

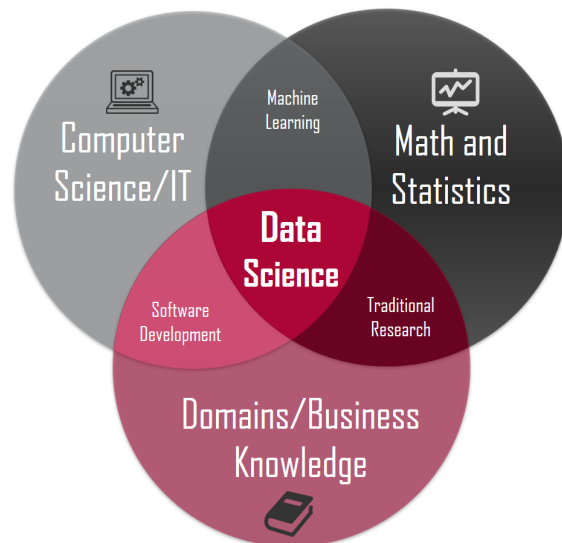
Modellazione statistica (statistical modeling)



Data Science

DS IS A MULTIDISCIPLINARY SUBJECT

- It starts from user-specified objectives
- It exploits algorithms to extract patterns and models, with a mathematic approach



Data Science life cycle

Outline

1) Regressione lineare semplice

STATISTICA NUMERICA, CAP. 6.4

1) Regressione lineare semplice

STATISTICA NUMERICA, CAP. 6.4

Modellazione statistica

La modellazione statistica, cioè la creazione di modelli su basi probabilistiche e non deterministiche,
Richiede una interazione con i dati.

Di solito si parte da una osservazione visuale dei dati, evidenziando relazioni e/o correlazioni fra di essi.

Come risultato si definisce un modello che possa rappresentare i dati.

Pacchetti Python per modellazione statistica:

- [statsmodel](#)
- [Scikit-learn](#)

Correlazione

Date due variabili aleatorie, la **correlazione** misura la relazione esistente fra di esse.

La **regressione lineare** è invece un modello che permette di fare previsione di una variabile conoscendo l'altra.

Il coefficiente di correlazione è un numero che risponde alla domanda:

«Sono le due variabili in relazione fra di loro? Se sì, quando cambia una, come cambia l'altra?»

Correlazione

Coefficiente di correlazione di Pearson fra le variabili aleatorie X e Y, dati i campioni X_i e Y_i

$$r = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

Il valore è nell'intervallo $[-1,1]$.

Correlazione

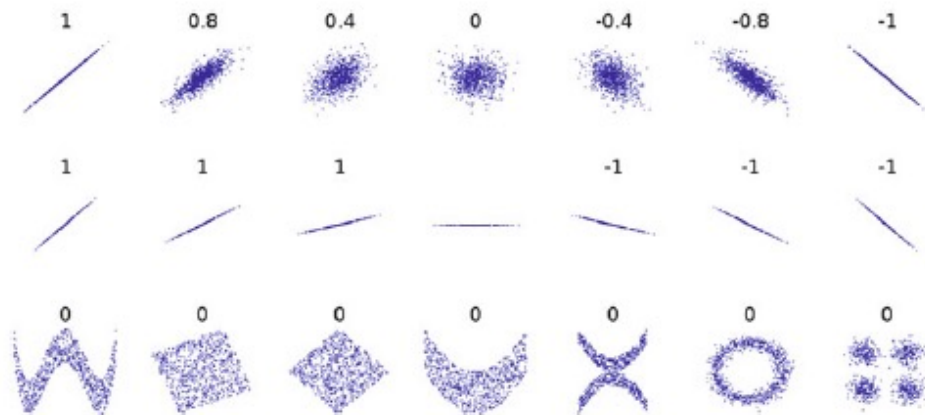


Fig. 11.1 Several sets of (x, y) points, with the correlation coefficient of x and y for each set. Note that the correlation reflects the noisiness and direction of a linear relationship (*top row*), but not the slope of that relationship (*middle*), nor many aspects of nonlinear relationships (*bottom*). N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero. (In Wikipedia. Retrieved May 27, 2015, from http://en.wikipedia.org/wiki/Correlation_and_dependence.)

