

Introduction to Data Science with Python

Week 0, Section 1: Linear Regression

Book: An introduction to Statistical Learning, 2nd

Our main goal in data science is to predict outcomes and analyze data.

When dealing with this kinds of problems there is a common **asymmetry**; the variable we want to predict may be difficult, more important than others, or might be influenced by others. Which leaves us with 2 types of variables:

- 1 values we **want** to predict. (y)
- 2 values we **use** to make a prediction. (X)

Response vs. Predictor Variables			y	
n observations				
	TV	Radio		
	230.1	37.8	69.2	
	44.5	39.3	45.1	
	17.2	45.9	69.3	
	151.5	41.3	58.5	
	180.8	10.8	58.4	

$\underbrace{\hspace{10em}}$
 p predictors

Capital letters = matrix

lower case = vectors

Pandas Correlation:

- $X.\text{shape}$ = returns dimensions in (n, p) format
- $Y.\text{shape}$ = returns $(n,)$ or $(n, 1)$
- $\text{df}[[\cdot]]$ vs $\text{df}[\cdot]$ = double returns a pd.DataFrame, and single returns pd.Series.

Statistical Model:

Assume that variable Y relates to predictor/s X , through an unknown $f(X)$ function and a random amount ε , that describes Y from the rule $f(x)$

$$Y = f(X) + \varepsilon$$

When we use a set of measurements (x_{i1}, \dots, x_{ip}) to predict a value we use

$$\hat{y} = \hat{f}(x_{i1}, \dots, x_{ip})$$

when we care about F 's form. We call it **inference problems**. for some problems we don't care about the specific F form, these are prediction problems

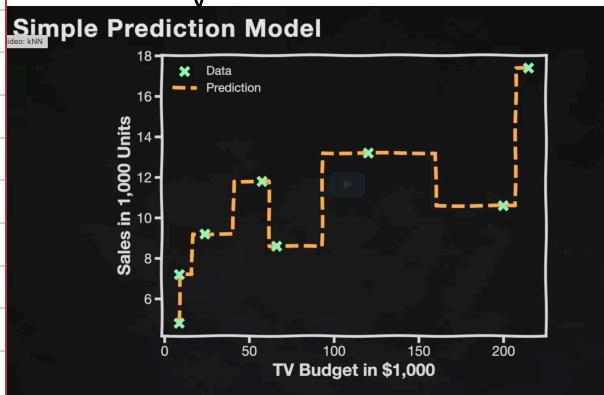
We can take the mean of all y_i 's with $\frac{1}{n} \sum_{i=1}^n y_i$:

In our specific dataset we can find the distances to all other points $D(x_q, x_i)$. We are looking to answer the sales for TV budget of 175,000.

We find the nearest neighbor (x_p, y_p) and predict

$$\hat{y}_q = y_p$$

Prediction for any budget:



1.2 Error evaluation and model comparison

first we're gonna learn how to determine how to choose models $k=1, k=10, k=70$.

To determine this we need to calculate the errors each model does.

Error Evaluation:

- first, we randomly hide some data from the model. This is called train-validation, that is done to check how the model does with unseen data.
- we use the train set to estimate \hat{y} , and the validation set to evaluate the model.
- we bring back the validation data, an error is the difference between the prediction and the true value. $r_i = y_i - \hat{y}_i$
- then to quantify how well the model performs,

we aggregate the errors, and that's called
loss, error, or cost function

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Model Comparison

- After getting all MSE 's, we choose the one closer to 0.

Model fitness

- To know if your model is good enough, we compare to a really good and bad models. In this case we have $\hat{y} = \bar{y}$ $\forall i \sum_i y_i$ for the worst, and $\hat{y}_i = y_i$ as the best one. Were using $k=3$, $MSE = 5$.
- then we create the entity R-squared

$$R\text{-squared} = \frac{\sum_i (\hat{y}_i - \bar{y})^2}{\sum_i (\bar{y} - y_i)^2}$$

if the model gets a 1, the model is perfect, and a 0

if it is as good as the mean value \bar{y} , if it is a negative value then it's worse than average.

1.3 Linear Regression

What happens when now we ask "How will sales go if we double the TV budget?"

Linear Models

To build one, we must a simple form of f

$$f(x) = \beta_0 + \beta_1 x$$

And then we add the estimate

$$\hat{y} = \hat{f}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$$

$\hat{\beta}_1$ and $\hat{\beta}_0$ are estimates of β_1 and β_0

How to choose a line

We can use MSE as our loss function

$$L(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^n [y_i - (\beta_1 x_i + \beta_0)]^2$$

The optimal values of the betas should be the minimum of the loss functions with respect to beta 0 and 1. To minimize loss.

Video: Linear Regression

Recap Exercise

```
>>> from sklearn.linear_model import LinearRegression
>>> df = pd.read_csv('Advertising.csv')
>>> X = df[['TV']].values
>>> y = df['Sales'].values
>>> reg = LinearRegression()
>>> reg.fit(X, y)
>>> reg.coef_
array([[0.04665056]])
>>> reg.intercept_
array([7.08543108])
>>> reg.predict(np.array([[100]]))
array([[11.75048733]])
```

HX

reg.fit (X, y)

chain rule = if a function depends on another function that depends on a variable.

Section 2: Multiple and Polynomial Regression

2.1 Multiple Regression

The goal now is for our model to take as much data as possible with multiple variables.

