

# Curve Fitting

Introduction to Numerical Problem Solving, Spring 2017

CC BY-NC-SA Sakari Lukkarinen

Helsinki University of Applied Sciences

# Motivation

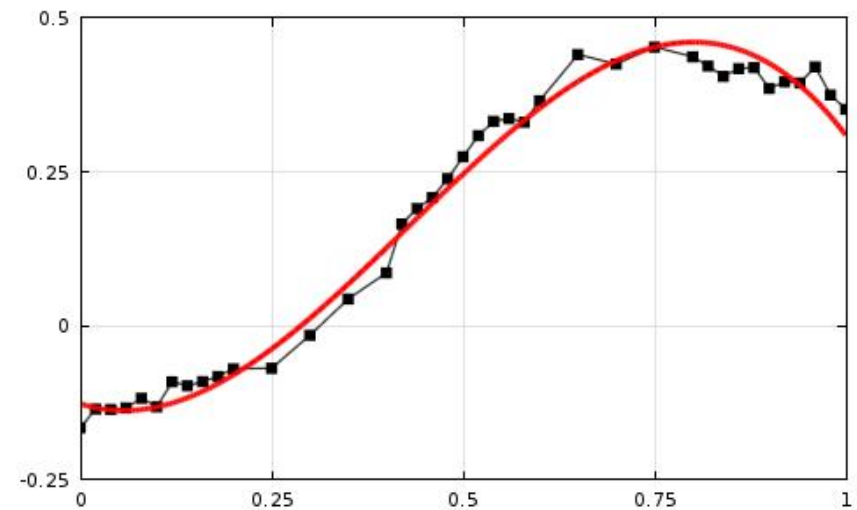
Given the  $n + 1$  data points  $(x_i, y_i)$   $i = 0, 1, 2, \dots, n$   
estimate  $y(x)$ .

# Motivation

If you have discrete data points (black dots) and you want to know the values between, you have two options:

Curve fitting – you make a mathematical model that models the data as good as possible (red line)

Interpolation – you interpolate the values between the points using some algorithm (black line)



# Example

(a) Least-squares linear regressions

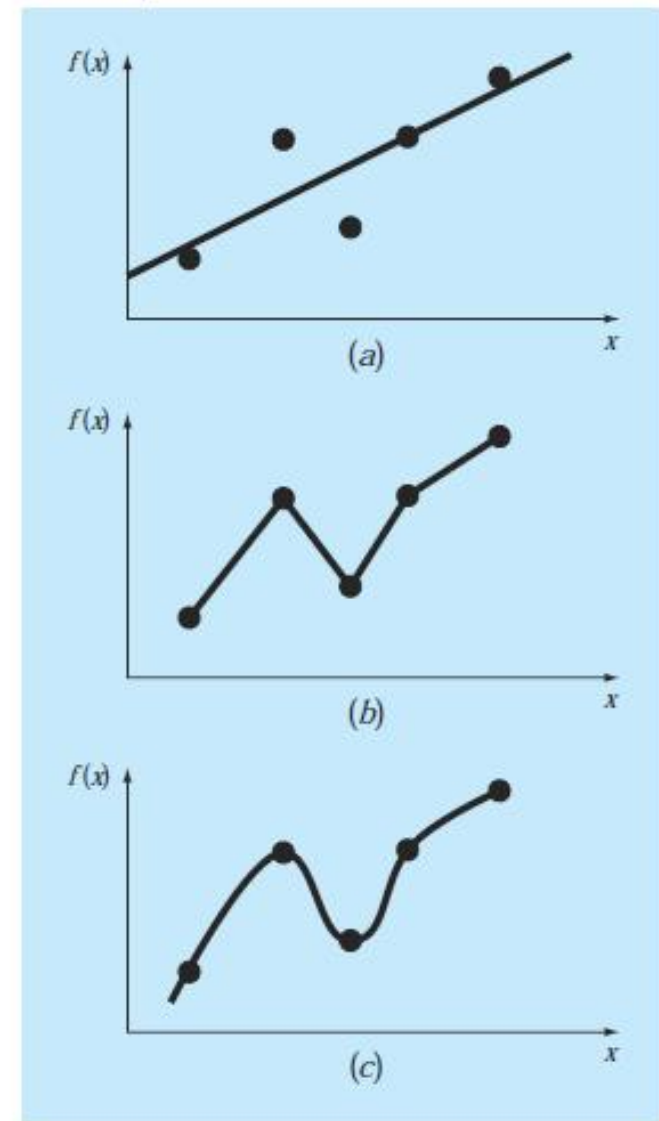
$$y = ax + b$$

(b) Linear interpolation

$$y_1 = a_1x + b_1, y_2 = a_2x + b_2, \dots$$

(c) Curvilinear interpolation, for example second order polynomials

$$a_1x^2 + b_2x + c, \dots$$



# Background - Simple statistics

- Mean
- Total sum of the squared residuals
- Standard deviation
- Variance
- Coefficient of variation

