

ISE 589/OR591

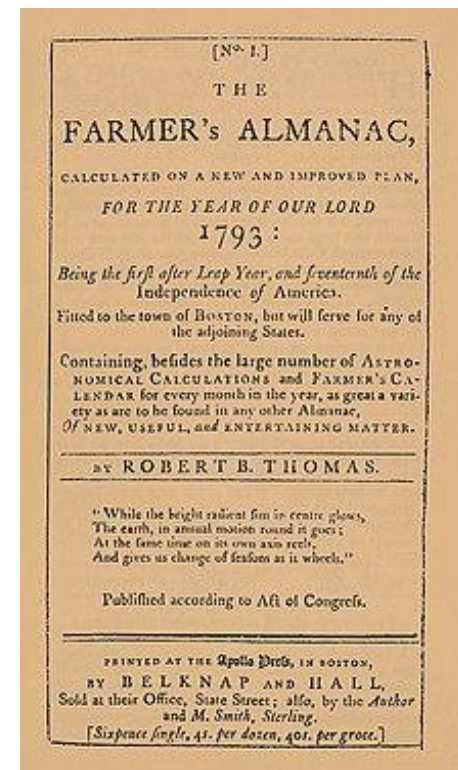
LECTURE 4 - PREDICTION

Unconstrained optimization models for machine learning

1. Least squares estimation
2. Linear regression

Prediction

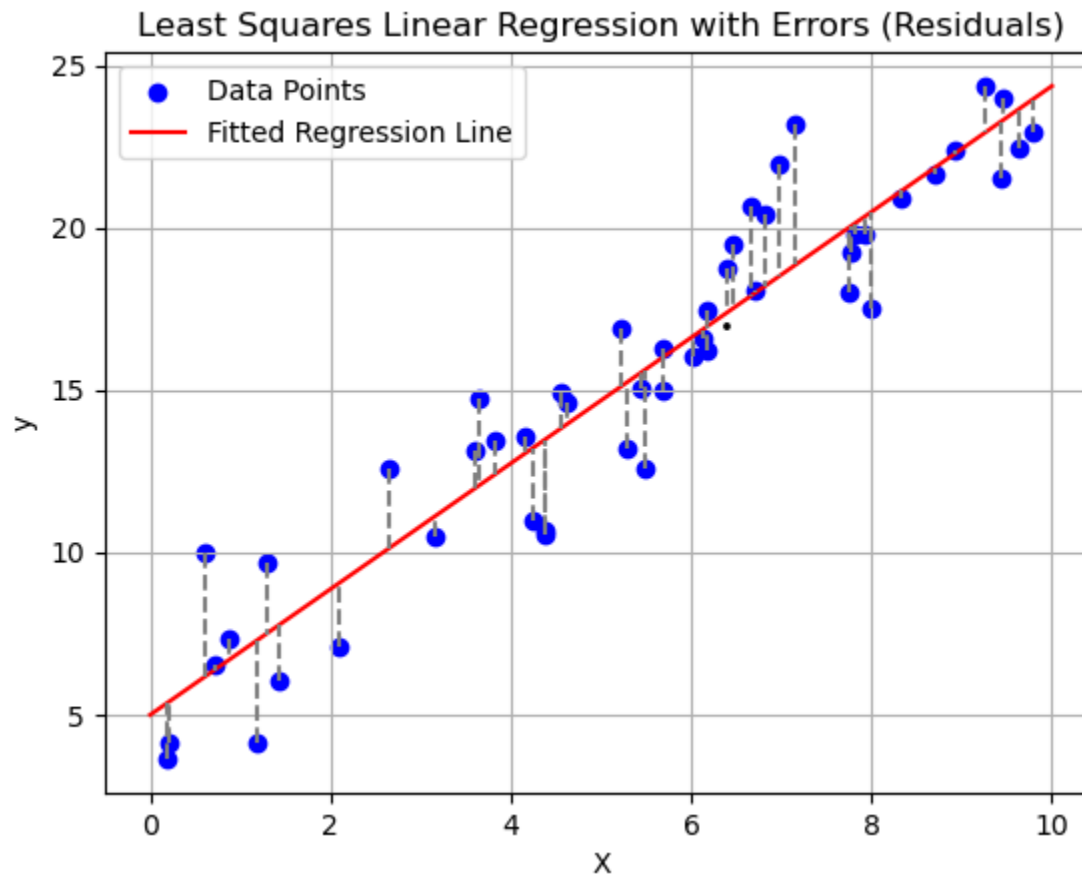
- Prediction (forecasting) is a process of finding the hidden relationship between Input data and output data for future telling
- It is often understood as learning from past experiences to tell outcomes of future events