When you add one more predictor, we have a 3D model

$$Sales_n = \beta_0 + \beta_1 * TV_n + \beta_2 * Radio_n + \beta_3 * Newspaper_n$$

$$Sales_1 = \beta_0 + [TV_1 * Radio_1 * Newspaper_1] \cdot \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

Video: Multilinear Regression

Multilinear Model: Example

For our data:

$$\begin{aligned} Sales_1 &= \begin{bmatrix} 1 & TV_1 & Radio_1 & Newspaper_1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \\ Sales_2 &= \begin{bmatrix} 1 & TV_2 & Radio_2 & Newspaper_2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \\ Sales_n &= \begin{bmatrix} 1 & TV_n & Radio_n & Newspaper_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \end{aligned}$$

$Y = X\beta$

HX

We can transpose a Matrix, which is moving it from vertical to horizontal or viceversa, using Numpy's `np.transpose()`

Video: Multilinear Regression

Recap: Inverse of a Matrix

```
In [16]: x = np.array([[1,2],[3,4]])
...
.... #Inverse array x
.... invX = np.linalg.inv(x) ←
.... print(invX)
...
.... #Verifying
.... print(np.dot(x, invX))
[[ -2.   1. ]
 [ 1.5 -0.5]]
[[1.00000000e+00  1.11022302e-16]
 [0.00000000e+00  1.00000000e+00]]
```

`numpy.linalg.inv()` is used to calculate the inverse of a matrix (if it exists!)

HX

When doing multiple linear regression, we still use MSE as a loss function to keep $\hat{Y} = X\beta$, then we take the partial derivative and equal it to 0

$$X^T(Y - X\beta) = 0$$

$$X^T X \beta = X^T Y$$

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

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

Practice:

We used Advertising.csv to predict sales based on 3 predictors

Linear Models and R²:

we used fit_and_plot_linear, which returns graphs, then we used fit_and_plot_multi to use multiple variables, this returned R² values.

Predictions:

Learned how to use itertools.combinations to get all possible combinations of predictors and then append to a list, this uses a for loop in a range. Then to start evaluating losses we set x = df [:] and y = df ['Sales'], then using train-test-split to assign train size to 80%, assign linearRegression to a variable

Predict using `predict(x-test)` and you can then get the MSE using mean-square-error

2.2 Techniques for Multilinear Modeling

Qualitative Predictors

Video: Multilinear Techniques

Qualitative Predictors										
Income	Limit	Rating	Cards	Age	Education	Sex	Student	Married	Ethnicity	Balance
14.890	3606	283	2	34	11	Male	No	Yes	Caucasian	333
106.02	6645	483	3	82	15	Female	Yes	Yes	Asian	903
104.59	7075	514	4	71	11	Male	No	No	Asian	580
148.92	9504	681	3	36	11	Female	No	No	African American	964
55.882	4897	357	2	68	16	Male	No	Yes	Caucasian	331

If the predictor only takes two variables, we create a dummy variable. For sex we can equal 1 to male and 0 to female.

β_0 = males, $\beta_0 + \beta_1$ = females, β_1 = difference between males and females. We create dummy variables for ethnicity we need an extra one. The idea is to find correlations.

Scaling

In an issue where we need a parabola but our code only gives straight lines, we need to scale the predictor, that fixes this situation but why? Generally, it works well when the predictors are not in the right magnitude, but it can also affect because of how sklearn solves a regression, when inverting a matrix it can fail, for which scaling can be helpful.

Normalize = linear scaling into 0-1 range

Standardize = Gaussian scaling around origin

Practice

We learned about categories and how to set strings to booleans by using dummy's, then you include them in R^2 's

2.3 Polynomial Regression

Sometimes our models might not be linear. Basically, $y = f_B(x)$ where f is not a linear function and B is a parameter of f

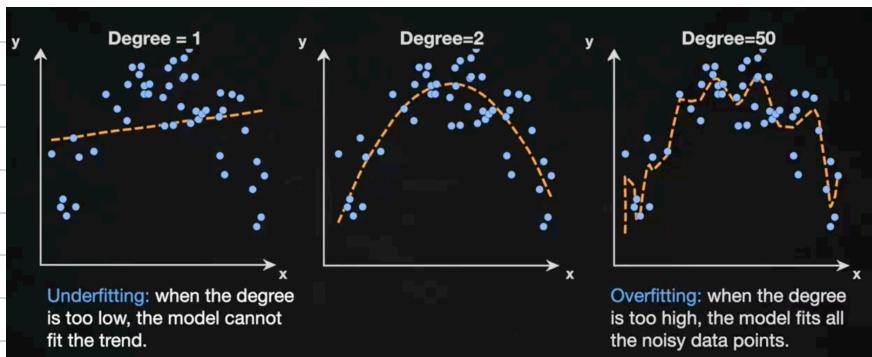
Polynomial Regression

$$\mathbf{X} = \begin{pmatrix} 1 & x_1^1 & \dots & x_1^M \\ 1 & x_2^1 & \dots & x_2^M \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^1 & \dots & x_n^M \end{pmatrix}$$

This looks a lot like multi-linear regression where the predictors are powers of x :

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,J} \\ 1 & x_{2,1} & \dots & x_{2,J} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,J} \end{pmatrix}$$

$\mathbf{H}\mathbf{X}$



Future Scaling

normally when the range of x is small or large

can cause some problems that can be solved by scaling

Feature Scaling

It is always a good idea to **scale** X when considering polynomial regression:

$$X^{stand} = \frac{X - \bar{X}}{\sigma_X} \quad X^{norm} = \frac{X - X_{min}}{X_{max} - X_{min}}$$

Hx

→ This are the ways to scale a matrix