

Machine Learning :

* Machine Learning refers to a set of tools for understanding data. (or)

Machine learning is programming computers to optimize a performance criteria using example data (or) faster experience. (or)

Machine learning is all about using right features to build the right models that achieves the right tasks. (or)

Machine learning is all about extracting knowledge from data.

* Machine learning is also known as statistical learning (or). Predict to analytics.

History of Statistical Learning :

* At the beginning of 19th century the programs called least squares method was developed.

* Least square method is used for the implementation of linear regression.

* Linear regression is used to predict the quantitative values.

* In 1936, the linear Discriminant Analysis was developed to predict qualitative values.

* In 1940, Logistic Regression was developed.

- It is used for classification.

* In 1970, there have generalized linear models is the term and developed to describe the entire class of statistical learning methods.

* Linear models are used to generalize the linear relation.

* In 1980, there have developed the algorithm called the Classification and Regression tree (CART).

* In same 1980, neural networks are developed.

* In 1990, support vector machines were developed.

- * we ca
- * Medic
- * Astro
- * public
- * Statist
- * Ident
- * class
- * SNO ram
- * Predi
- the bas
- * cust
- * Ide
- * cla
- classes
- * Esta
- * variab
- * clas
- * Types
- * statis
- . underst
- * these
- 1) su
- 2) Un
- 3) Re
- 1) sup
- for
- more
- * on
- vari
- * ve

Machine learning:

- * Machine Learning refers to a set of tools for understanding data. (or) machine learning is programming computers to optimize a performance criteria using example data (or) faster experience. (or)
- * Machine learning is all about using right features to build the right models that achieves the right tasks.

Machine learning is all about extracting knowledge from data.

- * Machine learning is also known as statistical learning (or).

Predict to analytics.

History of Statistical Learning:

- * At the beginning of 19th century the programs called least squares method was developed.
- * Least square method is used for the implementation of linear regression.
- * Linear regression is used to predict the quantitative values.

In 1936, the linear Discriminant Analysis was developed to predict qualitative values.

In 1940, logistic regression was developed.

- It is used for classification.

In 1970, that have generalized linear models is the term and it is developed to describe the entire class of statistical learning methods.

Linear models are used to generalize the linear relation.

+ linear models are used to predict the algorithm called the

In 1980, there have developed the algorithm called the Classification and Regression tree (CART).

In same 1980, neural networks are developed.

In 1990, support vector machines were developed.

- ## Applications
- * we can use in medicine
 - * astro physics
 - * public policy
 - * Statistical
 - * Identity
 - * classify
 - * ram.
 - * Predict
 - * the basis of
 - * customize
 - * Identity
 - * classify
 - * classes to
 - * Establish variables
 - * classify
 - * Types of
 - * statistical understanding
 - * these to
 - supervise
 - Unsuper
 - Reinforce
 - 1.) Supervise
 - Supervise for pred more in
 - * outcome variable
 - * vector

- Applications of Machine learning:
- * we can use machine learning in business
 - * Medicine field - diagnosis of V. many diseases
 - * Astro physics - relationship between galaxies etc.
 - * public policy, etc..
 - * Statistical learning problems:
 - * identify the risk factors for prostate cancer.
 - * classify a recorded phoneme based on a log-periodogram.
 - * Predict whether someone will have a heart attack, are the basis of demographic, diet and clinical measurements.
 - * customize an e-mail spam detection system.
 - * Identify the numbers in a hand written zip code.
 - * classify a tissue sample into one of several cancer classes based on a gene expression profile.
 - * Establish the relationship between salary and demographic variables in population survey data.
 - * classify the pincels in a satellite image.

* Types of learning :-

- * statistical learning refers to a vast set of tools for understanding the data.
- * these tools is classified into 2 categories.
 - 1) supervised statistical learning.
 - 2) unsupervised statistical learning.
 - 3) Reinforcement Learning.
- 1) Supervised Learning: building a statistical model for predicting (or) estimating the output based on one (or) more inputs with the supervising outputs.
- * outcome measurement (y) is also called as dependent variable, response variable (or) target variable.
- * vector of P Predictors measurements ' x ' also called as

input variables, features, independent variables, predictors, covariates, regressors.

- * In regression problem, y is quantitative.
- * In the classification problem, y is qualitative.
- * Examples for classification.
- * Hand-written digits classification.

Unsupervised learning

Unsupervised learning involves building a statistical model for predicting y , estimating the outputs based on one or more inputs without supervising outputs.

- * what is statistical learning?

Notation.

Here, sales is a response (or) target that we wish to predict. Sales generally refers to the response as y . TV is a feature, or input or predictor, we name it x_1 , likewise name Radio as x_2 , and Soson as x_3 . We can refer to these inputs vector collectively as x .

$$x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}$$

Now we write our model as

$$y = f(x) + \epsilon$$

where ϵ captures measurement errors and other discrepancies.

