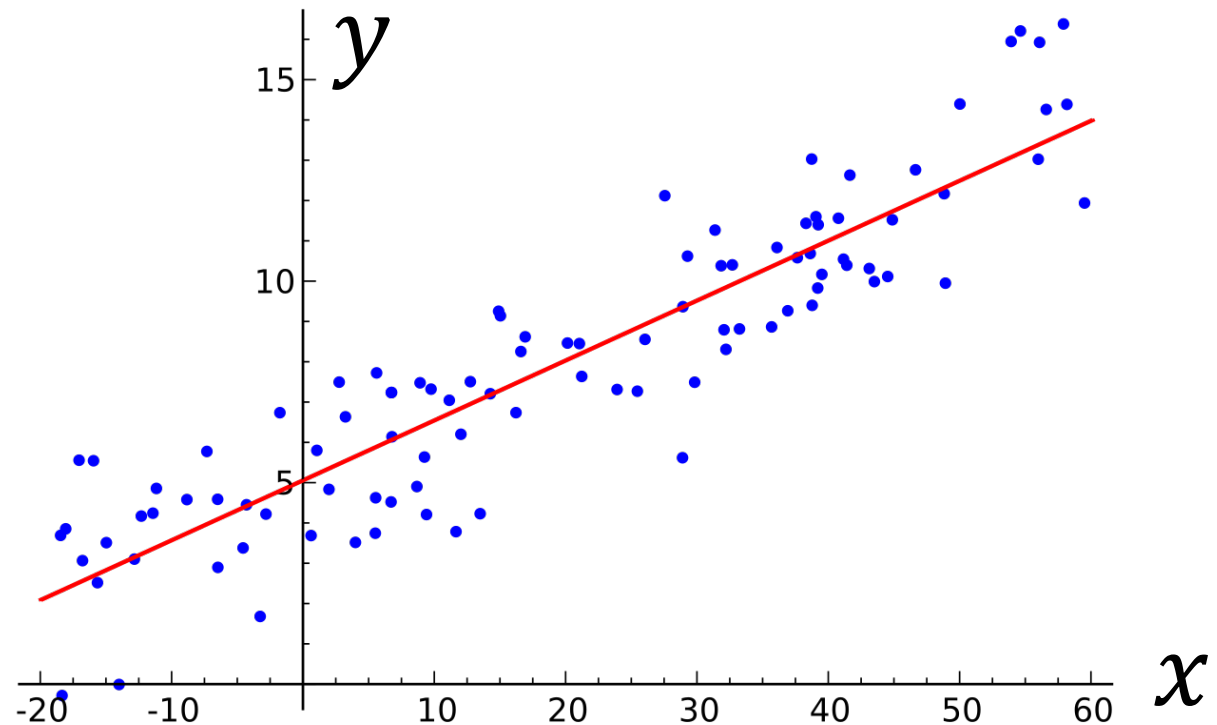**Input:**   vectors $\mathbf{x}_1, \cdots, \mathbf{x}_n \in \mathbb{R}^d$  and  labels $y_1, \cdots, y_n \in \mathbb{R}$

**Output:** a vector $\mathbf{w} \in \mathbb{R}^d$ and scalar $b \in \mathbb{R}$ such that  $\mathbf{x}_i^T \mathbf{w} + b \approx y_i$.

**Method**: least squares regression.

- The optimization model:

$$\min_{\mathbf{w}, b} \sum_{i=1}^{n} \left( \mathbf{x}_i^T \mathbf{w} + b - y_i \right)^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w}, b} \sum_{i=1}^{n} \left( \mathbf{x}_i^T \mathbf{w} + b - y_i \right)^2$$

Intercept (or bias)

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w},b} \sum_{i=1}^{n} \left( \mathbf{x}_i^T \mathbf{w} + b - y_i \right)^2$$

$$\bar{\mathbf{x}}_i = \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix}$$

