

4.1 - Linear and Logistic Regression

Giovanni Della Lunga
giovanni.dellalunga@unibo.it

Introduction to Machine Learning for Finance

Bologna - February, 2022

Subsection 1

Back to Linear Regression

Linear Regression

- Linear regression is a very popular tool because once you have made the assumption that the model is linear you do not need huge amount of data. In ML we refer to the constant term as the *bias* and the coefficients as *weights*
- Assume n observations and m features. Model is

$$Y = a + b_1X_1 + b_2X_2 + \cdots + b_mX_m + \epsilon$$

- Standard approach is to choose a and the b_i to minimize the mean square error (mse):

$$mse = \frac{1}{n} \sum_{j=1}^n [Y_j - (a + b_1X_{1,j} + b_2X_{2,j} + \cdots + b_mX_{m,j})]^2 \quad (1)$$

- This can be done analytically by inverting a matrix, in practice a numerical (gradient descent) is used.

Categorical Features (see chapter 2.1)

- Remember that categorical features are features where there are a number of non-numerical alternatives
- We can define a dummy variable for each alternative. The variable equals 1 if the alternative is true and zero otherwise. This is known as **one-hot encoding**
- But sometimes we do not have to do this because there is a natural ordering of variables, e.g.:
 - small=1, medium=2, large=3
 - assist. prof=1, assoc. prof=2, full prof =3

Dummy Variably Trap

- Suppose we have a constant term and a number of dummy variables (equal to 0 or 1)
- There is then no unique solution because, for any C , we can add C to the constant term and subtract C from each of the dummy variables without changing the prediction
- A side effect of regularization is that it solves this problem

Iowa House Price Case Study



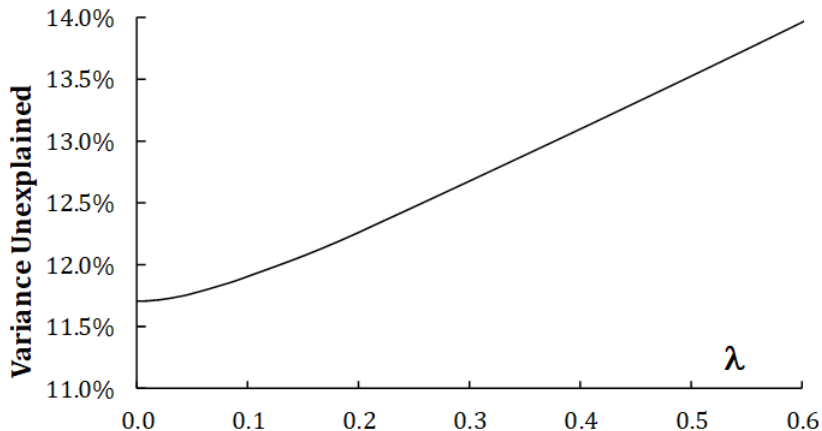
- The objective is to predict the prices of house in Iowa from features
- 800 observations in training set, 600 in validation set, and 508 in test set
- Here the original competition description:
<https://www.kaggle.com/c/house-prices-advanced-regression-techniques>

Iowa House Price Results (No regularization)

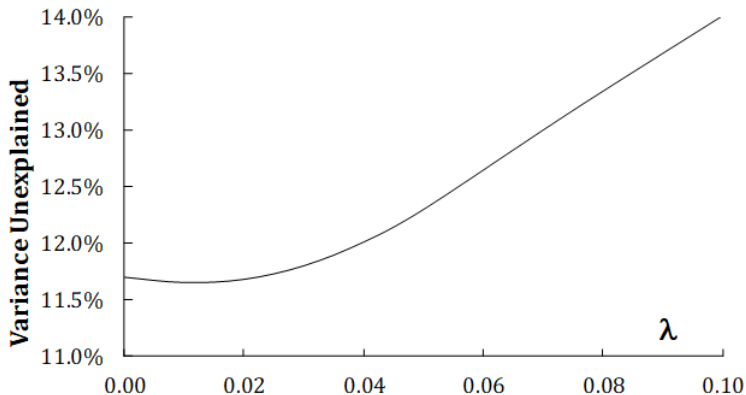
2 categorical variables included. Natural ordering for Basement quality. 25 dummy variables created for neighborhood

