

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

# Lecture 10: Least Squares

---

---

---

---

---

---

---

---

---

---

---

---

---

---

---

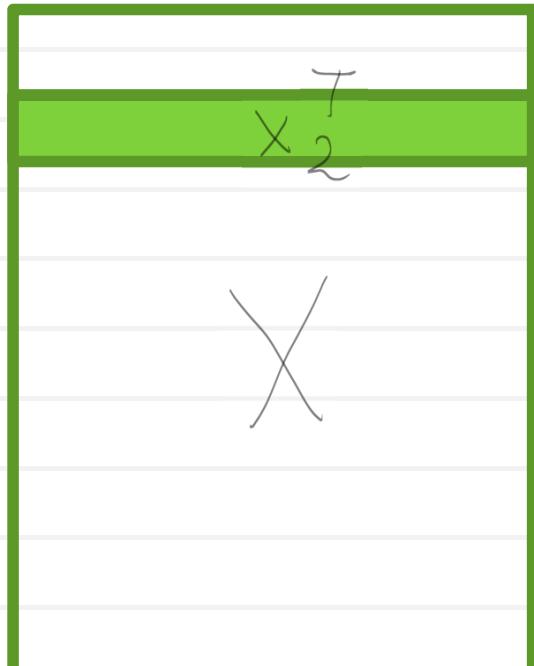
---

---

# Linear least squares setting

Least squares the most widely used data fitting model in statistics, science, engineering, etc.

Data matrix (aka observation matrix):



$$X = [x_1 \ x_2 \dots x_n]^T$$

$$x_i \in \mathbb{R}^m$$

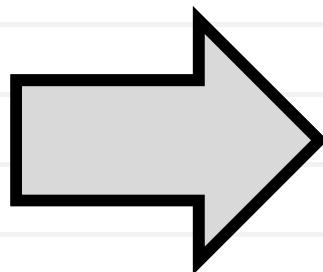
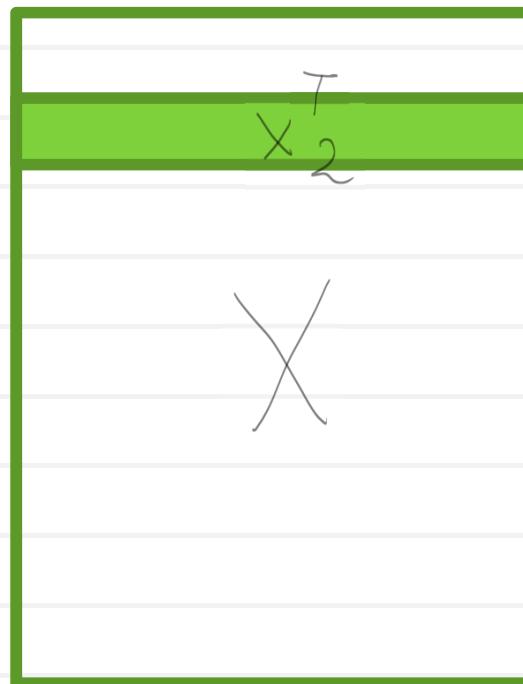
a vector of input variables  
*(observations, independent variables,  
features, predictors, regressors)*