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w},b} \sum_{i=1}^{n}\left(\mathbf{x}_i^T \mathbf{w} + b - y_i\right)^2$$

$$\bar{\mathbf{x}}_i = \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \qquad \bar{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w}, b} \sum_{i=1}^{n} \left( \mathbf{x}_i^T \mathbf{w} + b - y_i \right)^2$$

$$\bar{\mathbf{x}}_i = \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \qquad \bar{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

$$\bar{\mathbf{x}}_i^T \bar{\mathbf{w}} = \mathbf{x}_i^T \mathbf{w} + b$$

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w},b} \sum_{i=1}^{n} \left( \underbrace{\mathbf{x}_i^T \mathbf{w} + b}_{= \bar{\mathbf{x}}_i^T \bar{\mathbf{w}}} - y_i \right)^2$$

$$\bar{\mathbf{x}}_i = \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \qquad \bar{\mathbf{w}} = \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

$$\bar{\mathbf{x}}_i^T \bar{\mathbf{w}} = \mathbf{x}_i^T \mathbf{w} + b$$

# Least Squares Regression

- The optimization model:

$$\min_{\mathbf{w},b} \sum_{i=1}^{n} \left( \underbrace{\mathbf{x}_i^T \mathbf{w} + b}_{= \bar{\mathbf{x}}_i^T \bar{\mathbf{w}}} - y_i \right)^2$$

$$\updownarrow$$

$$\min_{\bar{\mathbf{w}} \in \mathbb{R}^{d+1}} \sum_{i=1}^{n} \left( \bar{\mathbf{x}}_i^T \bar{\mathbf{w}} - y_i \right)^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

$$\overline{\mathbf{X}} = \begin{bmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \mathbf{x}_3^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_n^T & 1 \end{bmatrix} \qquad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{bmatrix}$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

$$\overline{\mathbf{X}}\overline{\mathbf{w}} = \begin{bmatrix} \overline{\mathbf{x}}_1^T \overline{\mathbf{w}} \\ \overline{\mathbf{x}}_2^T \overline{\mathbf{w}} \\ \overline{\mathbf{x}}_3^T \overline{\mathbf{w}} \\ \vdots \\ \overline{\mathbf{x}}_n^T \overline{\mathbf{w}} \end{bmatrix}$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

$$\overline{\mathbf{X}}\overline{\mathbf{w}} = \begin{bmatrix} \overline{\mathbf{x}}_1^T \overline{\mathbf{w}} \\ \overline{\mathbf{x}}_2^T \overline{\mathbf{w}} \\ \overline{\mathbf{x}}_3^T \overline{\mathbf{w}} \\ \vdots \\ \overline{\mathbf{x}}_n^T \overline{\mathbf{w}} \end{bmatrix} \qquad \overline{\mathbf{X}}\overline{\mathbf{w}} - \mathbf{y} = \begin{bmatrix} \overline{\mathbf{x}}_1^T \overline{\mathbf{w}} - y_1 \\ \overline{\mathbf{x}}_2^T \overline{\mathbf{w}} - y_2 \\ \overline{\mathbf{x}}_3^T \overline{\mathbf{w}} - y_3 \\ \vdots \\ \overline{\mathbf{x}}_n^T \overline{\mathbf{w}} - y_n \end{bmatrix}$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

$$\left|\left| \overline{\mathbf{X}} \overline{\mathbf{w}} - \mathbf{y} \right|\right|_2^2 \;=\; \left|\left| \begin{bmatrix} \overline{\mathbf{x}}_1^T \overline{\mathbf{w}} - y_1 \\ \overline{\mathbf{x}}_2^T \overline{\mathbf{w}} - y_2 \\ \overline{\mathbf{x}}_3^T \overline{\mathbf{w}} - y_3 \\ \vdots \\ \overline{\mathbf{x}}_n^T \overline{\mathbf{w}} - y_n \end{bmatrix} \right|\right|_2^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

$$\left|\left| \overline{\mathbf{X}} \overline{\mathbf{w}} - \mathbf{y} \right|\right|_2^2 = \left|\left| \begin{bmatrix} \overline{\mathbf{x}}_1^T \overline{\mathbf{w}} - y_1 \\ \overline{\mathbf{x}}_2^T \overline{\mathbf{w}} - y_2 \\ \overline{\mathbf{x}}_3^T \overline{\mathbf{w}} - y_3 \\ \vdots \\ \overline{\mathbf{x}}_n^T \overline{\mathbf{w}} - y_n \end{bmatrix} \right|\right|_2^2 = \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2.$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \sum_{i=1}^{n} \left( \overline{\mathbf{x}}_i^T \overline{\mathbf{w}} - y_i \right)^2$$

Matrix form: 
$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}} \, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \quad \left\|\overline{\mathbf{X}} \, \overline{\mathbf{w}} - \mathbf{y}\right\|_2^2$$

