

MATH 829: Introduction to Data Mining and Analysis
Linear Regression: old and new

Dominique Guillot

Departments of Mathematical Sciences
University of Delaware

February 10, 2016

3/35

Linear regression: classical setting

$p = \text{nb. of variables}$, $n = \text{nb. of observations}$.

Classical setting:

- $n \gg p$ (n much larger than p). With enough observations, we hope to be able to build a good model.
- Note: even if the "true" relationship between the variables is not linear, we can include **transformations** of variables.
- E.g.

$$X_{p+1} = X_1^2, X_{p+2} = X_2^2, \dots$$

- Note: adding transformed variables can increase p significantly.
- A complex model requires a lot of observations.

Modern setting:

- In modern problems, it is often the case that $n \ll p$.
- Requires supplementary assumptions (e.g. sparsity).
- Can still build good models with very few observations.

3/35

Linear Regression: old and new

- Typical problem: we are given n observations of p variables X_1, \dots, X_p and Y .
- Goal: Use X_1, \dots, X_p to try to predict Y .
- Example: Cars data compiled using Kelley Blue Book ($n = 805$, $p = 11$).

Price	Mileage	Make	Model	Year	Type	Cylinder	Linear	Others	Crashes	Sound	Leather
17244.200	8623	BMW	Carbary	Sedan	Sedan	6	3.1	4	3	1	0
17642.000	9100	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
18238.800	10200	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
18306.913	10242	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
20286.17	10602	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
19709.000	12700	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
15230	22576	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
18348.542	22664	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
14862.004	24023	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
15209.010	27820	BMW	Carbary	Sedan	4D Sedan	6	3.1	4	3	1	0
20305.052	10207	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	0
20508.000	10600	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	0
20512.004	10600	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	0
18924.020	10600	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	1
18774.040	20500	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	1
18344.000	20700	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	0
18900.000	20700	BMW	Licenses	CX Sedan	Sedan	6	3.6	4	3	1	0

- Find a linear model $Y = \beta_1 X_1 + \dots + \beta_p X_p$.
- In the example, we want:
price = $\beta_1 \cdot \text{mileage} + \beta_2 \cdot \text{cylinder} + \dots$

2/35

Classical setting

Idea:

$$Y \in \mathbb{R}^{n \times 1} \quad X \in \mathbb{R}^{n \times p}$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} | & | & \dots & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \dots & \mathbf{x}_p \\ | & | & \dots & | \end{pmatrix},$$

where $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^{n \times 1}$ are the observations of X_1, \dots, X_p .

- We want $Y = \beta_1 X_1 + \dots + \beta_p X_p$.
- Equivalent to solving

$$Y = X\beta \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}.$$

4/35

Classical setting (cont.)

We need to solve $Y = X\beta$.

- Obviously, in general, the system has no solution.
- A popular approach is to solve the system in the least squares sense:

$$\hat{\beta} = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|Y - X\beta\|^2.$$

- How do we compute the solution?

Calculus approach:

$$\begin{aligned} \frac{\partial}{\partial \beta_i} \|Y - X\beta\|^2 &= \frac{\partial}{\partial \beta_i} \sum_{k=1}^n (y_k - X_{k1}\beta_1 - X_{k2}\beta_2 - \dots - X_{kp}\beta_p)^2 \\ &= 2 \sum_{k=1}^n (y_k - X_{k1}\beta_1 - X_{k2}\beta_2 - \dots - X_{kp}\beta_p) \times (-X_{ki}) \\ &= 0. \end{aligned}$$

Therefore,

$$\sum_{k=1}^n X_{ki}(X_{k1}\beta_1 + X_{k2}\beta_2 + \dots + X_{kp}\beta_p) = \sum_{k=1}^n X_{ki}y_k$$

1/35

Calculus approach (cont.)

Now

$$\sum_{k=1}^n X_{ki}(X_{k1}\beta_1 + X_{k2}\beta_2 + \dots + X_{kp}\beta_p) = \sum_{k=1}^n X_{ki}y_k \quad i = 1, \dots, p,$$

is equivalent to:

$$X^T X \beta = X^T y \quad (\text{Normal equations}).$$

We compute the Hessian:

$$\frac{\partial^2}{\partial \beta_i \partial \beta_j} \|Y - X\beta\|^2 = 2X^T X.$$

If $X^T X$ is invertible, then $X^T X$ is positive definite and

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

is the unique minimum of $\|Y - X\beta\|^2$.

6/35

Linear algebra approach

Want to solve $Y = X\beta$.

Linear algebra approach: Recall: If $V \subset \mathbb{R}^n$ is a subspace and $w \notin V$, then the best approximation of w by a vector in V is

$$\operatorname{proj}_V(w).$$

"Best" in the sense that:

$$\|w - \operatorname{proj}_V(w)\| \leq \|w - v\| \quad \forall v \in V.$$



Here:

$$X\beta \in \operatorname{col}(X) = \operatorname{span}\{\mathbf{x}_1, \dots, \mathbf{x}_p\}.$$

If $Y \notin \operatorname{col}(X)$, then the best approximation of Y by a vector in $\operatorname{col}(X)$ is

$$\operatorname{proj}_{\operatorname{col}(X)}(Y).$$

7/35

Linear algebra approach (cont.)

