# Parameter estimation & optimisation
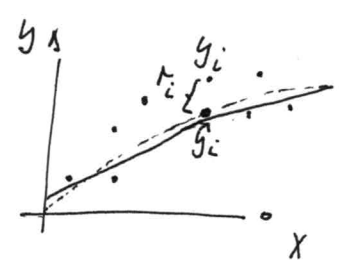
- estimation of parameters appearing in models / functions
- efficient use of data for mathematical modelling
- simplest model linear regression



- let $\vec{p}$ be the parameter vector ~~containing the par of~~
  parameters to be optimised

  $m$ :is number of parameters $p_1, \ldots, p_m$

Intuition: $y_i = p_0 + p_1 x_i + \varepsilon_i \quad \varepsilon_i \in \mathcal{N}(0, \sigma^2)$

$p_0, p_1$ so wählen, dass $y(x) = \hat{p}_0 + \hat{p}_1 x$
möglichst nah an den Daten liegt
d.h. Abstand minimieren

- let $\vec{z}$ be the ~~measurement~~ vector which is the "ideal" output of the system to
  be modelled

  the system in the noise-free case is described by a vector function $\vec{f}$
  which relates $\vec{z}$ to $\vec{p}$ ~~such~~ such that

$$\vec{f}(\vec{p}, \vec{z}) = 0$$

- In practice measurements $\vec{y}$ are only available for system output $\vec{z}$ with noise

$$\vec{y} = \vec{z} + \vec{\varepsilon}$$

- We take multiple measurements of the system $\{y_i\} \; i = 1, \ldots, n$
  and want to estimate $\vec{p}$ using $\{y_i\}$

  $\rightarrow$ due to noise $\vec{f}(\vec{p}, y_i) = 0$ is not valid anymore

~~the~~ Solution: We write cost function or objective function $F$
describing the error between measurement and system output for
given parameters

$$F(\vec{p}, y_1, \ldots, y_n)$$

and minimize the cost

if there are no constraints on $\vec{p}$ and function $F$ has first and
second order partial derivatives, necessary conditions for a minimum are

$$\frac{\partial F}{\partial \vec{p}} = 0 \qquad \text{and} \qquad \frac{\partial^2 F}{\partial \vec{p}^2} > 0$$

## Least-square optimisation

minimise the error of sum of squares

$$\min \; F = \sum_{i=1}^{n} \left( y_i - \hat{y}_i \right)^2 \qquad \text{with} \quad \hat{y}_i = \vec{f}(x_i, \vec{p}) \; \text{model predictions}$$

$$F(\vec{p}) = \sum (y_i - f(x_i, \vec{p}))^2 \qquad \text{for specific } x_i \text{ using the parameters}$$

if the _errors are normally distributed_
the least square estimates are also the _maximum likelihood estimates_

## How to find $\vec{\hat{p}}$ which minimizes $F$:

1. Finding an analytical solution by
   - differentiating $F$ with respect to $p_1, ..., p_m$
   - setting partial derivatives zero $\quad \dfrac{\partial F}{\partial p_i} = 0$
   - solving the resulting $m$ normal equations
   - only working for very few nonlinear models

2. Numerical solutions
   - try different values for parameters $\vec{p}_g$
   - calculate $F(\vec{p}_g)$ and work towards smaller $F$;

   3 main procedures (sensitivity based approaches)

   - _Steepest descent method (gradient descent)_

     

     Searches minimum $F$ by iteratively determining the direction in which the parameter estimate should change
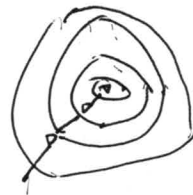
     $F(\vec{p})$ is defined and differentiable
     $F(\vec{p})$ decreases fastest in point $\vec{p} = \vec{a}$ if one goes
     from $\vec{a}$ in the direction of negative gradient of $F$ at $\vec{a}$

     $$- \nabla F(\vec{a})$$

     $$\vec{a}_{k+1} = \vec{a}_k - \int \nabla F(a_k) \qquad \int \text{ is the step size}$$
     $$\text{(allowed to change)}$$

     _Gauss-Newton method (GNA)_ or _linearization_
     - uses taylor series expansion to approximate the nonlinear model with
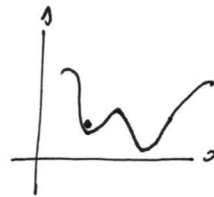       linear terms.
       terms are used in linear regression to come up with New terms

     (_Levenberg_-Marquardt algorithm or _damped least squares_
     interpolates between Gauss Newton & steepest descent algorithm)

   _Starting values:_
   - all iterative procedures require starting values
     - risk of local minima (multiple start methods, particle swarm)

     

## Important

- the simpler the model, the better the behavior in the estimation process.

- Over parameterization often leads to convergence problems
  - may have multiple solutions
  - high correlation between parameter estimates

Parametrization can have large influences

$$y_i = \frac{\beta_0 x_i}{\beta_1 x_i + \beta_2} + \varepsilon_i \qquad vs. \qquad y_i = \frac{x_i}{c_0 x_i + c_1}$$
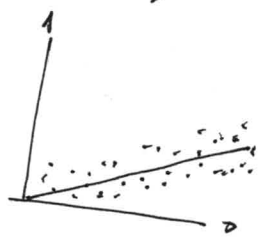
## Practical example:

linear regression

$$y_i = \overbrace{\beta_1 x_i}^{t_i} + \varepsilon_i \qquad \varepsilon_i \sim N(0, \sigma^2)$$

$$F(\beta_1) = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$\min \sum_{i=1}^{n} (y_i - \beta_1 x_i)^2$$



$$\frac{\partial F}{\partial \beta_1} = \sum_{i=1}^{n} 2(y_i - \beta_1 x_i) \cdot (-x_i) \quad * \quad \psi$$

$$\sum_{i=1}^{n} \left(-2 x_i y_i + 2 \beta_1 x_i^2\right) = 0$$

$$\sum x_i y_i = \beta_1 \sum x_i^2$$

$$\beta_1 = \frac{\sum x_i y_i}{\sum x_i^2}$$

$$\frac{\partial^2 F}{\partial^2 \beta} = 2 \sum_i x_i^2 > 0$$

$$\beta_0 = \langle \bar{y} \rangle - b_1 \langle x \rangle$$

$$\beta_1 = \frac{\sum_{i=1}^{n} (x_i - \langle x \rangle)(y_i - \langle y \rangle)}{\sum_{i=1}^{n} (x_i - \langle x \rangle)^2}$$

$$= \frac{cov(x, y)}{var(x)}$$

From where come the ~~derivations~~ derivatives from (for sensitivity based methods)

$$\frac{\partial F}{\partial p_i} =$$

- often finite differences, but very numerical intensive

$$\frac{\partial F}{\partial p_i}\bigg|_{\vec{a}} \sim \frac{F(\vec{a}) - F(\vec{a} + \Delta p_i)}{p_i - \Delta p_i}$$

- based on sensitivity equations
  analytical; e.g. use ODE's + $F$ and perform derivatives (symbolic math)

· <u>Stochastic approaches:</u>
  · simulated annealing
  · particle swarm — population of candidate solutions
· <u>Bayes approaches:</u> posterior distributions
  maximum likelihood
  + Gibbs sampling & MCMC sampling
· <u>Identifiability</u>: Profile likelihood
· Local vs global maximum: <u>Waterfall plots</u>