(with a good f).

What is $f(x)$ good for?

* with a good f we can make predictions of y at new points $x = x^*$.

* we can understand which components of $x = (x_1, x_2, \dots, x_p)$ are important in explaining y , and which are irrelevant.

e.g., Seniority and Years of Education have a big impact on income, but Marital status typically does not.

* Depending on the complexity of f , we may be able to understand how each component x_j of x affects y .

* Points often x have many components.

The Regression Function $f(x)$

* Is there an ideal $f(x)$? In particular, what is a good value for $f(x)$ at any selected value of x , say $x=4$? There can be many y values at $x=4$. A good value is

$$f(4) = E(Y|X=4)$$

$E(Y|X=4)$ means expected value (average) of y given $X=4$.
 $E(Y|X=x)$ is called "the regression function".

* The regression function $f(x)$

* It is also defined for vectors X : e.g. $f(x_1, x_2, x_3) = E(Y|X_1=x_1, X_2=x_2, X_3=x_3)$. But $f(x) = f(x_1, x_2, x_3) = E(Y|X_1=x_1, X_2=x_2, X_3=x_3)$ is the ideal or optimal predictor of y with regard to mean-squared prediction error; $f(x) = E(Y|X=x)$ is the function that minimizes $E[(Y - g(x))^2 | X=x]$ over all functions g ; it all points x .

* How to estimate ' f '? 18/

* To estimate ' f ' it was classified into two approaches.

1) Linear Approaches

2) Non-linear Approaches

i) Training data,

ii) Statistical method on this data to learn the actual input and output.

To estimate the function ' f ' we may follow parametric methods and (a) non-parametric methods.

i) Parametric method

In this we follow 2 steps. In this functional form of ' f ' has to assume, if we have to assume, i.e., $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$. If it is linear, $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$.

ii) In this we can find the β values making well conditioned 'P' matrix, if the actual data is not a linear data.

* There is a disadvantage, i.e., the fitting process of fitting of linear data

* Advantage is we need not to bother about the functional form.

2) Non-Parametric form, (we must).

* How to estimate 'f'.

* To estimate function 'f' we may use linear approach or non-linear approach.

* We assume that we have observed a set of different data points (Training data). Because, we will use these observations to train our method how to estimate 'f'.

* Let x_{ij} represents the value of j th predictor for the input observation i , where $i = 1, 2, \dots, n$; $j = 1, 2, \dots, p$.

* y_i represents the response variable for i th observation, then, the training data consists of $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$.

$x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$, i.e., representation of the training data.

* Therefore, our goal is to apply a statistical learning method to training data inorder to estimate the unknown function 'f'.

$$y \approx f(x)$$

* Now the statistical method for the above task can be characterized as either parametric (or) non-parametric.

* Parametric method. Parametric methods involve a two step model based approach.

Parametric methods involve a two step model based approach.

i) we make an assumption about the functional form of 'f'.

that assumption is 'f' is linear.

$$f(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

x_1, x_2, \dots are predictors

Once we have assumed that 'f' is linear. Now, the problem is to estimate parameters.

E) * After the model uses the most least squares

* Dis-advantage

① the mode true unknown to far from estimate will

* Non-parametric

* Non-parametric

about the points

* It estimates

points

Ex: thin plate

* Assessing

* Inorder to

learning measure

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

* If MSE actual out

* If MSE

actual out

* MSE is

* Model mea

(or) statistical

* Always

* Training

* The MS

- * After the model has been selected, we need the procedure that uses the training data to fit (f) to train the model.
- * The most common approach to train the model is least squares approach.

* Dis-advantages of the parametric approach :-

- * The model we choose will usually not natural form to true unknown form of ' y '. If the chosen model is far from the true functional form of ' y ', then, our estimate will be very poor, which needs to overfitting.
- * Non-parametric methods: (do not) makes any assumption about the functional form of ' f '.
- * It estimates the ' f ', which is very close to the data points; make splines or piecewise splines.
- * Thin plate spline.

* Assessing the model :-

- * Assessing the performance of a statistical learning method on a given data set, we need to measure mean squared error, $\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$.
- * If MSE is small, the predictions are very close to the actual outputs.
- * If MSE is large, predictions are far away from the actual outputs.
- * Model means
- * Statistical learning method
- (or) statistical learning must be small.
- * Always MSE must be small.
- * Training MSE:
- * The MSE is computed using the training data that way

Used to fit the model.

- * we would like to select the model for which the test MSE is as small as possible

* The Bias - Variance Trade off :-

- * the expected test MSE for a given value x_0 can always be decomposed into 3 fundamental quantities,
 - 1) the variance of $\hat{f}(x_0)$
 - 2) the squared bias of $\hat{f}(x_0)$.
 - 3) The variance of error (residual error) ϵ .

$$E(Y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + (\text{Bias}(\hat{f}(x_0)))^2 + \text{Var}(\epsilon)$$

* Variance :-

- * Variance refers to about amount by which \hat{f} will change if we estimate it using a different training data set.