$$\bar{y} = \frac{\sum y_i}{n}$$

$$S_t = \sum (y_i - \bar{y})^2$$

$$s_y = \sqrt{\frac{S_t}{n-1}}$$

$$s_y^2 = \frac{\sum (y_i - \bar{y})^2}{n-1}$$

$$\text{c.v.} = \frac{s_y}{\bar{y}} 100\%$$

# Normal distribution and confidence interval(s)

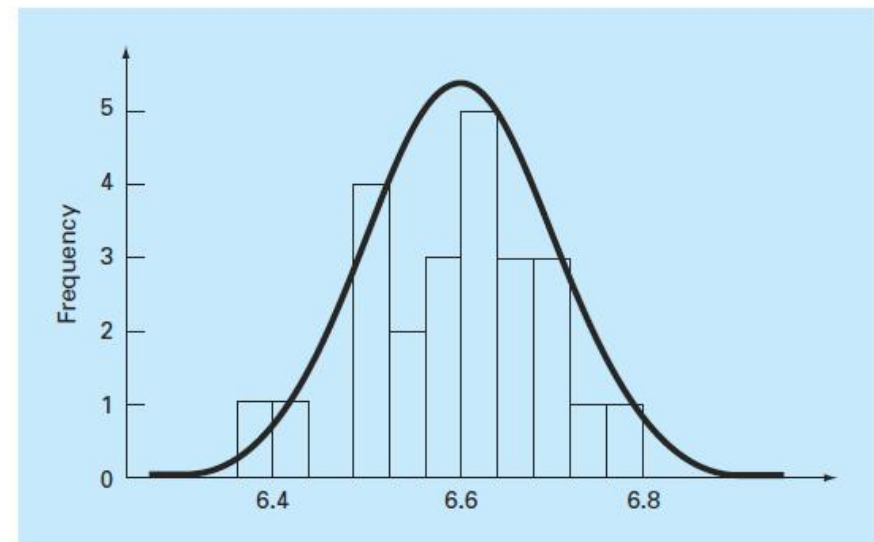
Mean value = 6.6

Standard deviation = 0.1

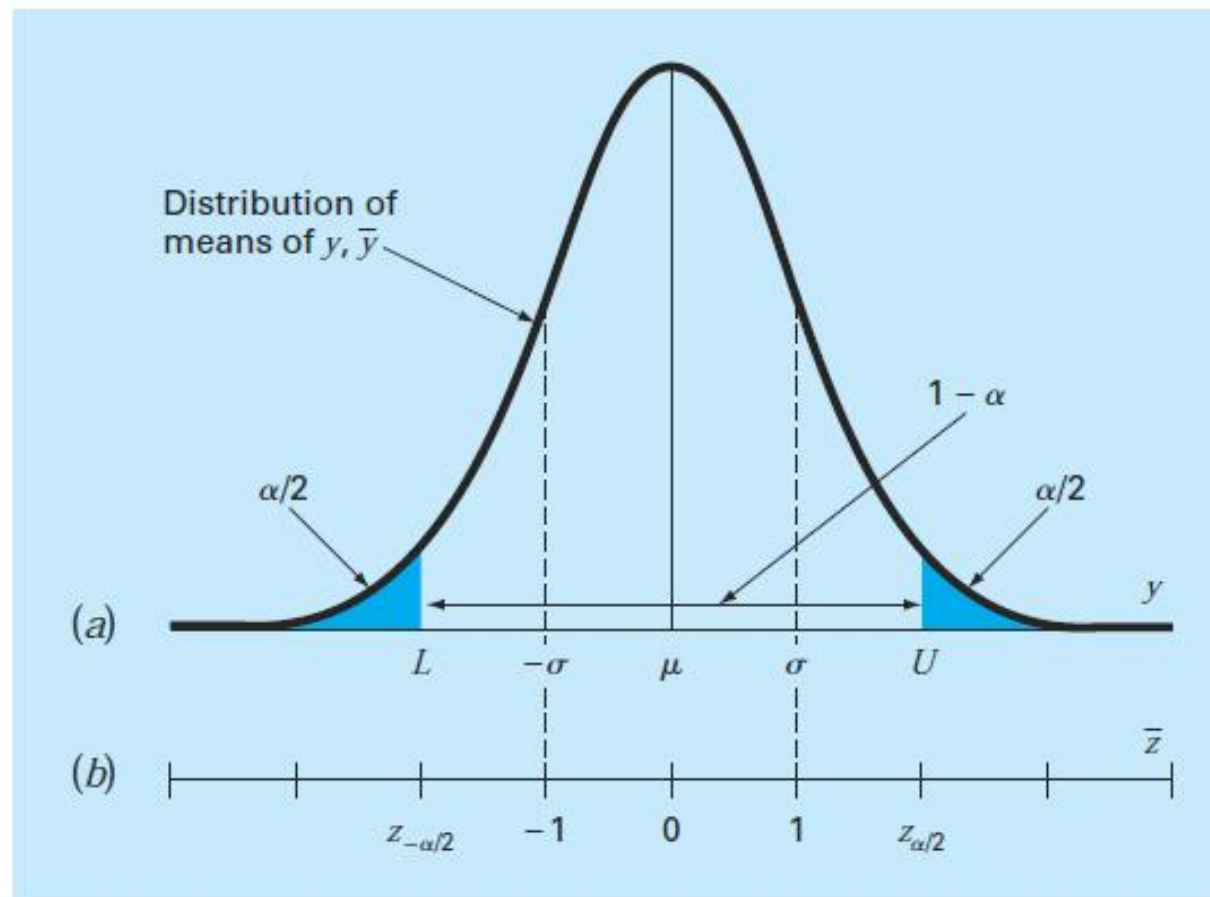
68 % of the values between:  
 $m - 1*s.d. \dots m + 1*s.d.$

