

10-301/601: Introduction to Machine Learning

Lecture 7 – Linear Regression

Henry Chai

5/30/23

Lecture 7 Polls

0 done

 **0 underway**

Which of the following best describes your use of SG2?

Did not review SG2 at all

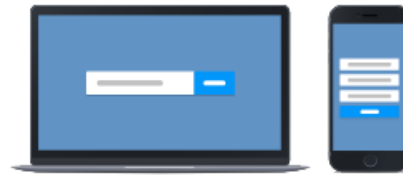
Read some or all of SG2 but did not attempt any problems

Completed or attempted some of the problems

Completed or attempted all of the problems

To the nearest integer, how many hours did you spend reviewing SG2? Please respond using digits [0-9], e.g., "12" instead of "twelve".

Join by Web



- 1 Go to **PolleV.com**
- 2 Enter **301601POLLS**
- 3 Respond to activity

i Instructions not active. **Log in** to activate

Front Matter

- Announcements:
 - PA2 released 5/25, due 6/01 at 11:59 PM
- Recommended Readings:
 - Murphy, Chapters 7.1-7.3

Recall: Regression

- Learning to diagnose heart disease

as a **(supervised)**

regression task

features

targets

data points

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000

Decision Tree Regression

Stump

- Learning to diagnose heart disease

as a (supervised)