In python:

```
from scipy.stats import pearsonr
```

File esempio:

[Example_correlation.py](#)

<https://docs.scipy.org/doc/scipy/reference/generated/scipy.stats.pearsonr.html>

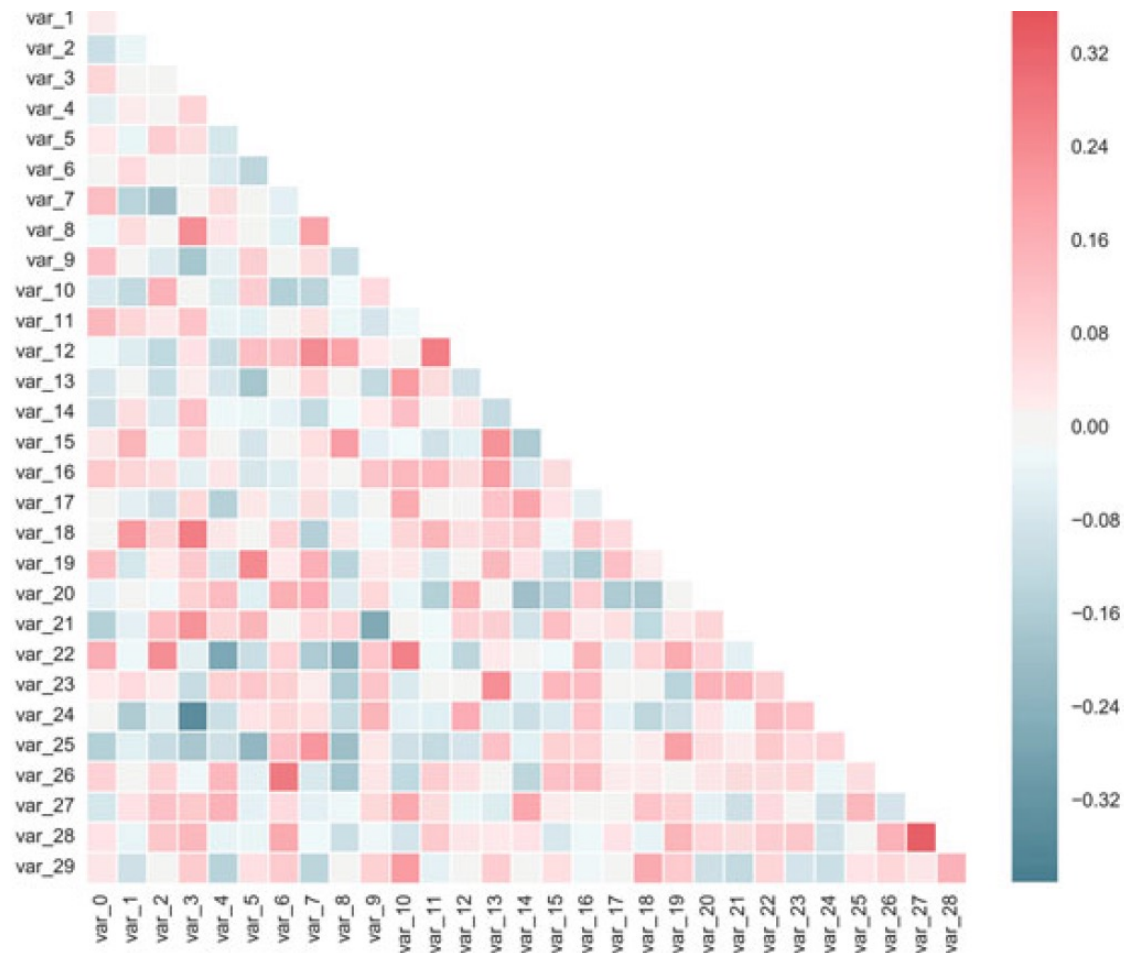
Matrice di correlazione

Un altro modo per visualizzare le correlazioni di coppie di variabili e' la **matrice di correlazione**.

```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="darkgrid")

rs = np.random.RandomState(33)
d = rs.normal(size=(100, 30))

f, ax = plt.subplots(figsize=(9, 9))
cmap = sns.diverging_palette(220, 10, as_cmap=True)
sns.heatmap(d, annot=False, sig_stars=False,
            diag_names=False, cmap=cmap, ax=ax)
f.tight_layout()
```



Matrice di correlazione

Un altro modo per visualizzare le correlazioni di coppie di variabili e' la **matrice di correlazione**.

