2695 Machine Learning Masters Program in Economics, Finance and Management





LINEAR REGRESSION

Different types of ML algorithms







Supervised learning: Classification vs Regression





Supervised learning: Classification vs Regression





Goal of Linear regression

Goal: Given input x we would like to compute output y

Data set

Living Area (Feet ²)	Price (\$)	
1180	221,900	
2570	538,000	
770	180,000	
1960	604,000	
1680	510,000	
5420	1,225,000	
1715	257,500	
1060	291,850	
1780	229,500	
1890	323,000	
3560	662,500	
1160	468,000	
1430	310,000	
1370	400,000	
1810	530,000	
x	у	

Example: Predict house price

x: living areay: price



Linear model



Goal: Given an input x we would like to compute an output y

Data set

Living Area (Feet ²)	Price (\$)
1180	221,900
2570	538,000
770	180,000
1960	604,000
1680	510,000
5420	1,225,000
1715	257,500
1060	291,850
1780	229,500
1890	323,000
3560	662,500
1160	468,000
1430	310,000
1370	400,000
1810	530,000
x	у

Example: Predict house's price



y: price



In linear regression we assume that y and x are related with the following equation:



- θ: parameter ٠
- ε: measurement or other noise ٠



iving area = 4876 feet2



The best possible line

Our goal is to estimate θ from a training data of $\langle x^i, y^i \rangle$ pairs in order to find the "best possible line".





What is the line with the *best fit*

- "Best" line depends on the objective (loss, cost) function: Objective function should reflect our goal.
- Cost function determines how much penalty should be assigned to an instance based on the error in the model's
 predicted value
- There exist many different cost functions, since there are many ways to compute the error between an estimated value and an actual value.



Mean squared error

- "Best" line depends on the objective (loss, cost) function: Objective function should reflect our goal.
- Cost function determines how much penalty should be assigned to an instance based on the error in the model's
 predicted value
- There exist many different cost functions, since there are many ways to compute the error between an estimated value and an actual value.

• $\frac{1}{m}\sum_i (y^i - \theta x^i)^2$

m: number of training data points



X (Size of House)



Objective functions

- "Best" line depends on the objective (loss, cost) function: Objective function should reflect our goal.
- Cost function determines how much penalty should be assigned to an instance based on the error in the model's
 predicted value
- There exist many different cost functions, since there are many ways to compute the error between an estimated value and an actual value.

m: number of training data points

- Why this objective function:
 - minimizes squared distance between measurements and predicted line: strongly penalizes very large errors
 - easy to compute



Finding the optimal value of the parameter

- Our goal is to estimate θ from a training data of $\langle x^i, y^i \rangle$ pairs to find the "best possible line" $y = \theta x + \varepsilon$
- We can find the optimal value for the parameters by minimizing our cost function:

$$\frac{1}{m}\sum_{i}(y^{i}- heta x^{i})^{2}$$

- Optimal value means:
 - For any other values of θ , the cost function would be higher for our **training data**.
- *Optimal value does not* mean:
 - that our model is generally a good model for the data
 - that our model will perform well on new instances.



Introducing the bias (intercept term)

So far, we assumed that the line passes through the origin.

- What if the line does not?
- No problem, simply change the model to:

$$y = \theta_0 + \theta_1 x + \varepsilon$$

• y intercept: value at which the fitted line crosses the y-axis





Training vs Inference

- **Training:** We learn the values of parameters θ_1 and θ_0 using pairs $\langle x^i, y^i \rangle$.
- Inference (Prediction): Using the model, we can predict the value of y, for a *new* x.





• What is the price of a house with this area?



From simple to multiple linear regression

- Generally, one dependent variable (y) depends on multiple factors (x₁, x₂,...).
- For example, the price of a house depends on many factors like the size of the house, number of rooms, attached facilities, distance of nearest station from it, distance of nearest shopping area from it...



The house's price depends on **multiple features** (variables)

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178



Multiple (multivariate) linear regression

This is quite like the simple linear regression model, but with multiple independent variables contributing to the dependent variable.

Each training sample has n attributes $(x_1, x_2, ..., x_n)$ and a corresponding single value of y.

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
example	<i>x</i> ₁	<i>x</i> ₂	<i>x</i> ₃	<i>x</i> ₄	У
	2104	5	1	45	460
	1416	3	2	40	232
	1534	3	2	30	315
	852	2	1	36	178



General linear regression



- What if your data is actually more complex than a simple straight line?
- Instead of the input variables (x) use some function of these values.
- Calculate functions of x, the regression model remains linear with respect to the coefficients.

Linear regression:General linear regression: $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$ $\hat{y} = \theta_0 + \theta_1 \phi_1(x) + \theta_2 \phi_2(x) + \dots + \theta_n \phi_n(x)$

 $\hat{y}=\Theta x$ $\hat{y}=\Theta x^2$



Polynomial regression

• Add powers of each feature as new features, then train a linear model on this extended set of features.

For lower degrees, the relationship has a specific name like quadratic.