95 % of the values between:  
 $m - 1.96*s.d \dots m + 1.96*s.d.$

99 % of the values between:  
 $m - 3.0*s.d \dots m + 3.0*s.d.$



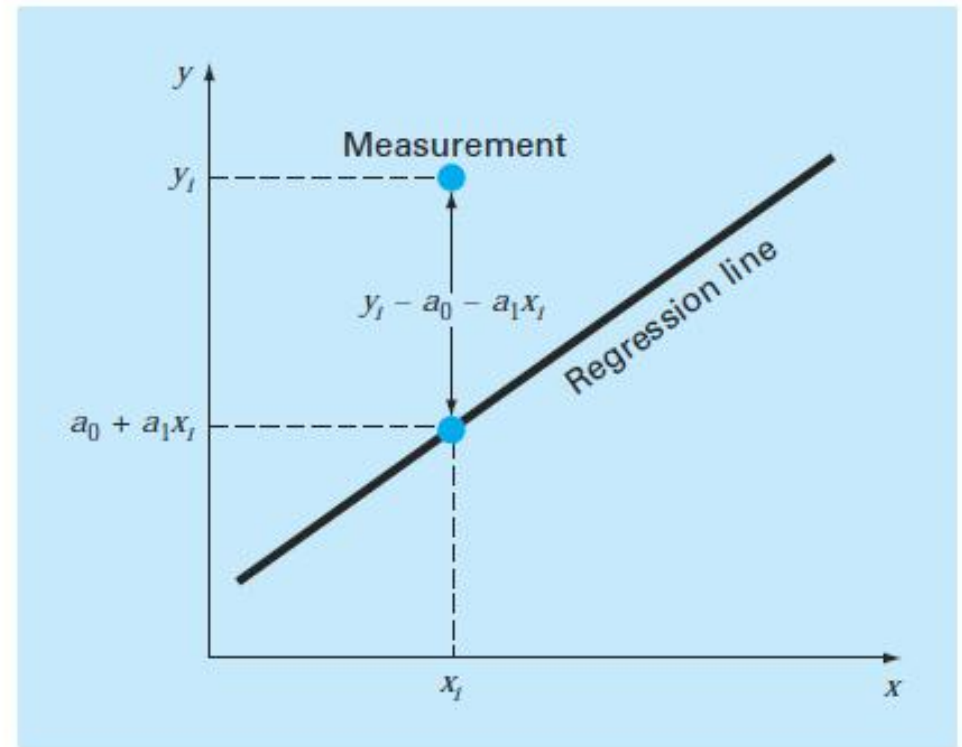
# Two-sided confidence interval



# Residuals and linear model

The residual represent the vertical distance between the data point and the linear model. The coefficients can be found by minimizing the total error.

$$y = a_0 + a_1 x + e$$





# What is the best criteria for fitting a linear model?

Sum of the residual of the errors:

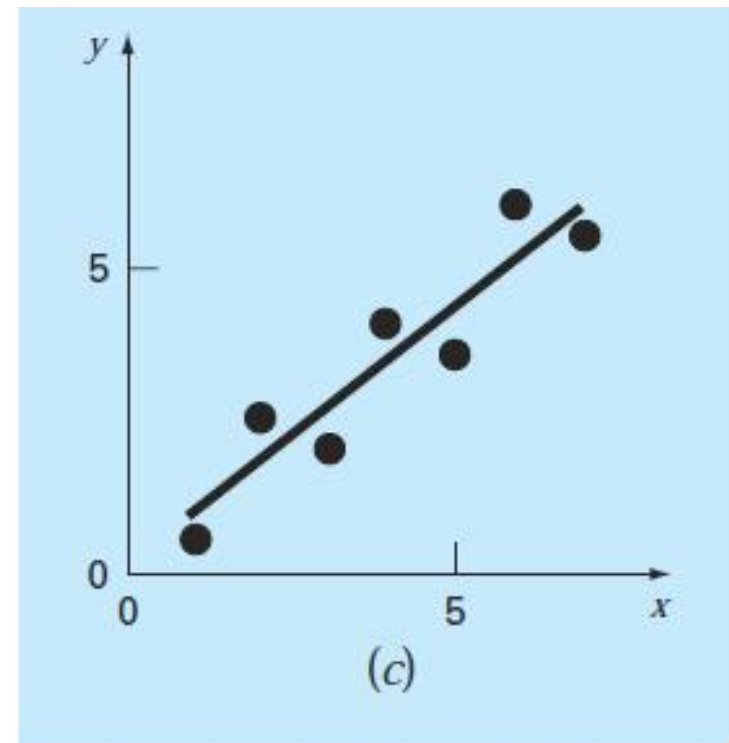
$$\sum_{i=1}^n e_i = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)$$

Sum of the absolute of the errors:

$$\sum_{i=1}^n |e_i| = \sum_{i=1}^n |y_i - a_0 - a_1 x_i|$$

Sum of the square of the errors:

$$\sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$



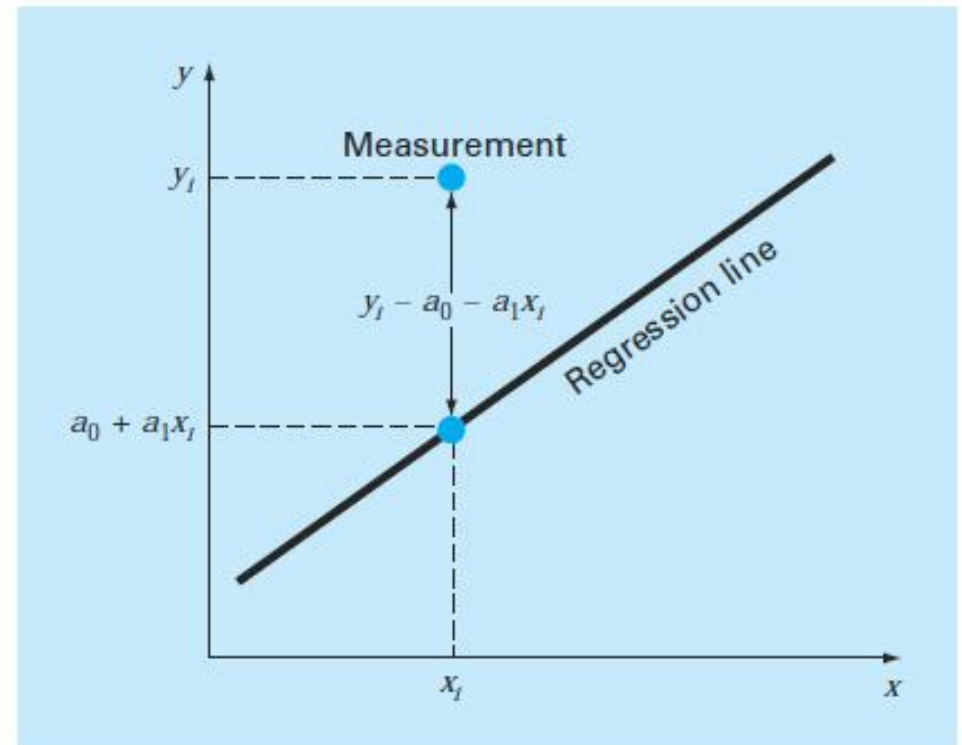
# Linear regression

If we use the sum of the squared residuals as the error criteria, the coefficients for the linear model can be found by minimizing the error.

$$a_1 = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

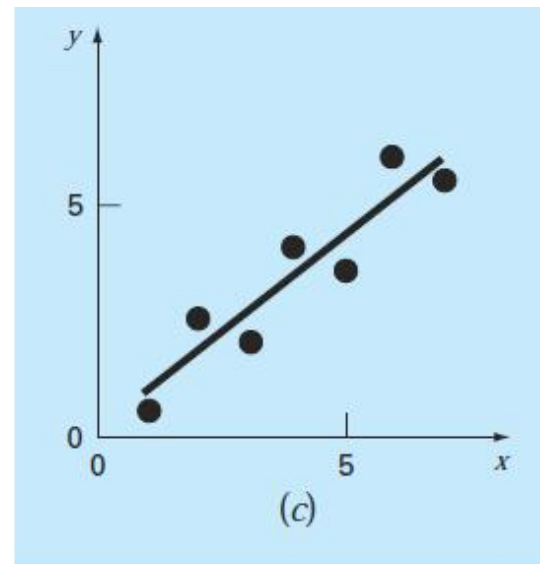


# Standard error of the estimate

Standard error of the estimate e.g. the spread around the regression line can be calculated from the sum of the squared residuals.