Regressione lineare: introduzione

La regressione lineare è il cuore della statistica.



Risponde alla domanda:

« come posso utilizzare i dati che ho misurato per fare previsioni su dati che non conosco? »

utilizzando in particolare un modello di tipo lineare.

La *regressione lineare* è la parte della statistica che studia la relazione fra due o più variabili, che sono legate in modo NON DETERMINISTICO, per fare inferenze sul modello.

In particolare, si usano relazioni fra due o più variabili in modo da potere avere informazioni su una di esse conoscendo i valori dell'altra. Esempi di variabili che non sono legate fra loro da una relazione deterministica: x = l'età di un bambino e Y =la sua altezza, x =il volume di un motore e Y =il suo consumo di carburante, x =tempo di studio e Y = voto all'esame, ecc. Poiché x non è una variabile casuale la indichiamo con la lettera minuscola, mentre Y , che è un variabile casuale, viene indicata con la lettera maiuscola. Nel caso lineare supponiamo una relazione appunto lineare fra le due variabili x e Y :

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

Regressione lineare: introduzione

Regressione lineare: introduzione

. Questa relazione, di per sè deterministica, viene generalizzata a una relazione probabilistica. Date quindi informazioni su x e Y , l'obiettivo è quello di *predire* un valore futuro di Y per un particolare valore di x .

In questo modello, x viene detta *variabile indipendente* e Y viene detta *variabile dipendente*.

Il modello viene costruito a partire da alcune osservazioni $(x_i, Y_i), i = 1, \dots, n..$

L'estensione al modello probabilistico è necessaria nel momento in cui le due variabili non hanno una relazione deterministica. In pratica, in corrispondenza di n variabili indipendenti x_1, x_2, \dots, x_n si hanno n valori Y_1, Y_2, \dots, Y_n , che sono legati dalla relazione:

$$Y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \quad i = 1, \dots, n$$

quindi differiscono, rispetto al modello lineare esatto, di una quantità ϵ_i .

I valori Y_i sono in generale variabili aleatorie.

Regressione lineare: introduzione

Regressione lineare: introduzione

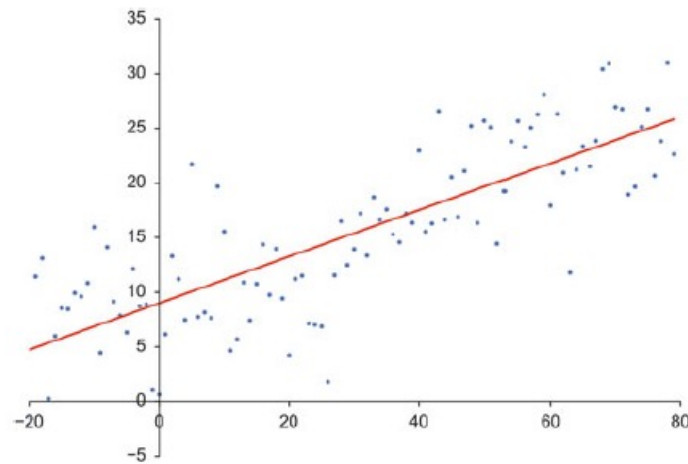


Fig. 11.2 Best-fit linear regression line to a given set of data

Stima dei parametri

‘Come calcolare stime dei parametri β_0 e β_1 della retta di regressione lineare assegnate le coppie $(x_i, Y_i), i = 1, \dots, n$? Cioé come determinare, fra le infinite rette del piano, una buona retta? Esiste una retta migliore delle altre?’