- n = 1, a first degree polynomial is simply a straight line. $\rightarrow \hat{y} = \theta_0 + \theta_1 x$
- n = 2, a second degree polynomial is called a **quadratic.** $\rightarrow \hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2$
- n = 3, it is **cubic.** $\rightarrow \hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$
- n = 4, it's called **quartic.** $\rightarrow \hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

 $\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \dots + \theta_n x^n$

where *n* is called the **degree** of the polynomial

The polynomial models can be used to approximate a complex nonlinear relationship



17



Linear Regression properties

- Supervised algorithm
- Training data point: one or more input values and a single output value
- Training: learning the line with the best fit
- Prediction: using the learnt line to predict the output value
- Simple to interpret, easy to train
- Used as a baseline for evaluating other, more complex models



2695 Machine Learning Masters Program in Economics, Finance and Management





ESTIMATING PARAMETERS OF LINEAR REGRESSION



Equation for calculating the parameters

 $\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

- \hat{y} is the predicted value.
- *n* is the number of features.
- x_i is the ith feature value.
- θ_j is the jth model parameter (including the bias term θ_0 and the feature weights $\theta_1, \theta_2, \dots, \theta_n$).



ith data point with n features

Normal Equation (Closed Form)

Solve θ analytically: a mathematical equation that gives the result directly

$$\hat{\theta} = \left(\mathbf{X}^T \cdot \mathbf{X} \right)^{-1} \cdot \mathbf{X}^T \cdot \mathbf{y}$$

- $\hat{\theta}$ is the value of θ that minimizes the cost function.
- **y** is the vector of target values containing $y^{(1)}$ to $y^{(m)}$.

It's suitable for small feature set (e.g. < 1000 features)



Gradient Descent

One of the most popular optimization methods in Machine Learning (*)



(*) used in training of neural networks

How does Gradient descent work



 $\boldsymbol{J}(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^{m} (y_i - \hat{y}_i)^2$

Gradient descent algorithm

The parameter are iteratively updated in the following equation: Learning rate (Step size) $\theta_{new} = \theta_{old} - \alpha \cdot \nabla_{\theta} J(\theta)$ Gradient $\nabla_{\theta} J(\theta) = \frac{d}{d\theta} J(\theta)$

- 1. Pick a value for the learning rate $\boldsymbol{\alpha}$
- 2. Start by filling θ_{old} with random values (random initialization)
- 3. Calculate the gradient at the point θ_{old} . Update new parameter θ_{new} by following the opposite direction
- 4. Repeat until cost converges to the minimum



Learning rate

Learning rate is a hyperparameter



If the learning rate is too large



24

Gradient Descent Pitfalls

Gradient Descent for convex function & non-convex function





The gradient descent is guaranteed to converge arbitrarily close to the global minimum when cost function is convex

If the random initialization starts the algorithm on the left, then it will converge to a local minimum, which is not as good as the global minimum.



2695 Machine Learning Masters Program in Economics, Finance and Management





EVALUATION METRICS REGRESSION



Criteria for a *good* model



Accurate

Are we making correct predictions?

Interpretable

How easy is it to explain how the predictions are made?



Fast

How long does it take to build a model and how long does the model take to make predictions?



How much longer do we have to wait if we build/predict using a lot of data?



Mean Squared Error

MSE is the mean of the squared value of the errors



- MSE: mean of the squares of the difference between the actual value and the predicted values for all data points.
- A smaller score is better.



Root Mean Squared Error



- The most popular evaluation metrics used in regression problems. It has same unit with "Y".
- A smaller score is better.

Mean Absolute Error



- MAE: Mean of the absolute difference between the actual value and the predicted values for all data points.
- A smaller score is better.

OVA SCHOOL OF

BUSINESS & ECONOMICS



Mean Absolute Percent Error



- MAPE: measures the size of the error in percentage terms.
- One of the most popular measures for forecasting time series error.
- A smaller score is better.



R² - coefficient of determination

• R² Score is based on comparing our model to the simplest possible model (mean of y)



Adjusted R^2



R²_{adj}

Adjusted R² penalizes the model for the inclusion of variables:

- more useless variables to a model, adjusted R² decreases
- more useful variables, adjusted R² increases

Adjusted R^2 is always less than or equal to R^2

- 0: model has no predictive power
- 1: model that perfectly predicts
- In the real world, adjusted R² lies between these values

$$R_{adj}^2 = 1 - \left[\frac{(1-R^2)(n-1)}{n-k-1}\right]$$

n is the sample size (i.e., number of points in your data sample)
k is s the total number of explanatory variables in the model
(not including the constant term)

2695 Machine Learning Masters Program in Economics, Finance and Management



2d

OVERFITTING AND REGULARIZATION



What do we want to achieve with our ML model?

Plots of polynomials having various orders n

The plot shows the real function that we want to approximate, which is a part of the cosine function





Very important concepts: Generalization, underfitting and overfitting

- Generalization: model performs well on new, unseen data
- **Underfitting**: model is too simple to capture the patterns in the data
- **Overfitting**: model learns noise from the training data, leading to poor performance on new data



Ensuring good generalization by evaluating on unseen data

Train model on the training set and then test it on the test set





Use of regularization to prevent overfitting

Regularization: introducing a preference for simpler models

New cost function = original cost function + **penalty**

L₁ regularization

L₂ regularization

Penalizes **absolute values of the coefficients** (not the bias term).

Penalizes sum of the squares of the coefficients .



L1 (Lasso) regularization

- L₁ reduces the impact of less important features, shrinks their coefficients all the way to 0.
- Performs feature selection, produces a sparse model (i.e., with few nonzero feature coefficients)



Penalizes absolute values of the coefficients

The optimal parameters are those that minimize the new cost function.

hyperparameter: regularization strength



L2 (Ridge) regularization

- Does not eliminate features completely (coefficients can be very very small, but they are not zero).
- Reduces the magnitude of all coefficients

New cost function = MSE(Θ)+ $\alpha \sum_{i=1}^{n} \theta_i^2$

Penalizes sum of the squares of the coefficients.

The optimal parameters are those that minimize the new cost function.

hyperparameter: regularization strength