m

# Linear least squares setting

$$X = [x_1 \ x_2 \dots x_n]^T$$

$$y = [y_1 \ y_2 \dots y_n]^T \in \mathbb{R}^n$$

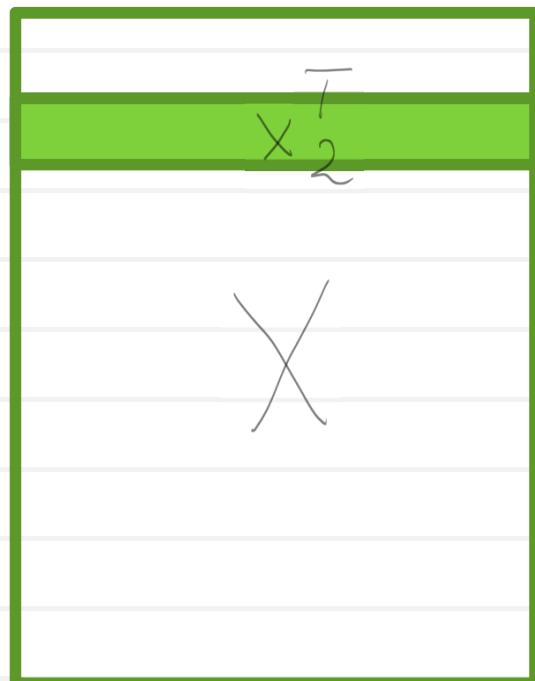


*Outcomes (dependent variables, response)*

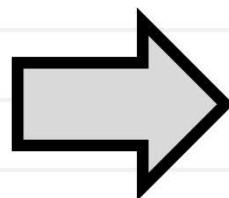
# Linear least squares setting

$$X = [x_1 \ x_2 \dots x_n]^T$$

$$y = [y_1 \ y_2 \dots y_n]^T \in \mathbb{R}^n$$



$n^*$



$n$

Model to be recovered by optimization.

$$x_1^T w = y_1$$

$$x_2^T w = y_2$$

$$\dots$$

$$x_n^T w = y_n$$

$X \ w = y + \text{"noise"}$

# Linear least squares setting

$$X \omega = y + \text{"noise"}$$

$\omega$  that gives the best fit (in the *least-quadratic sense*) is then recovered by:

$$\min_{\omega} \frac{1}{2} \| X\omega - y \|_2^2 = F(\omega)$$

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (x_i^\top \cdot \omega - y_i)^2$$

Residual

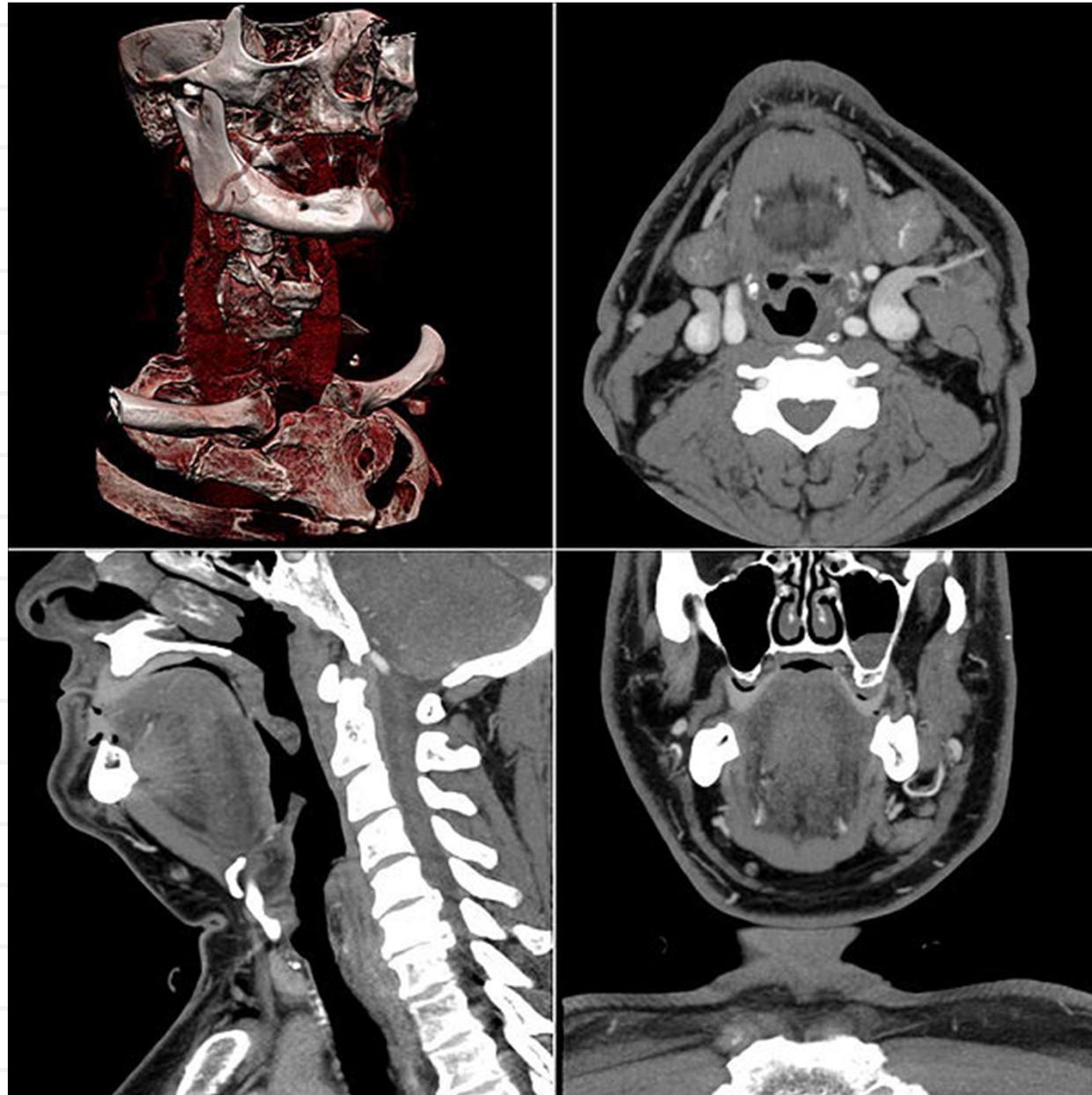
vector:

$$r(\omega) = \begin{bmatrix} x_1^\top \omega - y_1 \\ \vdots \\ x_n^\top \omega - y_n \end{bmatrix}$$

We can rewrite:

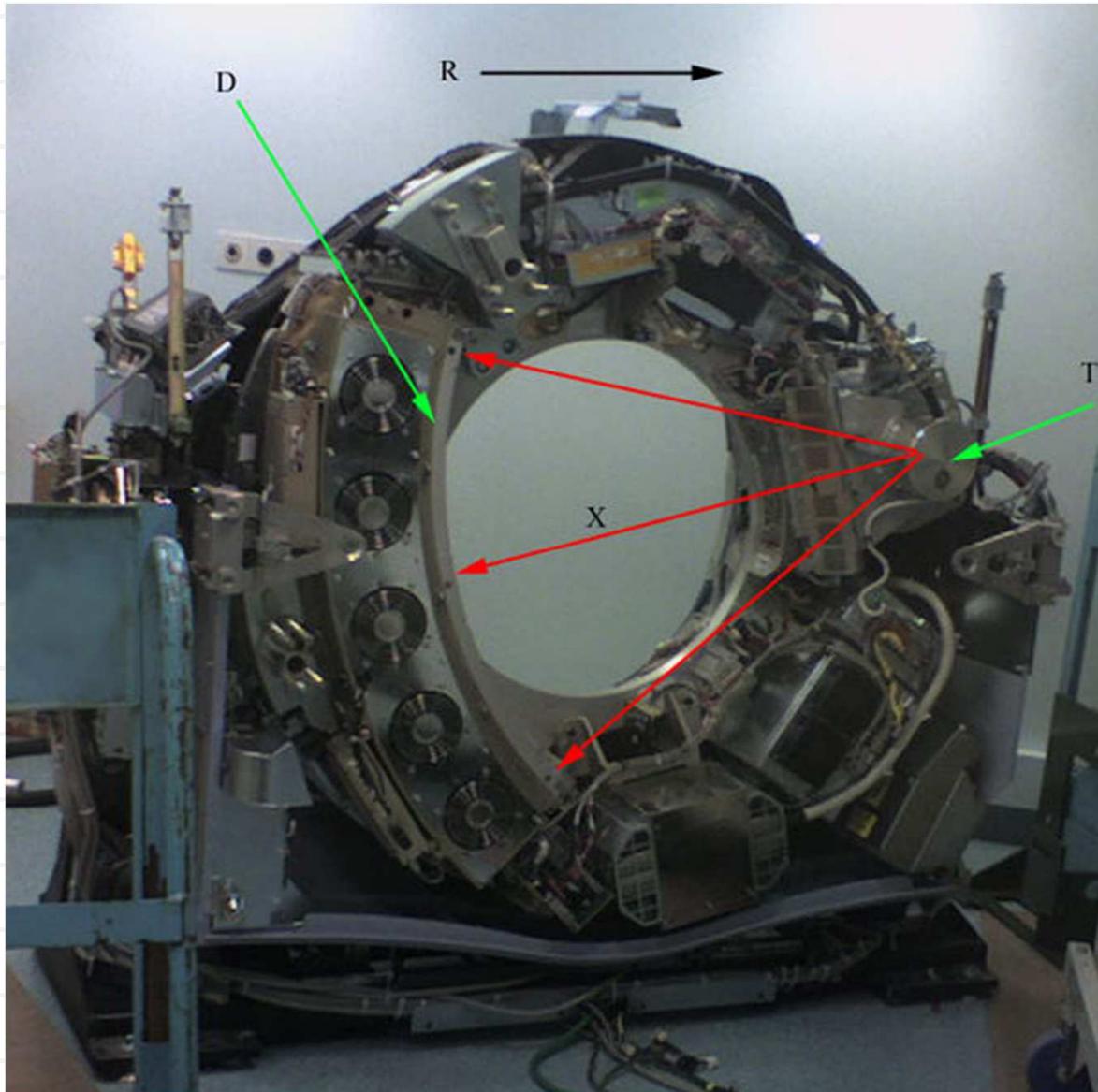
$$\min_{\omega} \frac{1}{2} \| r(\omega) \|_2^2$$