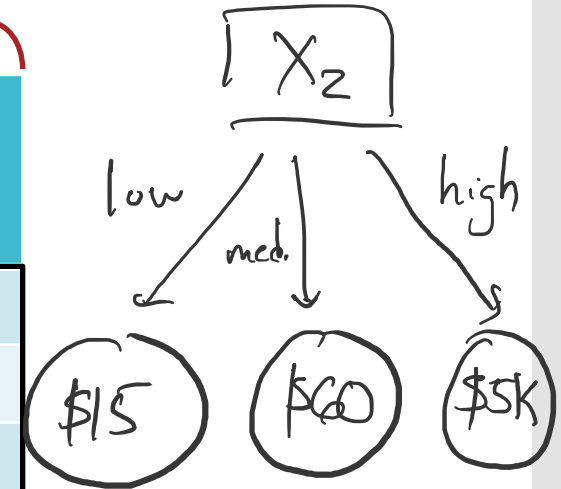
regression task

features

targets

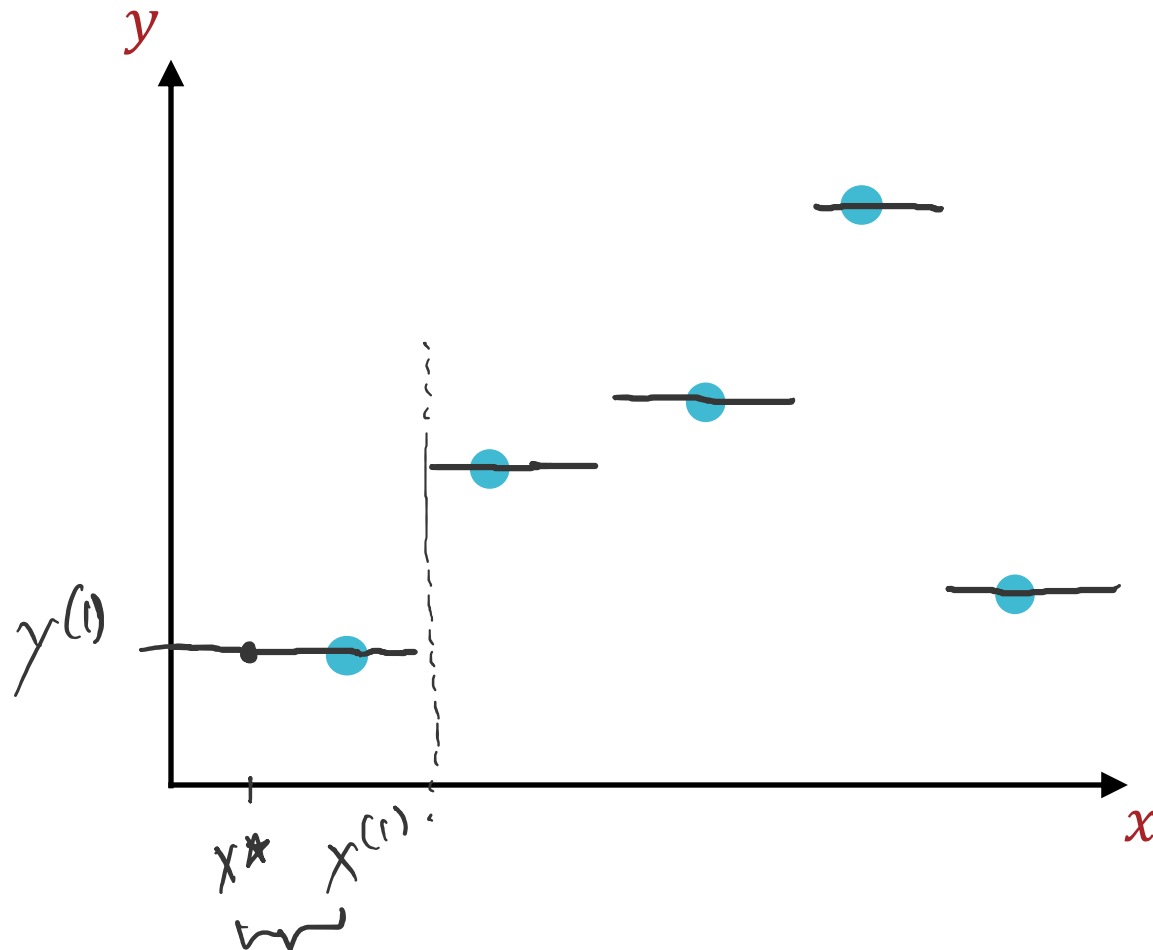
data points

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?
Yes	Low	Normal	\$0
No	Medium	Normal	\$20
No	Low	Abnormal	\$30
Yes	Medium	Normal	\$100
Yes	High	Abnormal	\$5000



1-NN Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and one-dimensional inputs $x \in \mathbb{R}$