So $\|Y - \operatorname{proj}_{\operatorname{col}(X)}(Y)\| \leq \|Y - X\beta\| \quad \forall \beta \in \mathbb{R}^p.$

Therefore, to find $\hat{\beta}$, we solve

$$X\hat{\beta} = \operatorname{proj}_{\operatorname{col}(X)}(Y)$$

(Note: this system always has a solution.)

With a little more work, we can find an explicit solution:

$$Y - X\hat{\beta} = Y - \operatorname{proj}_{\operatorname{col}(X)}(Y) = \operatorname{proj}_{\operatorname{col}(X)^\perp}(Y).$$

Recall

$$\operatorname{col}(X)^\perp = \operatorname{null}(X^T).$$

Thus,

$$Y - X\hat{\beta} = \operatorname{proj}_{\operatorname{null}(X^T)}(Y) \in \operatorname{null}(X^T).$$

That implies:

$$X^T(Y - X\hat{\beta}) = 0.$$

Equivalently,

$$X^T X \hat{\beta} = X^T Y \quad (\text{Normal equations}).$$

8/35

Theorem (Least squares theorem)

Let $A \in \mathbb{R}^{n \times m}$ and $b \in \mathbb{R}^n$. Then

- $Ax = b$ always has a least squares solution \hat{x} .
- A vector \hat{x} is a least squares solution iff it satisfies the normal equations

$$A^T A \hat{x} = A^T b.$$

- \hat{x} is unique \Leftrightarrow the columns of A are linearly independent $\Leftrightarrow A^T A$ is invertible. In that case, the unique least squares solution is given by

$$\hat{x} = (A^T A)^{-1} A^T b.$$

9/15

The file JSE_Car_Lab.csv:

```
1 Price,MPG,Make,Model,Trim,Type,Cylinder,Liter,Doors,Cruise,Smart,Leather
2 17181.00,28,Mercedes,2012,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
3 17542.00,30,Mercedes,2013,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
4 16126.00,30,Mercedes,2011,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
5 16136.00,32,Infiniti,2008,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
6 16126.00,30,Mercedes,2011,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
7 15795.00,32,Mercedes,2013,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
8 15126.00,30,Mercedes,2011,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
9 14986.00,32,Infiniti,2011,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
10 14986.00,30,Mercedes,2011,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,0
11 15275.00,30,Mercedes,2010,Basic,Century,Sedan,40,Sedan,4,1,1,4,1,1,1
```

Loading the data with the headers using Pandas:

```
import pandas as pd
data = pd.read_csv('./data/JSE_Car_Lab.csv', delimiter=',')
```

We extract the numerical columns:

```
y = np.array(data['Price'])
x = np.array(data['MPG'])
x = x.reshape(len(x), 1)
```

10/15

Building a simple linear model with Python (cont.)

The `scikit-learn` package provides a lot of very powerful functions/objects to analyse datasets.

Typical syntax:

- Create object representing the model.
- Call the `fit` method of the model with the data as arguments.
- Use the `predict` method to make predictions.

```
from sklearn.linear_model import LinearRegression
lin_model = LinearRegression(fit_intercept=True)
lin_model.fit(x,y)
```

```
print lin_model.coef_
print lin_model.intercept_
```

We obtain $\text{price} \approx -0.17 \cdot \text{mileage} + 24764.5$.

11/15

Measuring the fit of a linear model

How good is our linear model?

- We examine the *residual sum of squares*:

$$RSS(\hat{\beta}) = \|y - X\hat{\beta}\|^2 = \sum_{k=1}^n (y_k - \hat{y}_k)^2.$$

```
((y-lin_model.predict(x))**2).sum()
```

We find: 76855792485.91. Quite a large error... The average absolute error:

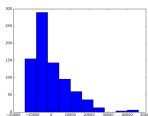
```
(abs(y-lin_model.predict(x))).mean()
is 7596.28. Not so good...
```

- We examine the distribution of the residuals:

```
import matplotlib.pyplot as plt
plt.hist(y-lin_model.predict(x))
plt.show()
```

12/15

Histogram of the residuals:

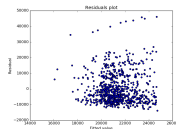


- Non-symmetric.
- Heavy tail.

- The heavy tail suggests there may be outliers.
- It also suggests transforming the response variable using a transformation such as \log , $\sqrt{\cdot}$, or $1/x$.

13/35

Plotting the residuals as a function of the fitted values, we immediately observe some patterns.



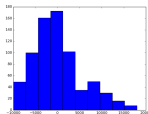
Outliers? Separate categories of cars?

14/35

Improving the model

- Add more variables to the model.
- Select the best variables to include.
- Use transformations.
- Separate cars into categories (e.g. exclude expensive cars).
- etc.

For example, let us use all the variables, and exclude Cadillacs from the dataset.



- Much more symmetric.
- Closer to a Gaussian distribution.

Average absolute error drops to 4241.21.

15/35