# Computer tomography



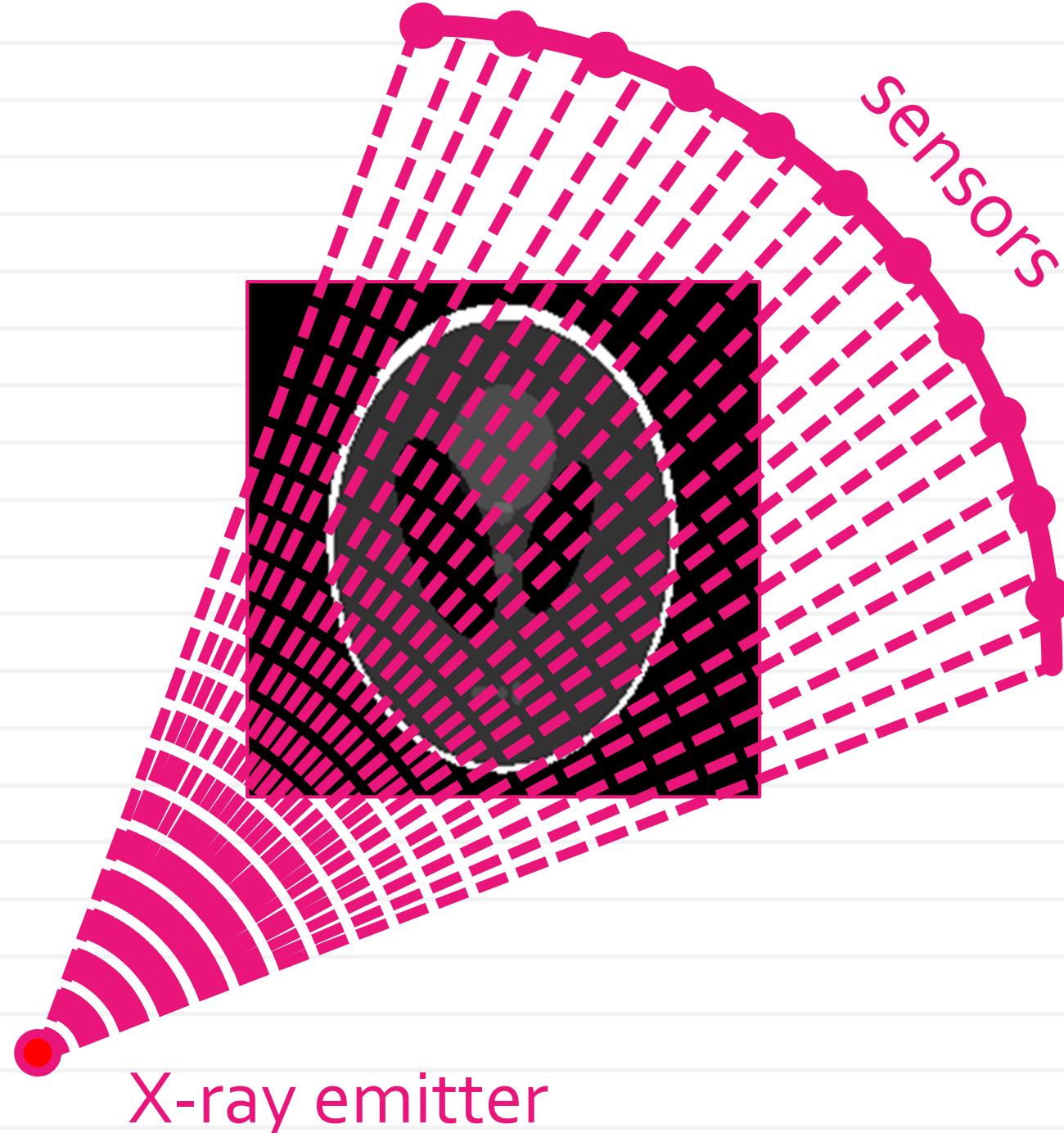
*Image from wikipedia*