Visto che gli errori ϵ_i hanno distribuzione `norm(mean=0, sd= σ)`, allora la variabile aleatoria $Y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ ha distribuzione normale con deviazione standard σ . La funzione di verosimiglianza è:

$$\begin{aligned} L(\beta_0, \beta_1) &= \prod_{i=1}^n f_{Y_i}(x_i) \\ &= \prod_{i=1}^n (2\pi\sigma^2)^{-1/2} \exp \left\{ \frac{-(Y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ \frac{-\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2} \right\}; \end{aligned}$$

Stima dei parametri

facendo il logaritmo naturale di $L(\beta_0, \beta_1)$ si ha:

$$F(\beta_0, \beta_1) = \ln(L(\beta_0, \beta_1)) = -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{\sum_{i=1}^n (Y_i - \beta_0 - \beta_1 x_i)^2}{2\sigma^2}.$$

Per minimizzare questa funzione rispetto alle variabili β_0 e β_1 :

$$\frac{\partial F}{\partial \beta_0} = 0, \quad \frac{\partial F}{\partial \beta_1} = 0.$$

Quindi:

$$\frac{\partial F}{\partial \beta_0} = -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(Y_i - \beta_0 - \beta_1 x_i)(-1),$$

da cui:

$$n\beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n Y_i.$$

Stima dei parametri

Per quanto riguarda l'altra derivata:

$$\begin{aligned}\frac{\partial \ln L}{\partial \beta_1} &= -\frac{1}{2\sigma^2} \sum_{i=1}^n 2(Y_i - \beta_0 - \beta_1 x_i)(-x_i) \\ &= \frac{1}{\sigma^2} \sum_{i=1}^n (x_i Y_i - \beta_0 x_i - \beta_1 x_i^2),\end{aligned}$$

da cui:

$$\beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i.$$

Stima dei parametri

Quindi devo risolvere il sistema costituito dalle due seguenti equazioni:

$$n\beta_0 + \beta_1 \sum_{i=1}^n x_i = \sum_{i=1}^n Y_i$$

$$\beta_0 \sum_{i=1}^n x_i + \beta_1 \sum_{i=1}^n x_i^2 = \sum_{i=1}^n x_i Y_i$$

che dà come soluzione:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i Y_i - (\sum_{i=1}^n x_i) (\sum_{i=1}^n Y_i) / n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 / n}$$

e

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

dove \bar{Y} e \bar{x} sono, rispettivamente, la media dei valori Y_i e x_i .

Significato dei parametri della retta



Significato dei parametri della retta

Se abbiamo trovato per esempio una retta che rappresenta il modello di crescita di un bambino (in centimetri) in funzione della sua età espressa in mesi. La retta ha equazione:

$$y(x) = 50 + 0.753x$$

- Il coefficiente angolare indica la pendenza della retta. In questo esempio rappresenta di quanto aumenta l'altezza in funzione dei mesi di età. Per ogni mese aumenta di 0.753 centimetri.
- L'intercetta rappresenta l'altezza a 0 mesi, cioè alla nascita. In base ai dati a disposizione è stata stimata in 50 cm.
- Cosa significa calcolare $y(5)$? L'altezza stimata dal modello a 5 mesi, cioè $50 + 0.753 \cdot 5$.

Inferenza sui parametri

I parametri $\hat{\beta}_0$ e $\hat{\beta}_1$ sono anch'essi variabili aleatorie e dipendono dal campione considerato. Su di essi, pertanto, si possono fare le inferenze di tipo statistico che abbiamo visto nei paragrafi precedenti, quali stime di intervalli di confidenza e test di ipotesi.

Per il test di verifica di ipotesi, il parametro più importante è sicuramente $\hat{\beta}_1$ rispetto a $\hat{\beta}_0$. Per $\hat{\beta}_1$ il test di ipotesi più frequente è quello che verifica se $\hat{\beta}_1 \neq 0$, cioè se la retta di regressione è parallela all'asse x oppure no. Una retta di regressione parallela all'asse x significa che il valore della variabile aleatoria Y NON cambia al variare della variabile indipendente x . Quindi il test di ipotesi è formulato come:

- $H_0: \beta_1 = 0$
- $H_a: \beta_1 \neq 0$

Valori predetti

Il modello serve per predire il valore della variabile aleatoria in corrispondenza di un o più valori della variabile indipendente.

I valori predetti possono essere:

- **In sample.** I valori della variabile indipendente in cui si fa la predizione sono nell'insieme dei dati a disposizione. Posso confrontare le predizioni con i valori osservati nelle stesse ascisse.
- **Out of sample.** I valori della variabile indipendente NON sono nell'insieme dei dati a disposizione. Non ho quindi nessun termine di confronto per i valori che vengono predetti.