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

**Tasks**

**Methods**

**Algorithms**

Linear Regression

Least Squares Regression

LASSO

Least Absolute Deviations

**?**

# Least Squares Regression

- The optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

## Tasks

Linear Regression

## Methods

Least Squares Regression

LASSO

Least Absolute Deviations

## Algorithms

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

Gradient: $\dfrac{\partial \left\| \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \mathbf{y} \right\|_2^2}{\partial \overline{\mathbf{w}}} = 2(\overline{\mathbf{X}}^T \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \overline{\mathbf{X}}^T \mathbf{y})$

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

Gradient: $\dfrac{\partial \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2}{\partial \overline{\mathbf{w}}} = 2(\overline{\mathbf{X}}^T \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \overline{\mathbf{X}}^T \mathbf{y}) \boxed{= \mathbf{0}}$

$1^{\text{st}}$-order optimality condition

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

Gradient: $\dfrac{\partial \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2}{\partial\, \overline{\mathbf{w}}} = 2(\overline{\mathbf{X}}^T \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \overline{\mathbf{X}}^T \mathbf{y}) = \mathbf{0}$

Normal equation: $\overline{\mathbf{X}}^T \overline{\mathbf{X}}\, \overline{\mathbf{w}}^\star = \overline{\mathbf{X}}^T \mathbf{y}$

Assume $\overline{\mathbf{X}}^T \overline{\mathbf{X}}$ is full rank.

Analytical solution: $\overline{\mathbf{w}}^\star = (\overline{\mathbf{X}}^T \overline{\mathbf{X}})^{-1} \overline{\mathbf{X}}^T \mathbf{y}$

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

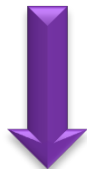# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

Gradient: $\dfrac{\partial \left\| \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \mathbf{y} \right\|_2^2}{\partial\,\overline{\mathbf{w}}} = 2(\overline{\mathbf{X}}^T \overline{\mathbf{X}}\,\overline{\mathbf{w}} - \overline{\mathbf{X}}^T \mathbf{y}) = \mathbf{0}$

Gradient descent repeats:

1. Compute gradient: $\mathbf{g}_t = \overline{\mathbf{X}}^T \overline{\mathbf{X}}\,\overline{\mathbf{w}}_t - \overline{\mathbf{X}}^T \mathbf{y}$
2. Update: $\overline{\mathbf{w}}_{t+1} = \overline{\mathbf{w}}_t - \alpha_t\,\mathbf{g}_t$

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left|\left|\overline{\mathbf{X}}\,\overline{\mathbf{w}} - \mathbf{y}\right|\right|_2^2$$

Convergence:  after $O\left(\kappa \log \frac{1}{\epsilon}\right)$ iterations,

$$\left|\left|\overline{\mathbf{X}}\,(\overline{\mathbf{w}}_t - \overline{\mathbf{w}}^\star)\right|\right|_2 \leq \epsilon \left|\left|\overline{\mathbf{X}}\,(\overline{\mathbf{w}}_0 - \overline{\mathbf{w}}^\star)\right|\right|_2.$$

$\kappa = \dfrac{\lambda_{\max}(\overline{\mathbf{X}}^T\overline{\mathbf{X}})}{\lambda_{\min}(\overline{\mathbf{X}}^T\overline{\mathbf{X}})}$ is the condition number.

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}}\, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

Convergence: after $O\left(\sqrt{\kappa}\, \log \frac{1}{\epsilon}\right)$ iterations,

$$\left\| \overline{\mathbf{X}}\, (\overline{\mathbf{w}}_t - \overline{\mathbf{w}}^\star) \right\|_2 \le \epsilon \left\| \overline{\mathbf{X}}\, (\overline{\mathbf{w}}_0 - \overline{\mathbf{w}}^\star) \right\|_2 .$$

$\kappa = \dfrac{\lambda_{\max}(\overline{\mathbf{X}}^T \overline{\mathbf{X}})}{\lambda_{\min}(\overline{\mathbf{X}}^T \overline{\mathbf{X}})}$  is the condition number.

The pseudo-code of CG is available at the Wikipedia.