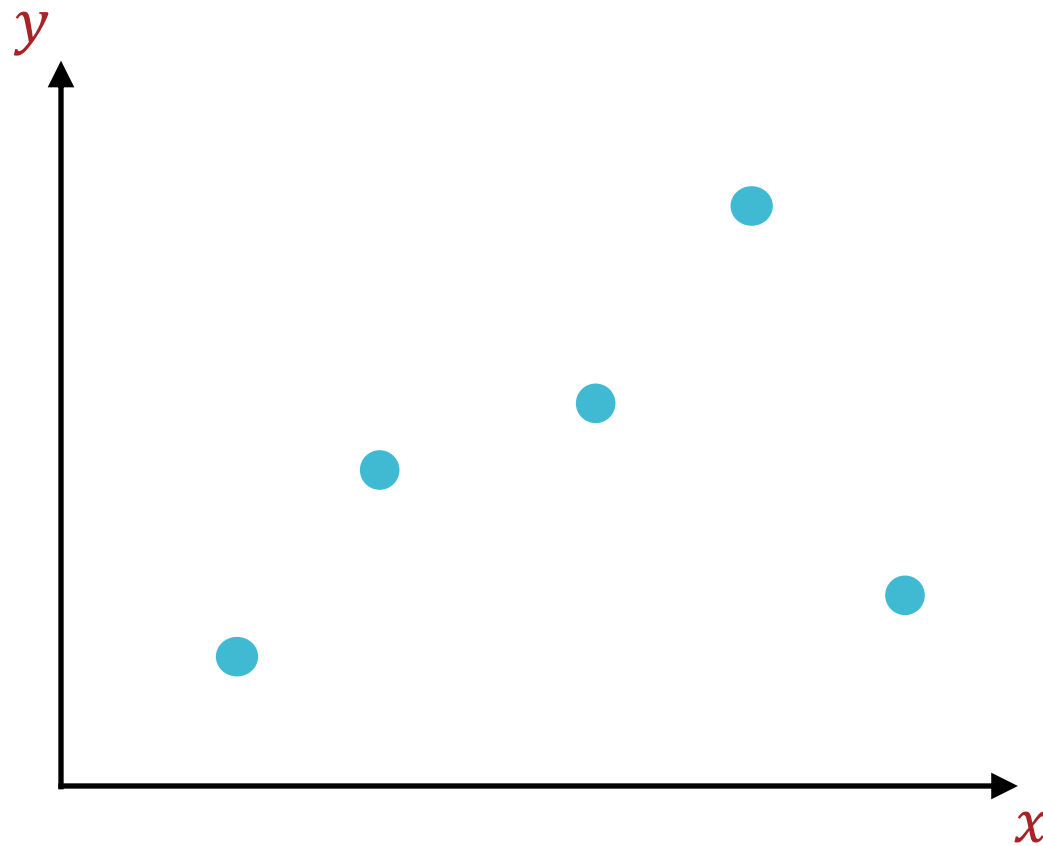


Distance-weighted?

2-NN

Regression?

- Suppose we have real-valued targets $y \in \mathbb{R}$ and one-dimensional inputs $x \in \mathbb{R}$



Linear Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and D -dimensional inputs $\mathbf{x} = [x_1, \dots, x_D]^T \in \mathbb{R}^D$

- **Assume**

$$y = \mathbf{w}^T \mathbf{x} + w_0$$

Linear Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and D -dimensional inputs $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$
- **Assume**

$$y = \mathbf{w}^T \mathbf{x}$$

General Recipe for Machine Learning

- Define a model and model parameters
- Write down an objective function
- Optimize the objective w.r.t. the model parameters

Recipe for Linear Regression

- Define a model and model parameters

Assume $y = w^T x$

Parameters: $w = [w_0, w_1, \dots, w_D]^T$

- Write down an objective function

Minimize the squared error

$$l_D(w) = \sum_{n=1}^N l^{(n)}(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2$$

- Optimize the objective w.r.t. the model parameters

Solve in "closed-form" by taking partial derivatives and finding "critical points"

Minimizing the Squared Error

$$l_D(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2$$

$$= \sum_{n=1}^N \left(\left(\sum_{d=0}^D w_d x_d^{(n)} \right) - y^{(n)} \right)^2$$

$$\frac{\partial l_D}{\partial w_d} = \sum_{n=1}^N \frac{\partial}{\partial w_d} \left(\left(\sum_{d=0}^D w_d x_d^{(n)} \right) - y^{(n)} \right)^2$$

$$= \sum_{n=1}^N 2 \left(\left(\sum_{d=0}^D w_d x_d^{(n)} \right) - y^{(n)} \right) \frac{\partial}{\partial w_d} \left(\left(\sum_{d=0}^D w_d x_d^{(n)} \right) - y^{(n)} \right)$$

$$= \sum_{n=1}^N 2 \left(\left(\sum_{d=0}^D w_d x_d^{(n)} \right) - y^{(n)} \right) x_d^{(n)}$$

↳ set equal to 0 and solve!

Recipe for Linear Regression

- Define a model and model parameters
 - Assume $y = \mathbf{w}^T \mathbf{x}$
 - Parameters: $\mathbf{w} = [w_0, w_1, \dots, w_D]$
- Write down an objective function
 - Minimize the squared error
$$\ell_{\mathcal{D}}(\mathbf{w}) = \sum_{n=1}^N \ell^{(n)}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y^{(n)})^2$$
- Optimize the objective w.r.t. the model parameters
 - Solve in *closed form*: take ~~partial derivatives~~ gradient, set to 0 and solve

Linear Regression

- Suppose we have real-valued targets $y \in \mathbb{R}$ and D -dimensional inputs $\mathbf{x} = [1, x_1, \dots, x_D]^T \in \mathbb{R}^{D+1}$

- **Assume**

$$y = \mathbf{w}^T \mathbf{x}$$

- Notation: given training data $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$

$$\bullet X = \begin{bmatrix} 1 & \mathbf{x}^{(1)T} \\ 1 & \mathbf{x}^{(2)T} \\ \vdots & \vdots \\ 1 & \mathbf{x}^{(N)T} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_D^{(1)} \\ 1 & x_1^{(2)} & \cdots & x_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(N)} & \cdots & x_D^{(N)} \end{bmatrix} \in \mathbb{R}^{N \times D+1}$$

is the *design matrix*

- $\mathbf{y} = [y^{(1)}, \dots, y^{(N)}]^T \in \mathbb{R}^N$ is the *target vector*