Il coefficiente R^2

È possibile avere un *singolo numero* che mi dà indicazioni sulla bontà del modello regressione lineare semplice rispetto al campione di dati a disposizione?

Il *coefficiente semplice di determinazione* viene calcolato appunto per questo scopo. Esso è definito dalla formula:

$$r^2 = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2},$$

dove, ricordiamo:

- $Y_i, i = 1, \dots, n$ sono i valori ‘osservati’ del campione;
- $\hat{Y}_i, i = 1, \dots, n$ sono i valori ‘fittati’, cioè i valori del modello di regressione lineare semplice in corrispondenza delle ascisse x_i ($\hat{Y}_i = \beta_0 + \beta_1 x_i, i = 1, \dots, n$);
- $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ è la media dei valori osservati.

Il coefficiente R^2

Si ha che $0 \leq r^2 \leq 1$. Tanto più r^2 è vicino a 1, tanto più il modello di regressione lineare è *buono*; tanto più r^2 è vicino a 0, tanto più il modello non è rappresentativo del campione dei dati. In quest'ultimo caso, l'analista cerca un modello differente da quello lineare per rappresentare i dati (una regressione non lineare o multivariata che coinvolga più di una variabile per esempio).

Associato al coefficiente semplice di determinazione r^2 si utilizza il *coefficiente semplice di correlazione* r che si ottiene come:

$$|r| = \sqrt{r^2}.$$

Per quanto riguarda il segno di r , si assume il segno della stima di β_1 calcolata.

Il coefficiente R^2

Il valore del coefficiente R^2 che indica un buon modello non si può definire a priori.

Dipende dalla disciplina, di solito nelle discipline scientifiche R^2 è maggiore rispetto alle discipline sociali.

In finanza e marketing, dipende da quali dati stiamo considerando.

Attenzione! Il coefficiente R^2 da solo non ha significato se non c'è effettivamente una dipendenza lineare fra i dati.

Librerie Python

Librerie Python per fare regressione lineare semplice:

Scipy.stats ([scipy.stats.linregress — SciPy v1.10.1 Manual](#))

Esempio file: *simul_linregress.py*

Statsmodel ([Introduction — statsmodels](#))-

Esempio file: *example_statsmodels.py*

Esempio file: *example_statsmodels_simul.py*

Scikit-learn ([sklearn.linear_model.LinearRegression — scikit-learn 1.2.2 documentation](#))

Esempio file: *example_sklearn.py*

Esempio file: *example_salary_sklearn.py*

Scipy.stats: esempio (file simul_linregress.py)

```
import matplotlib.pyplot as plt
```

```
from scipy import stats  
from numpy.random import randn  
from numpy.random import seed
```

```
seed(1)
```

```
x = randn(10)  
y = 1.6*x + randn(10)  
res = stats.linregress(x, y)  
print(f"R-squared: {res.rvalue**2:.6f}")
```

```
plt.plot(x, y, 'o', label='original data')  
plt.plot(x, res.intercept + res.slope*x, 'r', label='fitted line')  
plt.legend()  
plt.show()
```

Scikitlearn: esempio simulazione

```
from sklearn.linear_model import LinearRegression
```

```
# Generate sample data
```

```
x = np.array([1, 2, 3, 4, 5])
```

```
y = np.array([3, 5, 7, 9, 11])
```

```
# Reshape data
```

```
x = x.reshape(-1, 1)
```