* Bias :-

Bias refers to the error i.e., introduced by approximation of a real life problem.

- * If we want to minimize the expected test MSE we need to select the model that simultaneously achieves low variance and low bias.

21/1/22

* Simple Linear Regression using single predictor x :-

- * we assume a model

$$Y = \beta_0 + \beta_1 x + \epsilon$$

β_0, β_1 are two unknown constants that represent and slope also known as coefficients (or) parameters & ϵ is error term or ϵ -term.

- * Given some estimates $\hat{\beta}_0$ & $\hat{\beta}_1$ for model coefficients we predict future sales when using

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

- * where \hat{y} indicates that prediction of y is the basis of the regression is x .

21/1/22
Sup

is linear
supervised

x_1, x_2 :-

* True
simple U

* we are

whereas
the int
parameters

* Given
we pred

where

$x = x$.

Residual :-

Predicted

* Estimate

* Let \hat{y}_i
the i^{th}
residual

* we de

R_{SS}

or equivale

$R_{SS} = ($

* the less

the R_{SS}

$\hat{\beta}_1$

21/7/22 Supervised Learning

* Linear Regression: Linear Regression is a simple approach to supervised learning. It assumes that the dependencies of y on x_1, x_2, \dots, x_p is linear.

* True regression function never a linear.

* Simple linear regression using a single predictor x

* we assume a model

$$y = \beta_0 + \beta_1 x + \epsilon$$

where β_0, β_1 are two unknown constants that represent the intercept and slope, also known as coefficients or the parameters and ϵ is the error term.

* Given some estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ for the model coefficients, we predict future sales using

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

where \hat{y} indicates a prediction of y on the basis of $x = x_i$. then that symbol denotes an estimated value.

Residual: The difference between actual output and the predicted output caused by residual.

Predicted of the parameters by least squares

* Estimation: let $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ be the prediction for y_i based on

* let $e_i = y_i - \hat{y}_i$ represents the i th residual.

* we define the residual sum of squares (RSS)

$$RSS = e_i^2 + e_2^2 + \dots + e_n^2$$

or equivalently as

$$RSS = (y_1 - \hat{\beta}_0 - \hat{\beta}_1 x_1)^2 + (y_2 - \hat{\beta}_0 - \hat{\beta}_1 x_2)^2 + \dots + (y_n - \hat{\beta}_0 - \hat{\beta}_1 x_n)^2$$

* the least squares approach chooses $\hat{\beta}_0$ and $\hat{\beta}_1$ to minimize the RSS. The minimizing values can be shown to be

$$\hat{\beta}_0 + \hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

where, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ are the sample means.

* Assessing the accuracy of the coefficient estimates:

- * The standard error of an estimator reflects how it varies under repeated sampling we have,

$$SE(\hat{\beta}_1)^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}, SE(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]$$

where, $\sigma^2 = \text{var}(e)$

- * These standard error can be used to compute confidence intervals. A 95% confidence interval is defined as a range of values such that with 95% probability, the range will contain the true unknown value of the parameter. It has

$$\hat{\beta}_1 \pm 2 \cdot SE(\hat{\beta}_1)$$

that is, there is approximately a 95% chance that the interval lies between $[\hat{\beta}_1 - 2 \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot SE(\hat{\beta}_1)]$.

$$[\hat{\beta}_1 - 2 \cdot SE(\hat{\beta}_1), \hat{\beta}_1 + 2 \cdot SE(\hat{\beta}_1)]$$

- * For the advertising data, the 95% confidence interval for β_1 is $[0.042, 0.053]$.

* Hypothesis testing: It is used to test on the

- * Standard Error can also be performed hypothesis tests on the coefficients. The most common hypothesis tests involves testing the null hypothesis of

H_0 : there is no relationship between x and y , versus the alternative hypothesis.

H_A : There is some relationship between x and y .

- * Mathematically this corresponds to testing

$$H_0: \beta_1 = 0$$

versus

$$H_A: \beta_1 \neq 0$$

since, if $\beta_1 = 0$ then the model reduces to $y = \beta_0 + e$ and x is not associated with y .