Minimizing the Squared Error

$$l_D(w) = \sum_{n=1}^N (w^T x^{(n)} - y^{(n)})^2 = \sum_{n=1}^N \underbrace{(x^{(n)T} w - y^{(n)})^2}$$

$$= (Xw - y)^T (Xw - y) = \|Xw - y\|_2^2$$

$$= (w^T X^T - y^T)(Xw - y)$$

$$= w^T X^T X w - \underbrace{y^T X w - w^T X^T y}_{\text{like } (Xw)^2} + y^T y \quad \text{like } (y)^2$$

kind of
like
 $(Xw)^2$

$$\rightarrow = \overbrace{w^T X^T X w} - 2w^T X^T y + \overbrace{y^T y}$$

$$\nabla_w l_D(w) = \begin{bmatrix} \frac{\partial l_D}{\partial w_0} \\ \vdots \\ \frac{\partial l_D}{\partial w_p} \end{bmatrix} = 2X^T X w - 2X^T y + 0$$

Minimizing the Squared Error

$$\begin{aligned} \ell_{\mathcal{D}}(\mathbf{w}) &= \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (\mathbf{x}^{(n)T} \mathbf{w} - y^{(n)})^2 \\ &= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \text{ where } \|\mathbf{z}\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{\mathbf{z}^T \mathbf{z}} \\ &= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \\ &= (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) \\ \nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\hat{\mathbf{w}}) &= (2\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} - 2\mathbf{X}^T \mathbf{y}) = \vec{0} \\ \rightarrow \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} &= \mathbf{X}^T \mathbf{y} \\ \rightarrow \hat{\mathbf{w}} &= \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}} \end{aligned}$$

Minimizing the Squared Error

$$\ell_{\mathcal{D}}(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^{(n)} - y^{(n)})^2 = \sum_{n=1}^N (\mathbf{x}^{(n)T} \mathbf{w} - y^{(n)})^2$$

$$= \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \text{ where } \|\mathbf{z}\|_2 = \sqrt{\sum_{d=1}^D z_d^2} = \sqrt{\mathbf{z}^T \mathbf{z}}$$

$$= (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

$$= (\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y})$$

$$\nabla_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) = (2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y})$$

$$H_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w}) = 2\mathbf{X}^T \mathbf{X}$$

$H_{\mathbf{w}} \ell_{\mathcal{D}}(\mathbf{w})$ is positive semi-definite

\Rightarrow at a minimum

Beyond the scope of the class

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

1. Is $\mathbf{X}^T \mathbf{X}$ invertible?

2. If so, how computationally expensive is inverting $\mathbf{X}^T \mathbf{X}$?

Closed Form Solution

Is $X^T X$ always invertible?

Yes

No

Unsure

If $X^T X$ is invertible, how computationally expensive is it to invert?

$O(N^2)$

$O(D^2)$

$O(ND)$

$O(N^3)$

$O(D^3)$

Closed Form Solution

$$X \in \mathbb{R}^{N \times (D+1)}$$

$$\hat{w} = (X^T X)^{-1} X^T y$$

1. Is $X^T X$ invertible?
 - When $N \gg D + 1$, $X^T X$ is (almost always) full rank and therefore, invertible!
 - If $X^T X$ is not invertible (occurs when one of the features is a linear combination of the others), what does that imply about our problem?
2. If so, how computationally expensive is inverting $X^T X$?
 - $X^T X \in \mathbb{R}^{D+1 \times D+1}$ so inverting $X^T X$ takes $O(D^3)$ time... *can get*
 - Computing $X^T X$ takes $O(ND^2)$ time *$O(D^{2.781})$*
 - What alternative optimization method(s) can we use to minimize the mean squared error?

Key Takeaways

- Decision tree and k NN regression
- Closed form solution for linear regression
 - Setting partial derivative/gradients to 0 and solving for critical points
 - Potential issues with the closed form solution: invertibility and computational c.sts