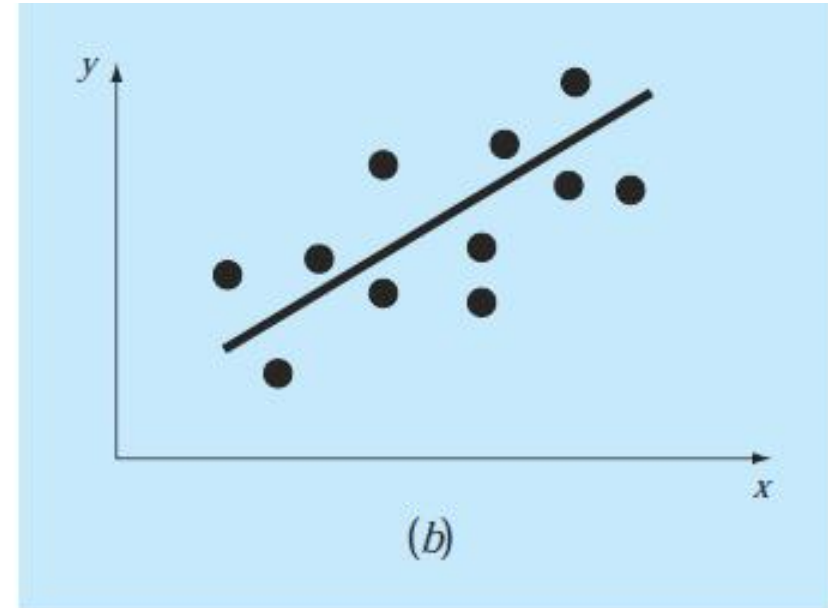
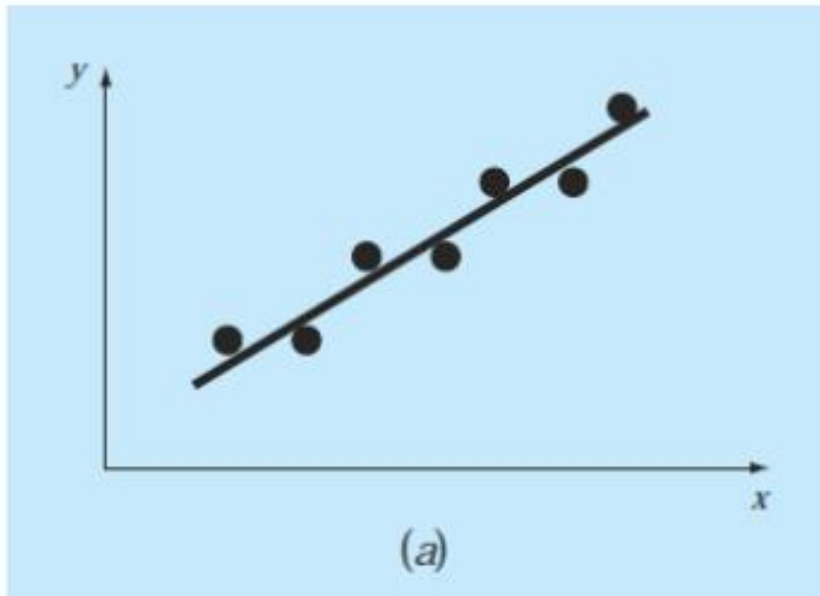
$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