- * To test
- $t = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)}$
- * This will assuming β_1
- * using stat of observing the p-value
- * R-Code

```
library(MASS)
library(Car)
Ad <- Ad
head(Ad)
lm.fit <- lm.fit(Ad ~ TV)
attach(lm.fit)
summary(lm.fit)
names(lm.fit)
coef(lm.fit)
confint(lm.fit)
predict(lm.fit)
plot(lm.fit)
abline(lm.fit)
```

- * To test the null hypothesis, we compute a t-statistic, given by
$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)}$$
- * This will have a t-distribution with $n-2$ degrees of freedom, assuming $\beta_1 = 0$.
- * Using statistical software, it is easy to compute the probability of observing any value equal to $|t|$ (or) larger we call this probability the p-value.
- * R-Code for Linear Regression :-

```

library(MASS)
library(ISLR2)
Ad <- read ("Advertising.csv", header = TRUE)
head(Ad)
lm.fit <- lm (sales ~ TV, data = Ad)
attach(Ad)
lm.fit <- lm(sales ~ TV)
lm.fit
summary(lm.fit)
names(lm.fit)
coef(lm.fit)
conf.int(lm.fit)
predict(lm.fit, data.frame(TV=c(200, 150, 90)), interval="confidence")
plot(TV, sales)
abline(lm.fit, lwd=3, col="red")

```

* To test the null hypothesis, we compute a t-statistic, given by

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{SE}(\hat{\beta}_1)}$$

* This will have a t-distribution with $n-2$ degrees of freedom, assuming $\beta_1 = 0$.

* Using statistical software, it is easy to compute the probability of observing any value equal to $|t|$ (or) larger we call this probability the p-value.

* R-Code for Linear Regression :-

library(MASS)

library(ISLR2)

Ad \leftarrow read ("Advertising.csv", header = TRUE)
head(Ad)

lm.fit \leftarrow lm(Sales ~ TV, data = Ad)

attach(Ad)

lm.fit \leftarrow lm(sales ~ TV)

lm.r.f

summary(lm.fit)
names(lm.fit)

coef(lm.fit)

conf.int(lm.fit)

predict(lm.fit, data.frame(TV=c(200, 150, 90)), interval = "Confidence")

plot(TV, sales)

abline(lm.fit, lwd = 3, col = "red")

* Assessing the overall Accuracy of the Model :-

* Assessing the Standard Error.

* We Compute the Residual Standard Error.

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

where, the residual sum-of-squares is RSS $= \sum_{i=1}^n (y_i - \hat{y}_i)^2$.

* R-Squared or fraction of variance explained is $R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}$

TSS

where, $TSS = \sum_{i=1}^n (y_i - \bar{y})^2$ is the total sum of squares.

- * It can be shown that in this simple linear regression setting that $R^2 = \rho^2$, where ρ is the correlation between x and y .
- * $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$

Quantity	Value
Residual standard error	3.26
R^2	0.612
F-statistic	312.1

* Multiple Linear Regression: It is like simple linear regression, but number of inputs will more than 1.

* Here our model is,

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

* Multiple linear regression is supervised learning approach, for predicting quantitative response based on more than one input.

* It can be model as $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$, we interpret β_j as the average effect on y of a one unit increase in x_j , holding all other predictors fixed. In the advertising example, the model becomes

$$\text{sales} = \beta_0 + \beta_1 \times \text{TV} + \beta_2 \times \text{radio} + \beta_3 \times \text{newspaper} + \epsilon$$

library(MASS)

library(ISLR2)

Ad <- read.csv("Advertising.csv", header = TRUE)

head(Ad)

lm.out <- lm(sales ~ TV + radio + newspaper, data = Ad)

attach(Ad)

(lm.fit <- lm(sales ~ TV))

(lm.fit)

summary(lm.fit)

output

coefficients

(intercept)

TV

radio

newspaper

Correlations

TV

radio

newspaper

sales

LS

atle

FOR

algm(TSS)

F =

RS

F-stat

If F-S

input and

output

Quantity

RS-E

R2

f-statistic

resid

predic

resid

(Im. fit)

summary(Im. fit)

Output

	estimate	std. error	t value	Pr(> t)
(Intercept)	2.938889	0.311908	9.422	<2e-16 ***
TV	0.045765	0.001395	32.809	<2e-16 ***
radio	0.188530	0.008611	21.893	<2e-16 ***
newspaper	-0.001037	0.005871	-0.177	0.86

Correlations

	TV	radio	newspaper	sales
TV	1.0000	0.0548	0.0567	0.7822
radio		1.0000	0.3541	0.5962
newspaper			1.0000	0.2283

Sales

- * Is at least one predictor useful?
- * For 1st one use F-statistic
- * $F = \frac{(TSS - RSS)/P}{RSS/(n-P-1)} \sim F_{P, n-P-1}$
- * F-statistic should be large.
- * If F-statistic = 1, then there is no relation between input and output variable.

Input and output variable

Quantity

value

RSS

R²

f-statistic

multiple regression this denotes a linear relationship between independent variables and the dependent variable.