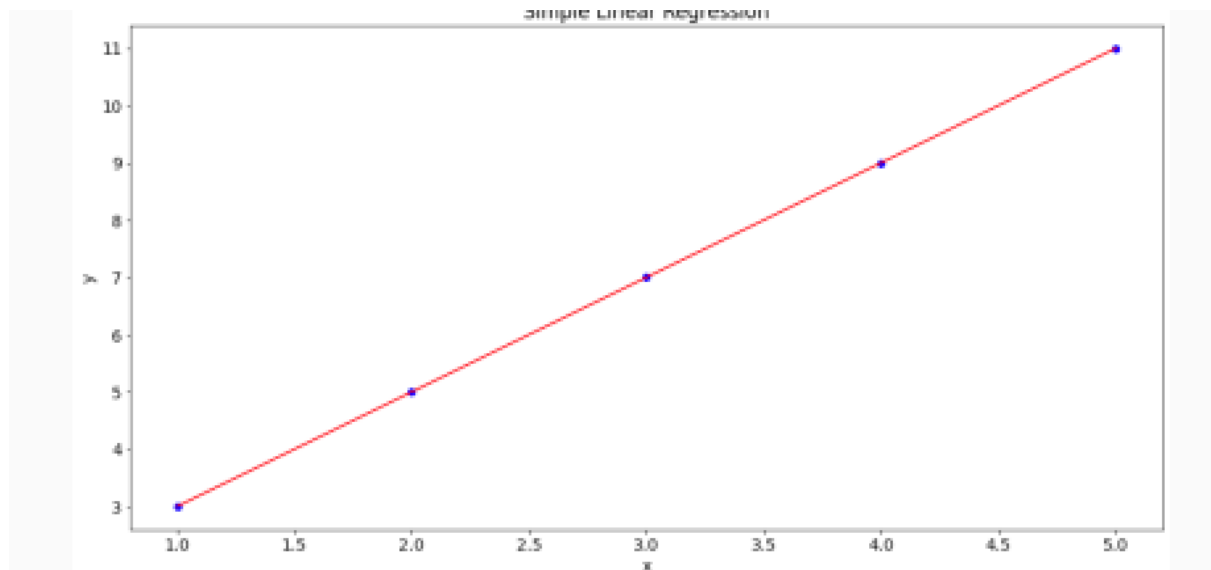
```
y = y.reshape(-1, 1)
```

```
# Create linear regression object and fit the model
```

```
reg = LinearRegression().fit(x, y)
```

```
# Predict the y-values using the trained model
```

```
y_pred = reg.predict(x)
```



Scikitlearn: esempio simulazione

Plot the data points and the linear regression line

```
plt.scatter(x, y, color='blue')  
plt.plot(x, y_pred, color='red')
```

```
# Add labels and a title to the plot  
plt.xlabel('x')  
plt.ylabel('y')  
plt.title('Simple Linear Regression')
```

```
# Display the plot  
plt.show()
```

Scikitlearn: esempio su data set

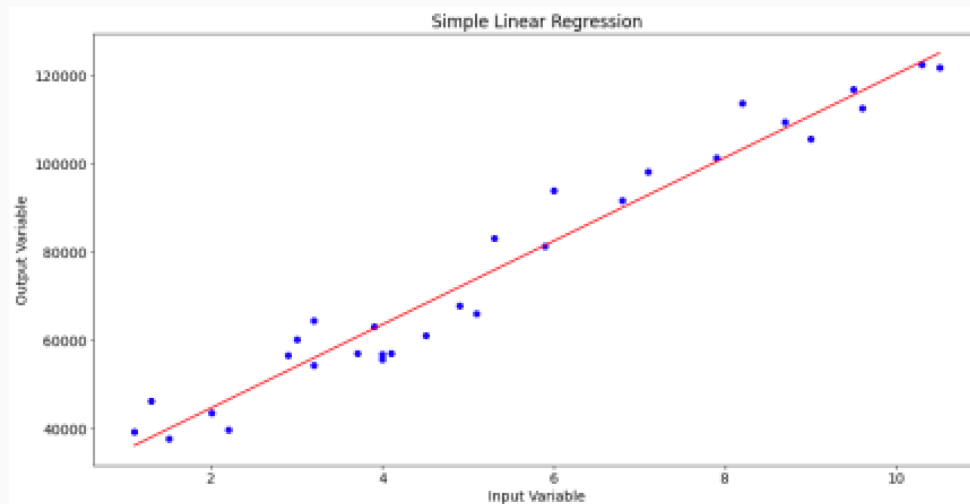
```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression

# Load data from Kaggle CSV file

#df = pd.read_csv('https://www.kaggle.com/vihansp/salary-
data-simple-linear-regression')
df = pd.read_csv("Salary_Data.csv")

# Extract input and output variables
x = df['YearsExperience'].values.reshape(-1, 1)

y = df['Salary'].values.reshape(-1, 1)
print(reg.intercept_, reg.coef_)
```