Prediction

- **Predictive analytics**, or **predictive AI**, encompasses a variety of [statistical](#) techniques from [data mining](#), [predictive modeling](#), and [machine learning](#) that analyze current and historical facts to make [predictions](#) about future or otherwise unknown events.
https://en.wikipedia.org/wiki/Predictive_analytics
- Some useful approaches
 - [Moving average](#) models (such as single moving average, autoregressive integrated moving average, [ARIMA](#))
 - [Time series](#) models
 - [Least-squares](#) models
 - [Linear regression](#)
 - [Support vector regression \(SVR\)](#)

Linear regression



<https://medium.com/physics-and-machine-learning/misconceptions-about-least-square-regression-1131841d240f>

Basics of data and linear algebra – vector norm

- (norm) A **vector norm** $\|\cdot\|$ is a function from \mathbb{R}^n to \mathbb{R} such that (i) $\|\mathbf{x}\| \geq 0$ for any $\mathbf{x} \in \mathbb{R}^n$, and $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$; (ii) $\|\mathbf{x}^1 + \mathbf{x}^2\| \leq \|\mathbf{x}^1\| + \|\mathbf{x}^2\|$ for any $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^n$; (iii) $\|\alpha\mathbf{x}\| = \alpha\|\mathbf{x}\|$ for any $\alpha \in \mathbb{R}$ and $\mathbf{x} \in \mathbb{R}^n$.
- (commonly used norm)
 - l_0 norm: $\|\mathbf{x}\|_0 = \sum_{i=1}^n \mathbb{1}(x_i \neq 0)$;
 - l_1 norm: $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$;
 - l_2 norm: $\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$;
 - l_p norm: $\|\mathbf{x}\|_p = \left(\sum_{i=1}^n |x_i|^p\right)^{\frac{1}{p}}$, $p \geq 1$;
 - l_∞ norm: $\|\mathbf{x}\|_\infty = \max_{i=1, \dots, n} |x_i|$.

Gradient functions with 2-norm functions

Gradient of function $f(\cdot)$ at point $\mathbf{x} \in \mathbb{R}^n$: $\nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix}$

1. $\nabla(\mathbf{c}^T \mathbf{x}) = \mathbf{c}, \quad \forall \mathbf{c}, \mathbf{x} \in \mathbb{R}^n$

2. $\nabla(\|\mathbf{x}\|_2) = \frac{\mathbf{x}}{\|\mathbf{x}\|_2}, \quad \forall \mathbf{x} \in \mathbb{R}^n$

3. $\nabla(\|\mathbf{x}\|_2^2) = 2\mathbf{x}, \quad \forall \mathbf{x} \in \mathbb{R}^n$

4. $\nabla(\mathbf{x}^T M \mathbf{x}) = (M + M^T)\mathbf{x}, \quad \forall M \in \mathbf{M}_{n \times n}$ and $\mathbf{x} \in \mathbb{R}^n$

5. $\nabla(\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2) = 2(\mathbf{A}^T \mathbf{A})\mathbf{x} - 2\mathbf{A}^T \mathbf{b}, \quad \forall \mathbf{A} \in M_{m \times n}, \mathbf{b} \in \mathbb{R}^m$
and $\mathbf{x} \in \mathbb{R}^n$

System of linear equations

- Given a system of m linear equations in n real variables:

$$\begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \iff A\mathbf{x} = \mathbf{b}$$

- Facts from linear algebra:

Case 1. When $m \leq n$ and A has full rank (rows of A are linear independent), the system has feasible solutions

Case 2. When $m > n$, too many equations to satisfy may lead to **infeasibility** of the system!