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Least Squares Regression

- Solve the optimization model:

$$\min_{\overline{\mathbf{w}} \in \mathbb{R}^{d+1}} \left\| \overline{\mathbf{X}} \, \overline{\mathbf{w}} - \mathbf{y} \right\|_2^2$$

**Tasks**

Linear Regression

**Methods**

Least Squares Regression

LASSO

Least Absolute Deviations

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient (CG)

# Solve Least Squares in Python

# 1. Load Data

```python
from keras.datasets import boston_housing

(x_train, y_train), (x_test, y_test) = boston_housing.load_data()

print('shape of x_train: ' + str(x_train.shape))
print('shape of x_test: ' + str(x_test.shape))
print('shape of y_train: ' + str(y_train.shape))
print('shape of y_test: ' + str(y_test.shape))
```

```
shape of x_train: (404, 13)
shape of x_test: (102, 13)
shape of y_train: (404,)
shape of y_test: (102,)
```

# 2. Add A Feature

```python
import numpy

n, d = x_train.shape
xbar_train = numpy.concatenate((x_train, numpy.ones((n, 1))),
                               axis=1)


print('shape of x_train: ' + str(x_train.shape))
print('shape of xbar_train: ' + str(xbar_train.shape))
```

```
shape of x_train: (404, 13)
shape of xbar_train: (404, 14)
```

# 3. Solve the Least Squares

Analytical solution: $\bar{\mathbf{w}} = (\bar{\mathbf{X}}^T \bar{\mathbf{X}})^{-1} \bar{\mathbf{X}}^T \mathbf{y}$

```python
xx = numpy.dot(xbar_train.T, xbar_train)
xx_inv = numpy.linalg.pinv(xx)
xy = numpy.dot(xbar_train.T, y_train)
w = numpy.dot(xx_inv, xy)
```

# 3. Solve the Least Squares

Analytical solution: $\overline{\mathbf{w}} = \boxed{(\overline{\mathbf{X}}^T \overline{\mathbf{X}})^{-1}} \overline{\mathbf{X}}^T \mathbf{y}$

```
xx = numpy.dot(xbar_train.T, xbar_train)
xx_inv = numpy.linalg.pinv(xx)
xy = numpy.dot(xbar_train.T, y_train)
w = numpy.dot(xx_inv, xy)
```

# 3. Solve the Least Squares

Analytical solution: $\overline{\mathbf{w}} = \boxed{(\overline{\mathbf{X}}^T \overline{\mathbf{X}})^{-1}} \boxed{\overline{\mathbf{X}}^T \mathbf{y}}$

```
xx = numpy.dot(xbar_train.T, xbar_train)
xx_inv = numpy.linalg.pinv(xx)
xy = numpy.dot(xbar_train.T, y_train)
w = numpy.dot(xx_inv, xy)
```

# 3. Solve the Least Squares

Training Mean Squared Error (MSE): $\frac{1}{n}\left\|\mathbf{y} - \overline{\mathbf{X}}\overline{\mathbf{w}}\right\|_2^2$

```python
y_lsr = numpy.dot(xbar_train, w)
diff = y_lsr - y_train
mse = numpy.mean(diff * diff)
print('Train MSE: ' + str(mse))
```

Train MSE: 22.00480083834814

# Linear Regression for Housing Price

Linear Regressor $\overline{\mathbf{w}}$

Train

label, $y_i$ ← **Price** — Selling price of the property ($1,000)

features, $\mathbf{x}_i$
- **Beds** — Number of bedrooms in the house
- **Baths** — Number of bathrooms in the house
- **Square Feet** — Size of the house in square feet
- **Miles to Resort** — Miles from the property to the downtown resort area
- **Miles to Base** — Miles from the property to the base of the ski resort's mountain
- **Acres** — Lot size in number of acres
- **Cars** — Number of cars that will fit into the garage
- **Years Old** — Age of the house, in years, at the time it was listed
- **DoM** — Number of days the house was on the market before it sold

# Linear Regression for Housing Price



Features of a House, $\mathbf{x}'$
➡ Extend it to $\bar{\mathbf{x}}'$

Linear Regressor $\bar{\mathbf{w}}$

Predict

Price:
$\bar{\mathbf{w}}^T \bar{\mathbf{x}}' = \$500\text{K}$

# 4. Make Prediction for Test Samples

- Add a feature to the test feature matrix: $\mathbf{X}_{\text{test}} \rightarrow \overline{\mathbf{X}}_{\text{test}}$.