29/10/22

* Depending on the important variables :-

- * the most direct approach is called all subsets or best subset regression : we compute the least squares fit for all possible subsets and then choose between them based on some criterion that balances training error with model size.
- * However we often can't examine all possible models, since they are 2^p of them; for example, when $p=40$ there are over a billion models!
- * Instead we need an automated approach that searches through a subset of them. We discuss two commonly used approaches next.
- * Forward Selection :-
 - * Begin with the null model - a model that contains an intercept but no predictors.
 - * Fit p simple linear regressions and add to the null model the variable that results in the lowest RSS.
 - * Add to that model the variable that results in the lowest RSS amongst all two-variable models.
 - * Continue until some stopping rule is satisfied, for example, when all remaining variables have a p-value above some threshold.
- * Backward Selection :-
 - * Start with all variables in the model.
 - * Remove the variables with the largest p-value — that is, the variable that is the least statistically significant.
 - * the new $(P-1)$ -variable model is fit, and the variable with the largest p-value is removed.
 - * Continue until a stopping rule is reached. For instance, we may stop when all remaining variables have a significant p-value defined by some significance threshold.

* Other consider

* Qualitative p

* Some pred taking a di

* These are variables.

(* See for ex card.)

Example: In males and a new variab

Resulting mod

$$y_i = \beta_0 + \beta_1 x_1 +$$

(Results for

* Qualitative p

* with more + variables. for two dummy

x_{11}

and the s

x_{12}

* then both

regression eq

$$y_i = \beta_0 + \beta_1 x_{11} +$$

- * Other considerations in the regression model:
- * Qualitative predictors:
 - * Some predictors are not quantitative but are qualitative, taking a discrete set of values.
 - * These are also called categorical predictors or factor variables.

(* See for example the scatterplot matrix of the credit card.).

Example: Investigate differences in credit card balance between males and females, ignoring the other variables, we create a new variable.

$$x_i = \begin{cases} 1 & \text{if } i\text{th person is female} \\ 0 & \text{if } i\text{th person is male.} \end{cases}$$

Resulting model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ if i th person is female.

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i = \begin{cases} \beta_0 + \epsilon_i & \text{if } i\text{th person is male.} \end{cases}$$

(Results for gender model):

* Qualitative predictors with more than 2 variables

* with more than two levels, we create additional dummy variables. For example, for the ethnicity variable we create two dummy variables. The first could be

$$x_{i1} = \begin{cases} 1 & \text{if } i\text{th person is Asian} \\ 0 & \text{if } i\text{th person is not Asian.} \end{cases}$$

and the second could be

$$x_{i2} = \begin{cases} 1 & \text{if } i\text{th person is caucasian} \\ 0 & \text{if } i\text{th person is not caucasian.} \end{cases}$$

* then both of these variables can be used in the regression equation, in order to obtain the model

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i = \begin{cases} \beta_0 + \epsilon_i & \text{if } i\text{th person is Asian} \\ \beta_0 + \beta_2 + \epsilon_i & \text{if } i\text{th person is caucasian} \\ \beta_0 + \epsilon_i & \text{if } i\text{th person is AA.} \end{cases}$$

* There will always be one or fewer dummy variable than the number of levels. The levels with no dummy variable, African American in this example - is known as the base line.

3/4/21 * Extensions of the linear model :-

* In linear models,

* The additive assumption implies that the effect of change in a predictor x_j on the response y is independent of the values of the other predictors.

* Linear assumption implies that the change in the response 'Y', due to a '1' unit change in x_j is constant, regardless of the value of x_j .

* Interactions :-

* In our previous analysis of the Advertising data, we assumed that the effect on Sales of increasing one, the advertising medium is independent of the amount spent on the other media.

* For example, the linear model

$$\text{Sales} = \beta_0 + \beta_1 x_{\text{TV}} + \beta_2 x_{\text{radio}} + \beta_3 x_{\text{newspaper}}$$

states that the average effect on sales of a one-unit increase in TV is always β_1 , regardless of the amount spent on radio.

* Modelling interactions

Model takes the form:

$$\text{Sales} = \beta_0 + \beta_1 x_{\text{TV}} + \beta_2 x_{\text{radio}} + \beta_3 x_{(\text{radio} \times \text{TV})} + \epsilon$$

$$= \beta_0 + (\beta_1 + \beta_3 x_{\text{radio}}) x_{\text{TV}} + \beta_2 x_{\text{radio}} + \epsilon.$$

Results :-

	Coefficient	Std. Error	t-statistic	P-value
Intercept	6.7502	0.248	27.23	<0.0001
TV	0.0191	0.002	12.90	<0.0001
radio	0.0289	0.009	3.24	0.0004
TV x radio	0.0041	0.000	20.73	<0.0001

* Hierarchy

If we

include the their coefficients

* Interaction

Consider the Predictor balance (qualitative)

* without a balance:

$$= \beta_1 x_{\text{TV}}$$

with interaction

balance is

$$= \begin{cases} (\beta_0 + \beta_1 x_{\text{TV}}) \\ \beta_2 x_{\text{radio}} \end{cases}$$

4/8/21

* Parametric

* A non

K-nearest
the K-n

* Given a

nearest ne
observations

f(x₀) is

of N(μ , σ^2)

* Hierarchy principle :-

If we include an interaction in a model, we should do include the main effects, even if the p-values associated with their coefficients are not significant.

