# Fitting different models

```
In [26]:   1  import numpy as np
           2  import matplotlib.pyplot as plt
           3  np.set_printoptions(precision=3)
           4  #generate the data
           5  x = np.arange(100)
           6  y = 150 + 3*x + 0.03*x**2 + 5*np.random.randn(len(x))
           7
```

# line fit

```
In [27]:   1
           2  #Design the design matrices
           3  M1 = np.vstack( (np.ones_like(x), x) ).T
           4  #print(M1)
           5  # Solve the equations
           6  p1 = np.linalg.lstsq(M1, y)
           7  #estimated parameters
           8  print('The coefficients from the linear fit: {0}'.format(p1[0]))
           9  #line fit
          10  y_pred1=  p1[0][1] * x  + p1[0][0]
```

```
The coefficients from the linear fit: [102.542    5.948]

C:\Users\jyostna\Anaconda3\lib\site-packages\ipykernel_launcher.py:6: FutureWa
rning: `rcond` parameter will change to the default of machine precision times
``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=No
ne`, to keep using the old, explicitly pass `rcond=-1`.
```

# quadratic fit

In [28]:

```python
#Design the design matrices
M2 = np.vstack( (np.ones_like(x), x, x**2) ).T
#print(M1)
# Solve the equations
p2 = np.linalg.lstsq(M2, y)
#estimated parameters
print('The coefficients from the quadratic fit: {0}'.format(p2[0]))


#quadratic  fit
y_pred2= p2[0][0] + p2[0][1] * x + p2[0][2] * x *x   #quadratic fit

```

The coefficients from the quadratic fit: [1.522e+02 2.907e+00 3.072e-02]

C:\Users\jyostna\Anaconda3\lib\site-packages\ipykernel_launcher.py:5: FutureWa
rning: `rcond` parameter will change to the default of machine precision times
``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=No
ne`, to keep using the old, explicitly pass `rcond=-1`.
  """

## Cubic fit

In [29]:

```python
#Design the design matrices
M3 = np.vstack( (np.ones_like(x), x, x**2, x**3) ).T
# Solve the equations
p3 = np.linalg.lstsq(M3, y)
#estimated parameters
print('The coefficients from the quadratic fit: {0}'.format(p3[0]))
#cubic  fit
y_pred3= p3[0][0] + p3[0][1] * x + p3[0][2] * x *x + p3[0][3] * x * x * x    #c

```
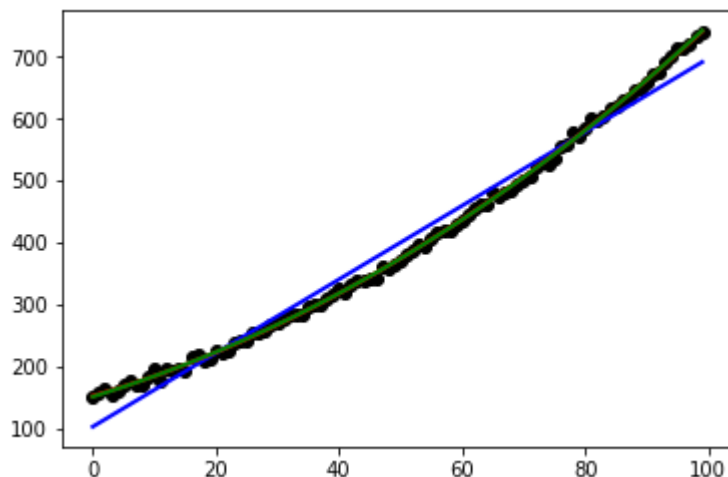
The coefficients from the quadratic fit: [1.511e+02 3.048e+00 2.714e-02 2.410e
-05]