Lot area (squ ft)	0.08	Number of half bathrooms	0.02
Overall quality (scale from 1 to 10)	0.21	Number of bedrooms	-0.08
Overall condition (scale from 1 to 10)	0.10	Total rooms above grade	0.08
Year built	0.16	Number of fireplaces	0.03
Year remodeled	0.03	Parking spaces in garage	0.04
Basement finished squ ft	0.09	Garage area (squ ft)	0.05
Basement unfinished squ ft	-0.03	Wood deck (squ ft)	0.02
Total basement squ ft	0.14	Open porch (squ ft)	0.03
1st floor squ ft	0.15	Enclosed porch (squ ft)	0.01
2 nd floor squ ft	0.13	Neighborhood (25 alternatives)	-0.05 to 0.12
Living area	0.16	Basement quality (6 natural ordering)	0.01
Number of full bathrooms	-0.02		

Ridge Results for validation set



Lasso Results for validation set



Iowa House Price Results

Non-zero weights for Lasso when $\lambda = 0.1$ (overall quality and total living area were most important)

Feature	Weight
Lot Area (square feet)	0.04
Overall quality (Scale from 1 to 10)	0.30
Year built	0.05
Year remodeled	0.06
Finished basement (square feet)	0.12
Total basement (square feet)	0.10
First floor (square feet)	0.03
Living area (square feet)	0.30
Number of fireplaces	0.02
Parking spaces in garage	0.03
Garage area (square feet)	0.07
Neighborhoods (3 out of 25 non-zero)	0.01, 0.02, and 0.08
Basement quality	0.02

Summary of Iowa House Price Results

- With no regularization correlation between features leads to some negative weights which we would expect to be positive
- Improvements from Ridge is modest
- Lasso leads to a much bigger improvement in this case
- Elastic net similar to Lasso in this case
- Mean squared error for test set for Lasso with $\lambda = 0.1$ is 14.7

Subsection 2

Logistic Regression

Logistic Regression

- The objective is to classify observations into a **positive outcome** and **negative outcome** using data on features
- Probability of a positive outcome is assumed to be a sigmoid function:

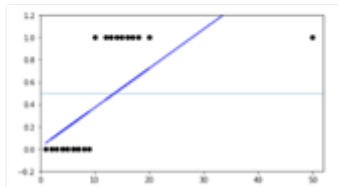
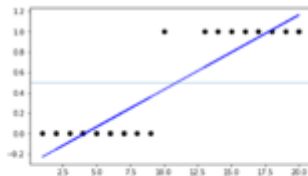
$$Q = \frac{1}{1 + e^{-Y}} \quad (2)$$

- where Y is related linearly to the values of the features:

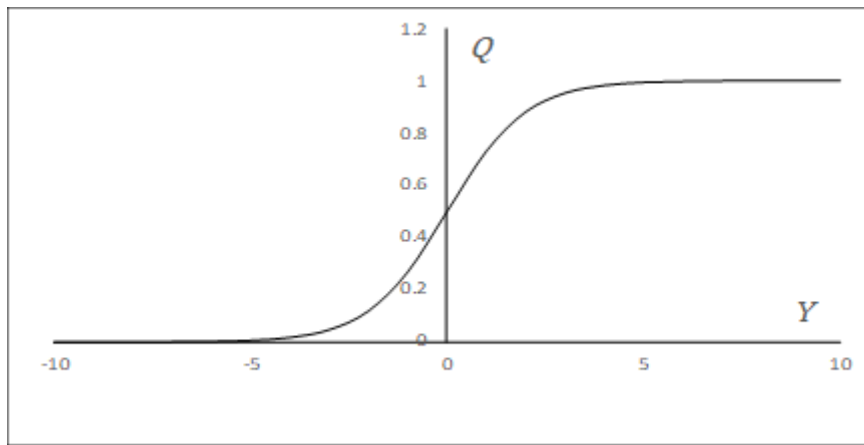
$$Y = a + b_1X_1 + b_2X_2 + \cdots + b_mX_m \quad (3)$$

- Can use regularization

Logistic Regression



The Sigmoid Function



Maximum Likelihood Estimation

- We use the training set to maximize
-

$$L = \sum_{\text{POS OUT}} \ln(Q) + \sum_{\text{NEG OUT}} \ln(1 - Q) \quad (4)$$

- This cannot be maximized analytically but we can use a gradient ascent algorithm

Confusion Matrix

	p' (Predicted)	n' (Predicted)
p (Actual)	True Positive	False Negative
n (Actual)	False Positive	True Negative