* Interaction between qualitative and quantitative variables:-

Consider the credit data set, and suppose that we wish to predict balance using income (quantitative) and student (qualitative).

* without an interaction term, the model takes the form:-
balance $\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 & \text{if } i^{\text{th}} \text{ person is a student} \\ \beta_0 + \beta_3 & \text{if } i^{\text{th}} \text{ person is not a student} \end{cases}$

$$= \beta_1 \times \text{income}_i + \begin{cases} \beta_0 + \beta_2 & \text{if } i^{\text{th}} \text{ person is a student} \\ \beta_0 & \text{if } i^{\text{th}} \text{ person is not a student.} \end{cases}$$

with interactions, it takes the form:-
balance $\approx \beta_0 + \beta_1 \times \text{income}_i + \begin{cases} \beta_2 + \beta_3 \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student.} \end{cases}$

$$= \begin{cases} (\beta_0 + \beta_2) + (\beta_1 + \beta_3) \times \text{income}_i & \text{if student} \\ \beta_0 + \beta_1 \times \text{income}_i & \text{if not student.} \end{cases}$$

4/8/22 * Parametric Methods versus Non-Parametric Methods :-

* A non parametric method akin to linear regression is K-nearest neighbors regression which is closely related to K-nearest neighbors classifier.

* Given a value for K and a prediction point x_0 , K-nearest neighbors regression first identifies the K -nearest neighbors that are closest to x_0 , represented by No. observations.

$f(x_0)$ is then estimated using the average of y_i 's like so

$$f(x_0) = \frac{1}{K} \sum_{x_i \in N_0} y_i.$$

* A parametric approach will out perform a non-parametric approach if the parametric form is close to true form of $f(x)$.

* the choice of a parametric approach versus a non-parametric approach will depend largely on bias-variance trade-off & shape of function $f(x)$.

* when the true relationship is linear, it is difficult for a non-parametric approach to compete with the linear regression because non-parametric approach incurs a cost in variance that is not offset by a reduction in bias.

* Additionally, in higher dimensions, K-nearest neighbors regression often performs worse than linear regression.

* This is often due to combining too small an n with too large a K , resulting in a given observation having no nearby neighbors.

* This is often called the curse of dimensionality. In other words, the K -observations nearest to an observation may be far away from x_0 in a p -dimensional space when p is large, leading to a poor prediction of $f(x_0)$ and a poor K-nearest neighbors regression fit.

* As a general rule, parametric models will tend to outperform non-parametric models when there are only a small no. of observations per predictor.

* classification
Set C, such eye color email $\in \{$

* Given a taking values build a vector X and

* Often we that X below

* For example of the probability than a car

10/8/22 * Logis

* Logistic regression particular cat

* Logistic regression Prediction formula

$P(X)$

* This yields

* The logistic

estimation

* $P(X) = \frac{1}{1+e^{-f(x)}}$

between

* As an since

* taking

* classification: qualitative variables takes values in an ordered set C , such as:

eye color $\in \{\text{brown, blue, green}\}$,

email $\in \{\text{spam, ham}\}$.

* Given a feature vector x and a qualitative response y taking values in the set C , the classification task is to build a function $c(x)$ that takes as input the feature vector x and predicts its value for y : i.e., $c(x) \in C$.

* Often we are more interested in estimating the probabilities that x belongs to each category in C .

* For example, it is more valuable to have an estimate of the probability that an insurance claim x is fraudulent, than a classification fraudulent or not.

10/8/22

Logistic Regression

* Logistic regression models the probability that y belongs to a particular category rather than modeling the response itself.

* Logistic regression uses the logistic function to ensure a prediction between 0 and 1. The logistic function takes the form

$$P(X) = \frac{e^{B_0 + P_1 X}}{1 + e^{B_0 + P_1 X}}$$

This yields a probability greater than 0 and less than 1.

* The logistic function can't be rebalanced to yield 0 or 1.

$$\text{otherwise } 0 + P(X) = \frac{e^{B_0 + P_1 X}}{1 - e^{B_0 + P_1 X}}, \quad 1 - P(X) = \frac{1}{1 + e^{B_0 + P_1 X}}$$

* $\frac{P(X)}{1 - P(X)}$ is known as the odds and takes on a value between 0 and infinity.

* As an example, a probability of 1 in 5 yields odds of 4 to 1.

