# Linear Regression

**Shusen Wang**

# Warm-up: Vector and Matrix

# Vector and Matrix

Vector ($n$-dim) $\qquad \mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$
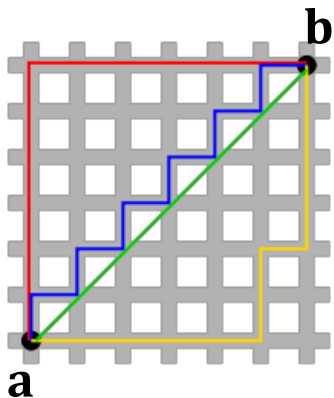
Matrix ($n{\times}d$) $\qquad \mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1d} \\ a_{21} & a_{22} & \cdots & a_{2d} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nd} \end{bmatrix}$

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

# 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):
  $$\left\|\mathbf{a} - \mathbf{b}\right\|_2 \text{ (green line)}$$

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

# 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.

# Warm-up: Optimization

# Optimization: Basics

Optimization problem: $\quad \min\limits_{\mathbf{w}} f(\mathbf{w}) \,; \quad$ 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: $\quad \min\limits_{\mathbf{w}} f(\mathbf{w})$ ; $\quad$ 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: $\min\limits_{\mathbf{w}} f(\mathbf{w})$ ; 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

- $\mathbf{w}^\star = \operatorname*{argmin}\limits_{\mathbf{w} \in \mathcal{C}} f(\mathbf{w})$ is the optimal solution to the problem
  - $f(\mathbf{w}^\star) \leq f(\mathbf{w})$ for all the vectors $\mathbf{w}$ in the set $\mathcal{C}$.

# Least Squares Regression

# The Linear Regression Task

**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 $\mathbf{x}_i^T \mathbf{w} + b \approx y_i$.

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

Solution:
$$y_i \approx 0.15 \, x_i + 5.0$$
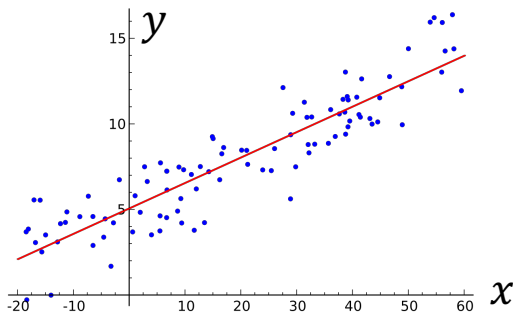
# The Linear Regression Task

**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):
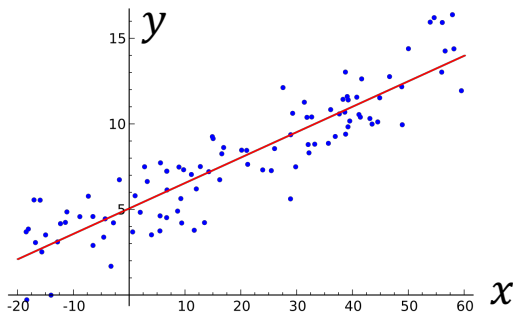how to compute $\mathbf{w}$ and $b$?

# The Linear Regression Task

**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:

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

# Least Squares Regression

- The optimization model:

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

# Least Squares Regression

- The optimization model:

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

Intercept

# Least Squares Regression

- The optimization model:

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

$$\updownarrow$$

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

- Define $\bar{\mathbf{x}}_i = [\mathbf{x}_i; 1] \in \mathbb{R}^{d+1}$
- Define $\overline{\mathbf{w}} = [\mathbf{w}, b] \in \mathbb{R}^{d+1}$
- ➔  $\mathbf{x}_i^T \mathbf{w} + b = \bar{\mathbf{x}}_i^T \overline{\mathbf{w}}$

# Least Squares Regression

- The optimization model:

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

$\updownarrow$

Matrix form:

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

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

$n \times (d+1)$ $\qquad\qquad$ $n \times 1$ $\qquad\qquad$ $n \times 1$

# Least Squares Regression

- The optimization model:

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

| Tasks | Methods | Algorithms |
|-------|---------|------------|
| Linear Regression | Least Squares Regression | ? |
| | LASSO | |
| | Least Absolute Deviations | |

# Least Squares Regression

- The optimization model:

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

| **Tasks** | **Methods** | **Algorithms** |
|---|---|---|
| Linear Regression | Least Squares Regression | Analytical Solution |
|  | LASSO | Gradient Descent (GD) |
|  | Least Absolute Deviations | Conjugate Gradient (CG) |

# Least Squares Regression

- Solve the optimization model:

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

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

| Algorithms |
| --- |
| Analytical Solution |
| Gradient Descent (GD) |
| Conjugate Gradient |

# Least Squares Regression

- Solve the optimization model:

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

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

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

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient

# Least Squares Regression

- Solve the optimization model:

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

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

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

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

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

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient

# Least Squares Regression

- Solve the optimization model:

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

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

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

# Least Squares Regression

- Solve the optimization model:

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

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

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

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

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient

# Least Squares Regression

- Solve the optimization model:

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

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

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

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

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

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient

# Least Squares Regression

- Solve the optimization model:

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

**Tasks**

Linear Regression

**Methods**

Least Squares Regression

LASSO

Least Absolute Deviations

**Algorithms**

Analytical Solution

Gradient Descent (GD)

Conjugate Gradient

# Linear Regression for Housing Price



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

Train

label, $y_i$

features, $\mathbf{x}_i$

| | |
|---|---|
| **Price** | Selling price of the property ($1,000) |
| **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: $\bar{\mathbf{w}}^T \bar{\mathbf{x}}' = \$0.5M$

Train

label, $y_i$

features, $\mathbf{x}_i$

| | |
|---|---|
| **Price** | Selling price of the property ($1,000) |
| **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 |