Scikitlearn: esempio su data set

```
# Create linear regression object and fit the model
```

```
reg = LinearRegression().fit(x, y)
```

```
# Predict the y-values using the trained model
```

```
y_pred = reg.predict(x)
```

```
# Plot the data points and the linear regression line
```

```
plt.scatter(x, y, color='blue')
```

```
plt.plot(x, y_pred, color='red')
```

```
# Add labels and a title to the plot
```

```
plt.xlabel('Input Variable')
```

```
plt.ylabel('Output Variable')
```

```
plt.title('Simple Linear Regression')
```


Analisi dei residui

I residui sono le differenze fra i valori osservati y_i e i valori predetti \hat{y}_i relativi alla stessa ascissa x_i .

I residui possono essere analizzati semplicemente graficamente per vedere il loro andamento, tramite *scatterplot* o *istogramma in frequenza*.

Funzioni Python: Esempio file [*example_residuals.py*](#)

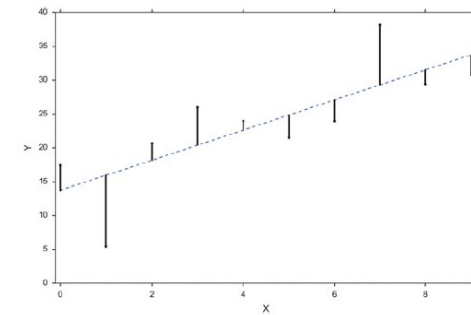
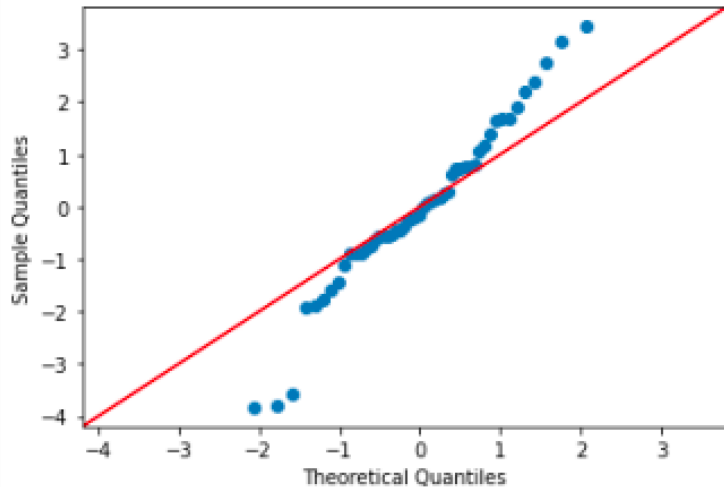


Fig. 11.3 Best-fit linear regression line (*dashed line*) and residuals (*solid lines*)



Python:

```
import statsmodels.api as sm  
fig = sm.qqplot(data, line='45')
```

Analisi dei residui

Sui residui viene fatta l'ipotesi di normalità, cioè si suppone, nel modello di regressione lineare, che i residui abbiano distribuzione normale con media 0 e deviazione standard σ

Non nota.

Per verificare l'ipotesi di normalità dei residui, posso fare il grafico QQ-plot.

Se le due distribuzioni sono simili, i punti devono stare molto vicini alla retta.

Test di ipotesi di normalità

Ci sono diversi test di ipotesi di normalità basati sul confronto della distribuzione stimata dei dati rispetto alla distribuzione normale.

Uno dei più famosi è il test di Shapiro-Wilk, che si basa sulla matrice di covarianza delle statistiche ordinate delle osservazioni e può essere utilizzato anche con un numero ridotto (< 50) di osservazioni.

H_0 : residui normali

H_a : residui non normali.

Test di ipotesi di normalità

Esempio: `from scipy.stats import shapiro`

```
gfg_data = randn(500)  
shapiro(gfg_data)
```

Output:

```
(0.9977102279663086, 0.7348126769065857)
```

Interpretazione dell'output: Poiché il p-value $0.73 > 0.5$ (livello di confidenza del test) **Non c'è evidenza x rigettare l'ipotesi nulla**, cioè per dire che i residui NON hanno distribuzione normale .

[scipy.stats.shapiro — SciPy v1.10.1 Manual](#)