Least squares optimization model

- **Case 1.** When $m \leq n$ and A has full rank, the system $Ax = b$ has **feasible solutions**

Question: **How to find a feasible solution?**

- Linear programming model (**how?**)
- Least squares optimization model

- **Case 2.** When $m > n$, too many equations to satisfy may lead to **infeasibility** of the system!

Question: **What to do?**

- Find an approximate solution

Least squares optimization model

- **Case 1.** When $m \leq n$ and A has full rank, the system $A\mathbf{x} = \mathbf{b}$ has **feasible solutions**

- **Task:** Find a feasible solution to the system

- **Reasoning:** (i) \mathbf{x} is feasible if and only if $A\mathbf{x} - \mathbf{b} = \mathbf{0}$
(ii) $\|\mathbf{v}\|_2 = 0$ if and if **vector $\mathbf{v} = \mathbf{0}$**

- **Optimization model:**

$$\begin{aligned} &\text{Minimize } \|A\mathbf{x} - \mathbf{b}\|_2^2 \\ &\text{s.t. } \mathbf{x} \in \mathbb{R}^n \end{aligned}$$

Least squares optimization model

- Optimization model:

$$\text{Minimize } \|A\mathbf{x} - \mathbf{b}\|_2^2$$

$$\text{s.t. } \mathbf{x} \in \mathbb{R}^n$$

Solution method:

- Optimal solutions of a unconstrained convex optimization problem achieved at the **point with zero gradient!**

$$\nabla \|A\mathbf{x} - \mathbf{b}\|_2^2 = 0 \Leftrightarrow 2(A^T A)\mathbf{x} - 2A^T \mathbf{b} = 0$$

$$\Leftrightarrow (A^T A)\mathbf{x} = A^T \mathbf{b} \quad (\text{Normal Equation})$$

If $(A^T A)^{-1}$ exists, then $\mathbf{x} = (A^T A)^{-1} A^T \mathbf{b}$.

Otherwise, **solve the normal equation**, e.g., Cholesky factorization method for $A^T A$ being positive semi-definite.

Least squares optimization model

- **Case 2.** When $m > n$, too many equations to satisfy may lead to **infeasibility** of the system!
- **Task:** Find an approximate solution of the system