$$s_{y/x} = \sqrt{\frac{S_r}{n-2}}$$



Small (a) and large (b)  
standard error of the estimate

$$s_{y/x} = \sqrt{\frac{S_r}{n-2}}$$



$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

# Goodness of the fit

(a) Total sum of the squares around the mean

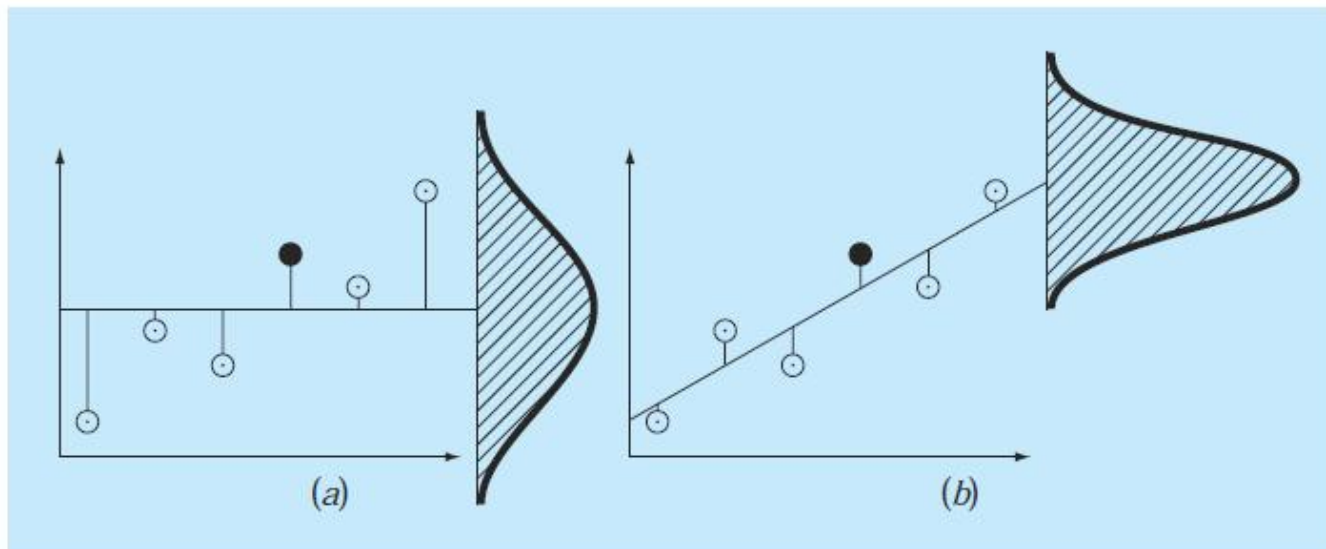
$$S_t = \sum (y_i - \bar{y})^2$$

(b) Total sum of the squares of residuals

$$S_r = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i)^2$$

Coefficient of determination

$$r^2 = \frac{S_t - S_r}{S_t}$$



## Statistics

# Numpy functions

### Order statistics

---

<code>amin(a[, axis, out, keepdims])</code>	Return the minimum of an array or minimum along an axis.
<code>amax(a[, axis, out, keepdims])</code>	Return the maximum of an array or maximum along an axis.
<code>nanmin(a[, axis, out, keepdims])</code>	Return minimum of an array or minimum along an axis, ignoring any NaNs.
<code>nanmax(a[, axis, out, keepdims])</code>	Return the maximum of an array or maximum along an axis, ignoring any NaNs.
<code>ptp(a[, axis, out])</code>	Range of values (maximum - minimum) along an axis.
<code>percentile(a, q[, axis, out, ...])</code>	Compute the qth percentile of the data along the specified axis.
<code>nanpercentile(a, q[, axis, out, ...])</code>	Compute the qth percentile of the data along the specified axis, while ignoring nan values.

### Averages and variances

---

<code>median(a[, axis, out, overwrite_input, keepdims])</code>	Compute the median along the specified axis.
<code>average(a[, axis, weights, returned])</code>	Compute the weighted average along the specified axis.
<code>mean(a[, axis, dtype, out, keepdims])</code>	Compute the arithmetic mean along the specified axis.
<code>std(a[, axis, dtype, out, ddof, keepdims])</code>	Compute the standard deviation along the specified axis.
<code>var(a[, axis, dtype, out, ddof, keepdims])</code>	Compute the variance along the specified axis.
<code>nanmedian(a[, axis, out, overwrite_input, ...])</code>	Compute the median along the specified axis, while ignoring NaNs.
<code>nanmean(a[, axis, dtype, out, keepdims])</code>	Compute the arithmetic mean along the specified axis, ignoring NaNs.
<code>nanstd(a[, axis, dtype, out, ddof, keepdims])</code>	Compute the standard deviation along the specified axis, while ignoring NaNs.
<code>nanvar(a[, axis, dtype, out, ddof, keepdims])</code>	Compute the variance along the specified axis, while ignoring NaNs.

### Correlating

---

<code>corrcoef(x[, y, rowvar, bias, ddof])</code>	Return Pearson product-moment correlation coefficients.
<code>correlate(a, v[, mode])</code>	Cross-correlation of two 1-dimensional sequences.
<code>cov(m[, y, rowvar, bias, ddof, fweights, ...])</code>	Estimate a covariance matrix, given data and weights.

## scipy.stats.linregress

**scipy.stats.linregress**(*x*, *y=None*)

[\[source\]](#)

Calculate a regression line

This computes a least-squares regression for two sets of measurements.

**Parameters:** *x, y* : *array\_like*

two sets of measurements. Both arrays should have the same length. If only *x* is given (and *y=None*), then it must be a two-dimensional array where one dimension has length 2. The two sets of measurements are then found by splitting the array along the length-2 dimension.

**Returns:**

**slope** : *float*

slope of the regression line

**intercept** : *float*

intercept of the regression line

**r-value** : *float*

correlation coefficient

**p-value** : *float*

two-sided p-value for a hypothesis test whose null hypothesis is that the slope is zero.