C:\Users\jyostna\Anaconda3\lib\site-packages\ipykernel_launcher.py:4: FutureWa
rning: `rcond` parameter will change to the default of machine precision times
``max(M, N)`` where M and N are the input matrix dimensions.
To use the future default and silence this warning we advise to pass `rcond=No
ne`, to keep using the old, explicitly pass `rcond=-1`.
  after removing the cwd from sys.path.

```
In [30]:   1  plt.scatter(x, y,  color='black')
           2  plt.plot(x, y_pred1, color='blue', linewidth=2)
           3
           4  plt.plot(x, y_pred2, color='red' , linewidth=2 )
           5
           6  plt.plot(x, y_pred3, color='green' , linewidth=2 )
```

Out[30]: [<matplotlib.lines.Line2D at 0x2851db92f28>]



# Fit models USING STATS MODELS

```
In [36]:   1  import statsmodels.api as sm
           2  Res1 = sm.OLS(y, M1).fit()
           3  Res2 = sm.OLS(y, M2).fit()
           4  Res3 = sm.OLS(y, M3).fit()
           5  print(Res1.summary())
```

| | | | |
|---|---|---|---|
| Model: | OLS | Adj. R-squared: | 0.982 |
| Method: | Least Squares | F-statistic: | 5264. |
| Date: | Tue, 16 Apr 2019 | Prob (F-statistic): | 5.52e-87 |
| Time: | 08:59:13 | Log-Likelihood: | -457.29 |
| No. Observations: | 100 | AIC: | 918.6 |
| Df Residuals: | 98 | BIC: | 923.8 |
| Df Model: | 1 | | |
| Covariance Type: | nonrobust | | |

```
====================================================================
              coef    std err          t      P>|t|      [0.025      0.975]
--------------------------------------------------------------------
```

# the AIC-value,

The Akaike Information Criterion, which can be used to assess the quality of the model: the lower the AIC value, the better the model. We see that the quadratic model has the lowest AIC value and therefore is the best model: it provides the same quality of fit as the cubic model, but uses fewer parameters to achieve that quality.

```
In [30]:  1  print('The AIC-value is \n {0:4.1f} for the linear fit,\n {1:4.1f} for the qua
          2  Res3.aic))
```

```
The AIC-value is
 909.5 for the linear fit,
 647.4 for the quadratic fit, and
 648.8 for the cubic fit
```

# The degrees of freedom (Df )

df of the model are the number of predictors, or explanatory, variables.

The Df of the residuals is the number of observations minus the degrees of freedom of the model, minus one (for the offset).

```
In [40]:  1  print( Res1.df_model, Res2.df_model, Res3.df_model)
          2
```

```
1.0 2.0 3.0
```

# R-Squared

$R^2$ value indicates the proportion of variation in the y-variable that is due to variation in the x-variables.

For simple linear regression, the $R^2$ value is the square of the sample correlation $r\_xy$.

For multiple linear regression with intercept (which includes simple linear regression), the $R^2$ value is defined as

$R^2$ = SS_mod / SS_tot

$R^2$ is close to 1 means the model is better

```
In [41]:  1  print( Res1.rsquared, Res2.rsquared, Res3.rsquared)
          2
```

```
0.9835452518838422  0.9988266561367175  0.9988334443953123
```

# Definitions for Regression with Intercept

The Sum-of Squares variables SSxx

n is the number of observations,

and k is the number of regression parameters.

$DF\_mod = k - 1$ is the (Corrected) Model Degrees of Freedom.

$DF\_res = n - k$ is the Residuals Degrees of Freedom $DF\_tot = n - 1$ is the (Corrected) Total Degrees of Freedom.

For multiple regression models with intercept, $DF\_mod + DF\_res = DF\_tot$.

$MS\_mod = SS\_mod/DFmod$ : Model Mean of Squares

$MS\_res = SS\_res/DF\_res$ : Residuals Mean of Squares.

$MS\_tot = SS\_tot/DF\_tot$ : Total Mean of Squares, which is the sample variance of the y-variable.

In [ ]:  | 1