Lending Club Case Study

- Data consists of loans made and whether they proved to be good or defaulted. (A restriction is that you do not have data for loans that were never made.)
- We use only four features
- Home ownership (rent vs. own)
- Income
- Debt to income
- Credit score
- Training set has 8,695 observations (7,196 good loans and 1,499 defaulting loans). Test set has 5,196 observations (4,858 good loans and 1,058 defaulting loans)

The Data

Home Ownership 1=owns, 0 =rents	Income (\$'000)	Debt to Income (%)	Credit score	1=Good, 0=Default
1	44.304	18.47	690	0
1	136.000	20.63	670	1
0	38.500	33.73	660	0
1	88.000	5.32	660	1

Results for Lending Club Training Set

- X_1 = Home Ownership
- X_2 = Income
- X_3 = Debt to income ratio
- X_4 = Credit score
-

$$Y = -6.5645 + 0.1395 \cdot X_1 + 0.0041 \cdot X_2 - 0.0011 \cdot X_3 + 0.0113 \cdot X_4$$

Decision Criterion

- The data set is imbalanced with more good loans than defaulting loans
- There are procedures for creating a balanced data set
- With a balanced data set we could classify an observation as positive if $Q > 0.5$ and negative otherwise
- However this does not consider the cost of misclassifying a bad loan and the lost profit from misclassifying a good loan
- A better approach is to investigate different thresholds, Z
- If $Q > Z$ we accept a loan
- If $Q \leq Z$ we reject the loan

Test Results

See Hull, Tables 3.10, 3.11, and 3.12

$Z = 0.75$:

	Predict no default	Predict default
Outcome positive (no default)	77.59%	4.53%
Outcome negative (default)	16.26%	1.62%

$Z=0.80$:

	Predict no default	Predict default
Outcome positive (no default)	55.34%	26.77%
Outcome negative (default)	9.75%	8.13%

$Z=0.85$:

	Predict no default	Predict default
Outcome positive (no default)	28.65%	53.47%
Outcome negative (default)	3.74%	14.15%

The Confusion matrix and common ratios

	Predict positive outcome	Predict negative outcome
Outcome positive	TP	FN
Outcome negative	FP	TN

$$\text{Accuracy} = \frac{TP + TN}{TP + FN + FP + TN}$$

$$\text{True Positive Rate (TPR also called sensitivity or recall)} = \frac{TP}{TP + FN}$$

$$\text{The True Negative rate(also called specificity)} = \frac{TN}{TN + FP}$$

$$\text{The False Positive Rate} = \frac{FP}{TN + FP}$$

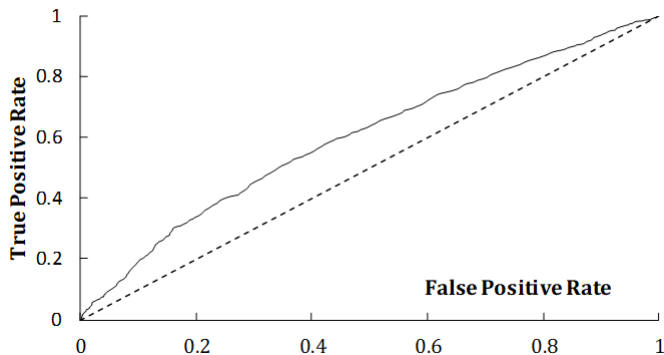
$$\text{Precision, P} = \frac{TP}{TP + FP}$$

$$\text{F score} = 2 \times \frac{P \times \text{TPR}}{P + \text{TPR}}$$

Test Set Ratios for different Z values

	$Z = 0.75$	$Z = 0.80$	$Z = 0.85$
Accuracy	79.21%	63.47%	42.80%
True Positive Rate	94.48%	67.39%	34.89%
True Negative Rate	9.07%	45.46%	79.11%
False Positive Rate	90.93%	54.54%	20.89%
Precision	82.67%	85.02%	88.47%
F-score	88.18%	75.19%	50.04%

As we change the Z criterion we get an ROC



Area Under Curve (AUC)

- The area under the curve is a popular way of summarizing the predictive ability of a model to estimate a binary variable
- When $AUC = 1$ the model is perfect.
- When $AUC = 0.5$ the model has no predictive ability
- When $AUC < 0.5$ the model is worse than random
- In this case $AUC = 0.6020$

Choosing Z

- The value of Z can be based on
- The expected profit from a loan that is good, P
- The expected loss from a loan that defaults, L
- We need to maximize: $P \times TP - L \times FP$