**stderr** : *float*

Standard error of the estimate

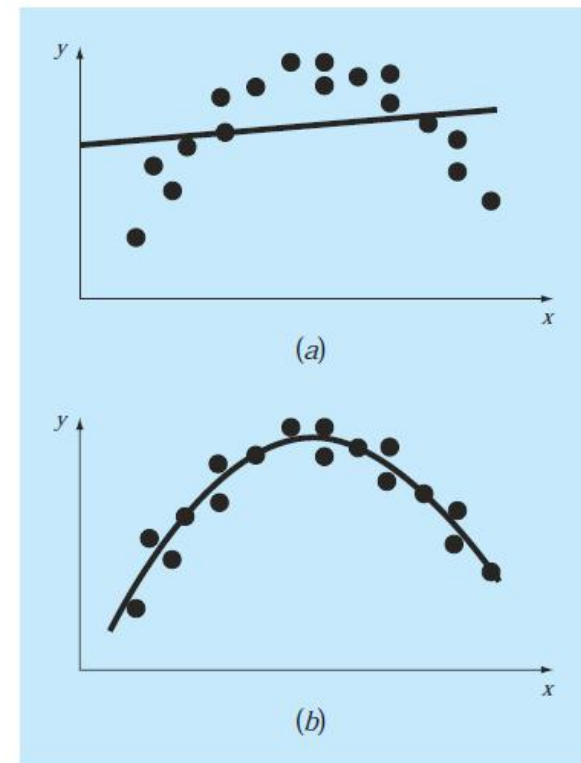
<https://docs.scipy.org/doc/scipy-0.14.0/reference/generated/scipy.stats.linregress.html>

# Nonlinear relationship

Example of data where linear ( $ax+b$ ) model doesn't work, but a parabolic model ( $ax^2 + bx + c$ ) would work better.

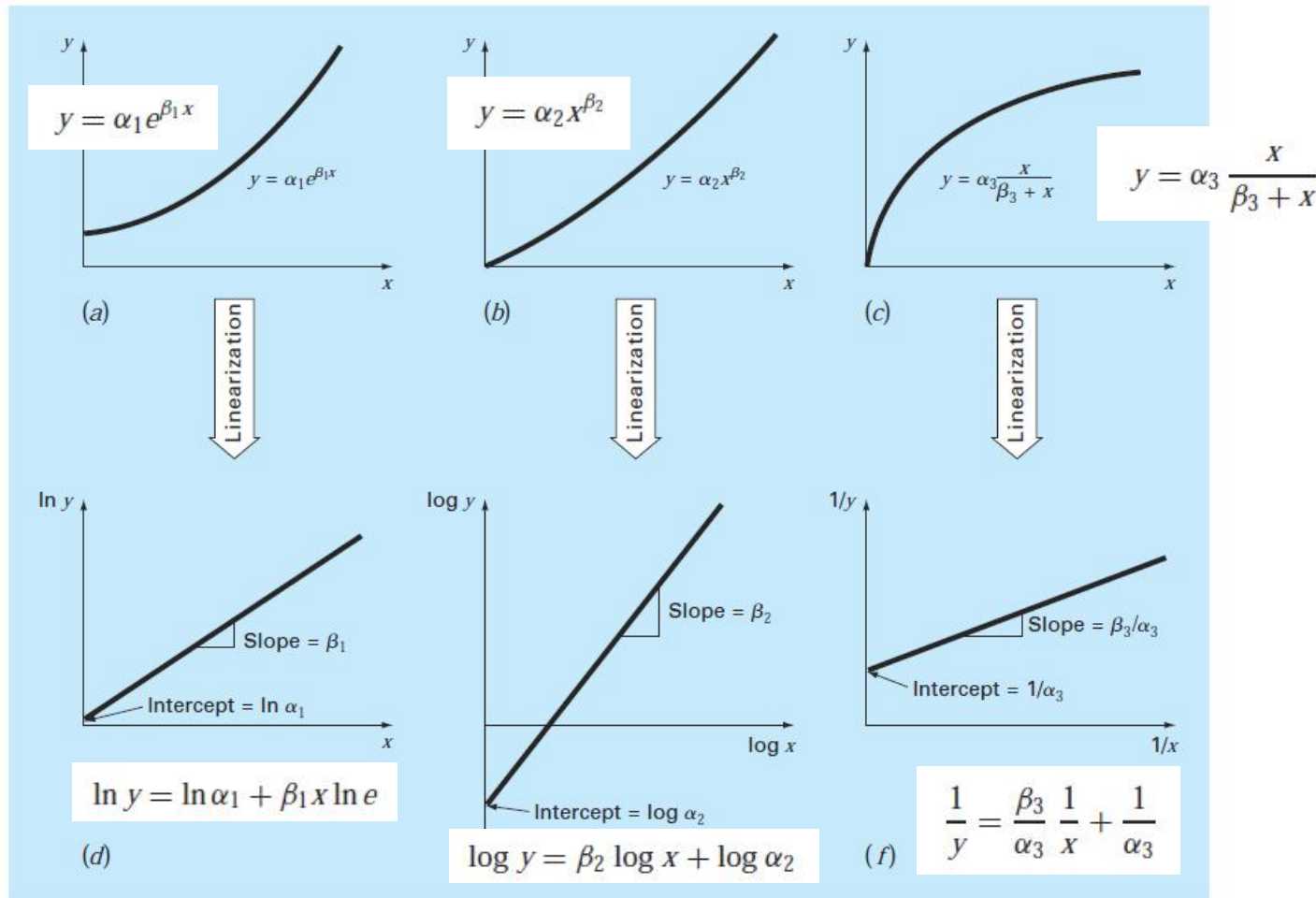
Solutions:

- Linearize the model and use linear regression
- Use nonlinear curve fitting





# Linearization of nonlinear relationship



## scipy.optimize.curve\_fit

[https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.curve\\_fit.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.curve_fit.html)

**scipy.optimize.curve\_fit**(*f*, *xdata*, *ydata*, *p0*=None, *sigma*=None, *absolute\_sigma*=False, *check\_finite*=True, *bounds*=(-inf, inf), *method*=None, *jac*=None, *\*\*kwargs*)

[\[source\]](#)

Use non-linear least squares to fit a function, *f*, to data.

Assumes  $ydata = f(xdata, *params) + \epsilon$

### Parameters:

**f** : callable

The model function,  $f(x, \dots)$ . It must take the independent variable as the first argument and the parameters to fit as separate remaining arguments.

**xdata** : An *M*-length sequence or an  $(k, M)$ -shaped array for functions with *k* predictors

The independent variable where the data is measured.

**ydata** : *M*-length sequence

The dependent data — nominally  $f(xdata, \dots)$

**p0** : None, scalar, or *N*-length sequence, optional

Initial guess for the parameters. If None, then the initial values will all be 1 (if the number of parameters for the function can be determined using introspection, otherwise a ValueError is raised).

**sigma** : None or *M*-length sequence or  $M \times M$  array, optional

Determines the uncertainty in *ydata*. If we define residuals as  $r = ydata - f(xdata, *popt)$ , then the interpretation of *sigma* depends on its number of dimensions:

- A 1-d *sigma* should contain values of standard deviations of errors in *ydata*. In this case, the optimized function is  $chisq = \sum((r / sigma) ** 2)$ .
- A 2-d *sigma* should contain the covariance matrix of errors in *ydata*. In this case, the optimized function is  $chisq = r.T @ inv(sigma) @ r$ .  
New in version 0.19.

None (default) is equivalent of 1-d *sigma* filled with ones.

## scipy.optimize.least\_squares

**scipy.optimize.least\_squares**(*fun*, *x0*, *jac*='2-point', *bounds*=(-inf, inf), *method*='trf', *ftol*=1e-08, *xtol*=1e-08, *gtol*=1e-08, *x\_scale*=1.0, *loss*='linear', *f\_scale*=1.0, *diff\_step*=None, *tr\_solver*=None, *tr\_options*={}, *jac\_sparsity*=None, *max\_nfev*=None, *verbose*=0, *args*=(), *kwargs*={}) [\[source\]](#)

Solve a nonlinear least-squares problem with bounds on the variables.

Given the residuals  $f(x)$  (an  $m$ -dimensional real function of  $n$  real variables) and the loss function  $\rho(s)$  (a scalar function), `least_squares` finds a local minimum of the cost function  $F(x)$ :

```
minimize  $F(x) = 0.5 * \sum(\rho(f_i(x))^2)$ ,  $i = 0, \dots, m - 1$   
subject to  $lb \leq x \leq ub$ 
```

The purpose of the loss function  $\rho(s)$  is to reduce the influence of outliers on the solution.

**Parameters:** *fun* : callable

Function which computes the vector of residuals, with the signature `fun(x, *args, **kwargs)`, i.e., the minimization proceeds with respect to its first argument. The argument  $x$  passed to this function is an ndarray of shape  $(n,)$  (never a scalar, even for  $n=1$ ). It must return a 1-d array\_like of shape  $(m,)$  or a scalar. If the argument  $x$  is complex or the function `fun` returns complex residuals, it must be wrapped in a real function of real arguments, as shown at the end of the Examples section.

**x0** : array\_like with shape  $(n,)$  or float

Initial guess on independent variables. If float, it will be treated as a 1-d array with one element.

[https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.least\\_squares.html](https://docs.scipy.org/doc/scipy/reference/generated/scipy.optimize.least_squares.html)