* since $\frac{0.2}{1 - 0.2} = \frac{1}{4}$, better fit with a prior estimate of 1/4.

* Taking a logarithm of both sides of the logistic odds equation yields an equation for the log-odds or logit.

$$\log \left[\frac{P(X)}{1-P(X)} \right] = \beta_0 + \beta_1 X$$

logistic regression has a logit that is linear in terms of X .

- * Unlike linear regression where β_1 represents the average change in Y with one-unit increase in X , for logistic regression, increasing X by one-unit yields a β_1 change in the log-odds which is equivalent to multiplying the odds by e^{β_1} .

* The relationship between $\log P(X)$ and X is not linear and because of this β_1 does not correspond to the change in $P(X)$ given one-unit increase in X . However, if β_1 is positive, increasing X will be associated with an increase in $P(X)$ and, similarly, if β_1 is negative, an increase in X will be associated with a decrease in $P(X)$. How much change will depend on the value of X .

Estimating Regression Coefficients

- * logistic regression uses a strategy called maximum likelihood to estimate regression coefficients.
- * Maximum likelihood plays out like this so to determine estimates for β_0 and β_1 such that the predicted probability of $\hat{P}(X_i)$ corresponds with the observed classes as closely as possible. Formally, this yield an equation called a likelihood function.

$$L(\beta_0, \beta_1) = \prod_{i:y_i=1} p(X_i) \times \prod_{j:y_j=0} (1-p(X_j))$$

- * Estimates for β_0 and β_1 are chosen so as to maximize this likelihood function.

- * linear regression's least squares approach is actually a special case of maximum likelihood.
- * logistic regression measures the accuracy of the coefficient estimates using a quantity called the z-statistic. The z-statistic is similar to the t-statistic. The z-statistic for β_1 is represented by

17/8/22

* Discriminant

	outlook
0	Rainy
1	Rainy
2	overcast
3.	sunny
4.	sunny
5,	sunny
6.	overcast
7.	Rainy
8.	Rainy
9.	Sunny
10.	Rainy
11.	overcast
12.	overcast
13.	sunny

Yes - 9

outlook:

Rainy - 5

overcast - 4

sunny - 5

$P(Y=x) =$

$$z\text{-statistic } (P_i) = \frac{\hat{P}_i}{SE(\hat{P}_i)}$$

17/8/22

* Discriminant Analysis :-

	outlook	temperature	humidity	windy	play golf
0.	Rainy	Hot	High	False	NO
1.	Rainy	Hot	High	True	NO
2.	overcast	Hot	High	False	YES
3.	sunny	Mild	High	False	YES
4.	sunny	Cool	Normal	False	YES
5.	sunny	Cool	Normal	True	NO
6.	overcast	Cool	Normal	True	YES
7.	Rainy	Mild	High	False	NO
8.	Rainy	Cool	Normal	False	YES
9.	Sunny	Mild	Normal	True	YES
10.	Rainy	Mild	Normal	True	YES
11.	overcast	Mild	High	False	YES
12.	overcast	Hot	Normal	True	NO
13.	sunny	Mild	High	True	YES

Yes - 9 ; No - 5

outlook :-

Rainy - 5 → 2
→ 3

overcast - 4 → 4
→ 0

sunny - 5 → 3
→ 2
 $P(Y=x) = \frac{P(x=4)}{P(x)}$

$x = \text{sunny}$.

$$P(Y = \text{Yes} | X) = \frac{P(X | \text{Yes}) * P(\text{Yes})}{P(X)}$$

$$= \frac{\frac{3}{9} \times \frac{9}{14}}{\frac{5}{14}} = 0.6$$

$$P(Y = \text{No} | X) = \frac{P(X | \text{No}) * P(\text{No})}{P(X)} = \frac{\frac{4}{5} \times \frac{5}{14}}{\frac{5}{14}} = 0.4$$

Humidity

High - 7 - $\frac{7}{14}$

3
4

Normal - 7 - $\frac{7}{14}$

6
1

$$P(\text{High} | \text{Yes}) = \left(\frac{3}{9} \times \frac{9}{14} \right) / \frac{9}{14} = \frac{3}{7} = 0.42857$$

$$P(\text{High} | \text{No}) = \left(\frac{4}{5} \times \frac{5}{14} \right) / \frac{5}{14} = \frac{4}{7} = 0.5714$$

- * let us consider we have data of 13 observations (0-12).
- * By using that data we need to predict the response of playgolf on features or predictors given. Yes - 9; No - 4.

outlook

Rainy - 5 - 2
4
3
N

temperature

hot - 4 - 2
2
N

overcast - 4 - 4

Mild - 5 - 4

0
N

4
1
N

sunny - 4 - 3

Cool - 4 - 3

1
N

4
1
N

Humidity:

High - 6 - 3
4
3
N

windy

True - 5 - 3
4
2
N

Normal - 7 - 6
4
1
N

Mild - 6 - 4
2
N

UNIT-11

* Resampling methods :-

(to get additional information
we are refitting).