Question: How to **define an approximate solution?**

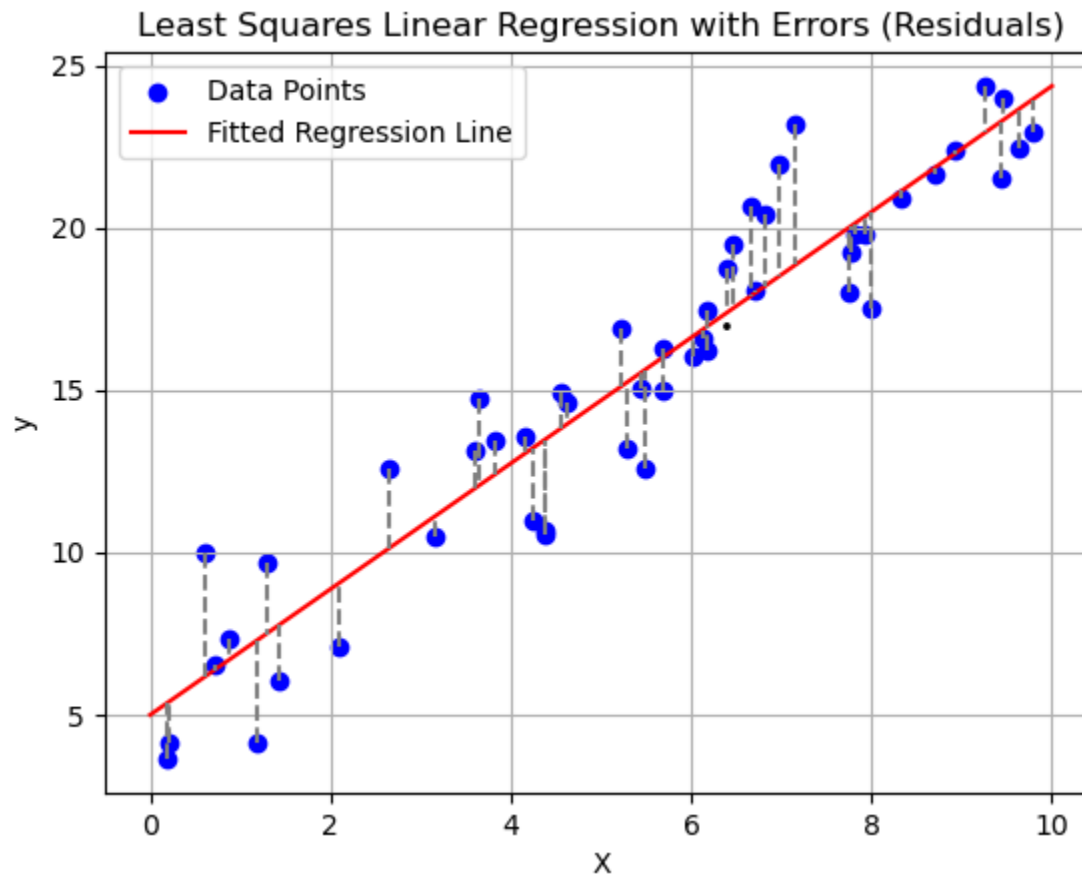
Find $\mathbf{y} \in \mathbb{R}^n$ such that

$$A\mathbf{y} - \mathbf{b} \approx \mathbf{0}$$

- **Fact:**

Different interpretations of “ \approx ” lead to different models!

Linear regression



<https://medium.com/physics-and-machine-learning/misconceptions-about-least-square-regression-1131841d240f>

Least squares linear regression

- Linear regression model

Dataset: $S = \{(\mathbf{x}^i, y^i) : i = 1, \dots, N\}$ for $\mathbf{x}^i \in \mathbb{R}^n$ and $y^i \in \mathbb{R}$.

Assumption:

- there is a hidden linear relationship between the input variables and output variables over S

Task: From the N observations (samples), find a linear hyperplane $y = \mathbf{m}^T \mathbf{x} + b$ ($\mathbf{m} \in \mathbb{R}^n, b \in \mathbb{R}$) in the $(n + 1)$ -dimensional space of (\mathbf{x}, y) , such that

$$y^i \approx \hat{y}^i = \mathbf{m}^T \mathbf{x}^i + b, i = 1, \dots, N$$

Least squares linear regression

- Hidden linear relationship

$$y^i \approx \hat{y}^i = \mathbf{m}^T \mathbf{x}^i + b, i = 1, \dots, N$$

- Least squares of errors:

1. define errors $e_i = |y^i - \hat{y}^i|$ for $i = 1, \dots, N$

2. minimize the squares of errors $\sum_{i=1}^N (y^i - \hat{y}^i)^2$

- Least squares optimization model

$$\text{Minimize } \sum_{i=1}^N (\mathbf{m}^T \mathbf{x}^i + b - y^i)^2$$

$$\text{s.t. } \mathbf{m} \in \mathbb{R}^n, b \in \mathbb{R}$$

Model reformulation

- Optimization model :

$$\begin{aligned} & \text{Minimize } \sum_{i=1}^N (\mathbf{m}^T \mathbf{x}^i + b - y^i)^2 \\ & \text{s.t. } \quad \mathbf{m} \in \mathbb{R}^n, b \in \mathbb{R} \end{aligned}$$