- Make prediction by: $\mathbf{y}_{\text{pred}} = \overline{\mathbf{X}}_{\text{test}}\,\overline{\mathbf{w}}$.

```
n_test, _ = x_test.shape
xbar_test = numpy.concatenate((x_test, numpy.ones((n_test, 1))), axis=1)
y_pred = numpy.dot(xbar_test, w)
```

# 4. Make Prediction for Test Samples

- Add a feature to the test feature matrix: $\mathbf{X}_{\text{test}} \rightarrow \overline{\mathbf{X}}_{\text{test}}$.
- Make prediction by: $\mathbf{y}_{\text{pred}} = \overline{\mathbf{X}}_{\text{test}}\overline{\mathbf{w}}$.
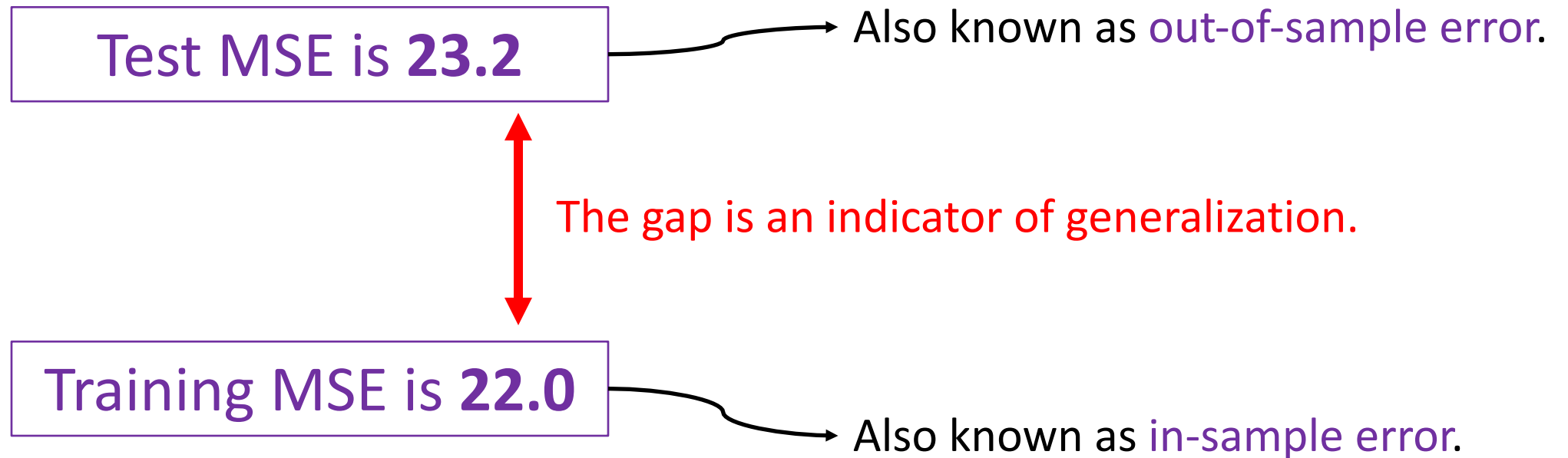- MSE (test): $\dfrac{1}{n_{\text{test}}}\left\|\mathbf{y}_{\text{pred}} - \mathbf{y}_{\text{test}}\right\|_2^2$

```python
# mean squared error (testing)

diff = y_pred - y_test
mse = numpy.mean(diff * diff)
print('Test MSE: ' + str(mse))
```

Test MSE: 23.195599256409857

Training MSE is **22.0**

# 4. Make Prediction for Test Samples

Test MSE is **23.2**

→ Also known as out-of-sample error.

↕ The gap is an indicator of generalization.

Training MSE is **22.0**

→ Also known as in-sample error.

# 5. Compare with Baseline

Trivial baseline:
- whatever the features are, the prediction is $\text{mean}(\mathbf{y})$.

```python
y_mean = numpy.mean(y_train)

diff = y_pred - y_mean
mse = numpy.mean(diff * diff)
print('Test MSE: ' + str(mse))
```

Test MSE: 57.38297638530044

Test MSE of least
squares is **23.19**

# Summary

- Linear regression problem.

- Least squares model.

- 3 algorithms for solving the model.

- Make predictions for never-seen-before test data.

- Evaluation of the model (training MSE and test MSE).

- Compare with baselines.