Resampling methods involves repeatedly drawing samples from the given data, and refit the model on each sample to get additional info about the model.

* There are 2 resampling methods,

1.) cross validation.

2.) Boot strap.

1.) Cross Validation: It measures the test error rate associated with a statistical learning method to evaluate the model performance.

2.) Boot strap: It measures the accuracy of parameter estimation (or) of given model.

* Validation set:- Before going to know about cross validation we need to know about validation set.

* the data is divided into 2 groups.

1.) train error

2.) testing error.

* (How to implement K-fold cross validation?)

* Cross validation:-

Training Error versus Test error:-

* Recall the distinction between the test error and the training error.

* The test error is the average error that results from using a statistical learning method to predict the response on a new observation, one that was not used in training the method.

* In contrast, the training error can be easily calculated by applying the statistical learning method to the observations used in its training.

* But the training error rate often is quite different from the test error rate, and in particular the former can dramatically underestimate the latter.

* More on prediction - error estimates:-

* Best solution: a large designated test set, often not available.

* Some methods make a mathematical adjustment to the training error rate in order to estimate the test error rate. These include the CP statistic, AIC and BIC. (Akaike's information criterion) - they are discussed elsewhere in this course.

* Here, we instead consider a class of methods that estimate the test error by holding out a subset of the training observations from the fitting process, and then applying the statistical learning method to those held out observations.

* validation

* Here, we two parts: set.

* The model model is in the val

* The resulting test error case of qu the case

* the valid

123

2

A random sight part is

* Drawbacks

* The validation variable, dependent included in included in

* In the v - those than in the model.

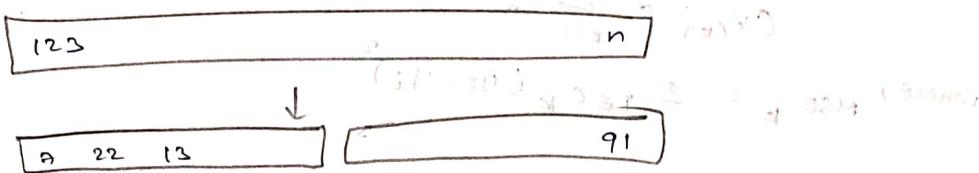
* This suggest overestimate entire data

* K-fold cross

* widely used

* Estimates an idea of

- * validation set approach: ~~cross-validation approach~~
- * Here, we randomly divide the 'available' set of samples into two parts: a training set and a validation or hold-out set.
- * The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.
- * The resulting validation-set provides an estimate of the test error. This is typically assessed using MSE, in the case of quantitative response and misclassification rate in the case of qualitative (discrete) response.
- * the validation process:



A random splitting into two halves: left part is training set, right part is validation set.

* Drawbacks of validation set approach:

- * The validation estimate of the test error can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set.
- * In the validation approach, only a subset of the observations - those that are included in the training set rather than in the validation set - are used to fit the model.
- * This suggests that the validation set error may tend to overestimate the test error for the model fit on the entire data set.

* K-fold Cross-validation:

- * widely used approach for estimating test error.
- * Estimates can be used to select best model, and to give an idea of the test error of the final chosen model.

* Idea is to randomly divide the data into k equal-sized parts. We leave out part k , fit the model to the other $k-1$ parts^{combined} and then obtain predictions for the left-out k^{th} part.

* This is done in turn for each part, $k=1, 2, \dots, K$ and then results are combined.

* The details?

* Let these K parts be C_1, C_2, \dots, C_K , where, C_k denotes the indices of the observations in part k . There are n_k observations in part k : if n is a multiple of K , then $n_k = n/K$.

* Compute:

$$CV_{CKN} = \sum_{k=1}^K \frac{n_k}{n} \text{MSE}_{C_k}$$

$$\text{where, } \text{MSE}_k = \sum_{i \in C_k} (y_i - \hat{y}_i)^2$$

The problem of leaving out a random part and fitting a model to the remaining $K-1$ parts is called K -fold cross-validation.

For example, consider a dataset consisting of n observations using two variables x_1 and x_2 to predict y . The observations can be partitioned into K folds. The first fold contains observations $i_1, i_2, \dots, i_{n/K}$. The second fold contains observations $i_{n/K+1}, i_{n/K+2}, \dots, i_{2n/K}$. The third fold contains observations $i_{2n/K+1}, i_{2n/K+2}, \dots, i_{3n/K}$. This process continues until all observations have been included in a fold.

For each fold, the remaining $K-1$ folds are used to train a model, and the fold being left out is used to test the model's performance. This process is repeated for all K folds. The final error estimate is the average of the errors from all K folds.