- Arrange data:

$$A = \begin{pmatrix} (\mathbf{x}^1)^T & 1 \\ \vdots & \vdots \\ (\mathbf{x}^N)^T & 1 \end{pmatrix} \in M_{N \times (n+1)}, \mathbf{z} = \begin{pmatrix} \mathbf{m} \\ b \end{pmatrix} \in \mathbb{R}^{n+1}, \mathbf{y} = \begin{pmatrix} y^1 \\ \vdots \\ y^N \end{pmatrix} \in \mathbb{R}^N$$

- Reformulated optimization model :

$$\begin{aligned} & \text{Minimize } \|\mathbf{Az} - \mathbf{y}\|_2^2 \\ & \text{s.t. } \quad \mathbf{z} \in \mathbb{R}^{n+1} \end{aligned}$$

A is given by input data, \mathbf{y} is given by output data, and \mathbf{z}^* determines the hidden linear relationship

Least squares linear regression

- Optimization model:

$$\begin{aligned} & \text{Minimize } \|A\mathbf{z} - \mathbf{y}\|_2^2 \\ & \text{s.t. } \mathbf{z} \in \mathbb{R}^{n+1} \end{aligned}$$

- Solution

$$\begin{aligned} \nabla \|A\mathbf{z} - \mathbf{y}\|_2^2 = 0 & \Leftrightarrow 2(A^T A)\mathbf{z} - 2A^T \mathbf{y} = 0 \\ & \Leftrightarrow (A^T A)\mathbf{z} = A^T \mathbf{y} \quad (\text{Normal Equation}) \end{aligned}$$

\mathbf{z}^* is a solution to a linear system of $n + 1$ variables, which is independent of the number of data points N !

Data manipulation – average effect

Data manipulation: Divide by N on both sides for average effect:

$$\frac{1}{N}A^T A = \frac{1}{N} \begin{pmatrix} x^1 & \dots & x^N \\ 1 & \dots & 1 \end{pmatrix} \begin{pmatrix} \frac{(x^1)^T}{1} \\ \vdots \\ \frac{(x^N)^T}{1} \\ \frac{1}{1} \end{pmatrix} = \frac{1}{N} \begin{pmatrix} \sum_{i=1}^N x^i (x^i)^T & \sum_{i=1}^N x^i \\ \sum_{i=1}^N (x^i)^T & N \end{pmatrix}$$

$$\frac{1}{N}A^T y = \frac{1}{N} \begin{pmatrix} x^1 & \dots & x^N \\ 1 & \dots & 1 \end{pmatrix} y = \frac{1}{N} \begin{pmatrix} \sum_{i=1}^N y^i x_1^i \\ \vdots \\ \sum_{i=1}^N y^i x_n^i \\ \sum_{i=1}^N y^i \end{pmatrix}$$

1. The **linear relationship** $\binom{m}{b}$ of the input vector x and output variable y **is determined by the average behavior of x^i , $x^i (x^i)^T$, y^i and $y^i x^i$.**
2. The **dimensionality of the underlying problem** depends solely on the number of features/attributes of the input and output variables ($n + 1$). **It is independent of the sample size (N).**

Discussion and extensions

- Other models for linear regression?
 - Least squares optimization model is commonly used
 - Support vector regression (SVR) model to be introduced
- Nonlinear regression?
 - increasing dimensionality for nonlinearity

Example: $y = ax^2 + bx + c$

quadratic relation between input $x \in \mathbb{R}$ and output y

$$y = (a, b, c) \begin{pmatrix} u_1 \\ u_2 \\ 1 \end{pmatrix} \text{ where } u_1 = x^2, u_2 = x$$

linear relation between input $u \in \mathbb{R}^2$ and output y !

Data regression - Support vector regression

- <https://medium.com/analytics-vidhya/support-vector-regression-svr-model-a-regression-based-machine-learning-approach-f4641670c5bb>