# Computer tomography

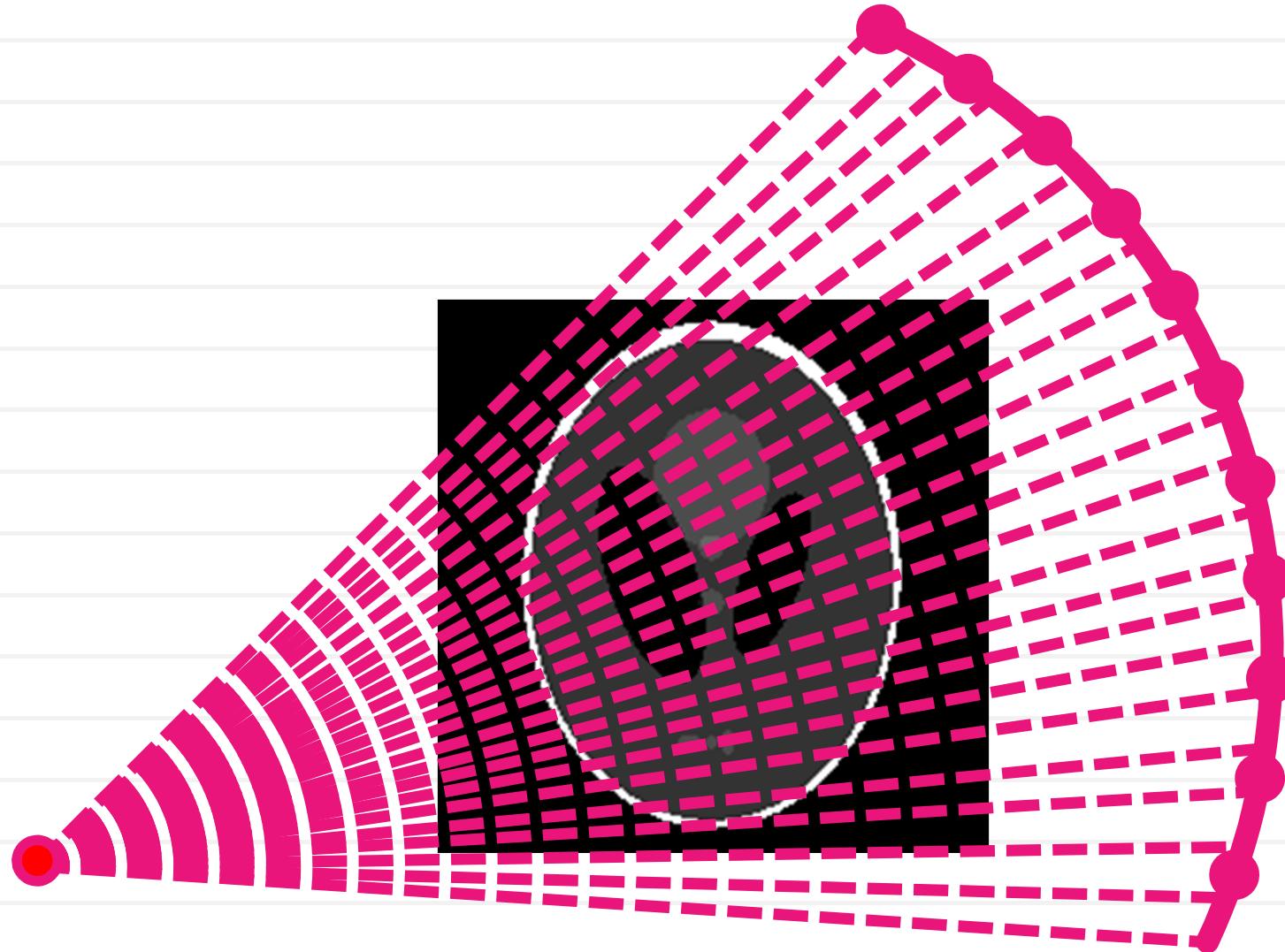


*Image from wikipedia*

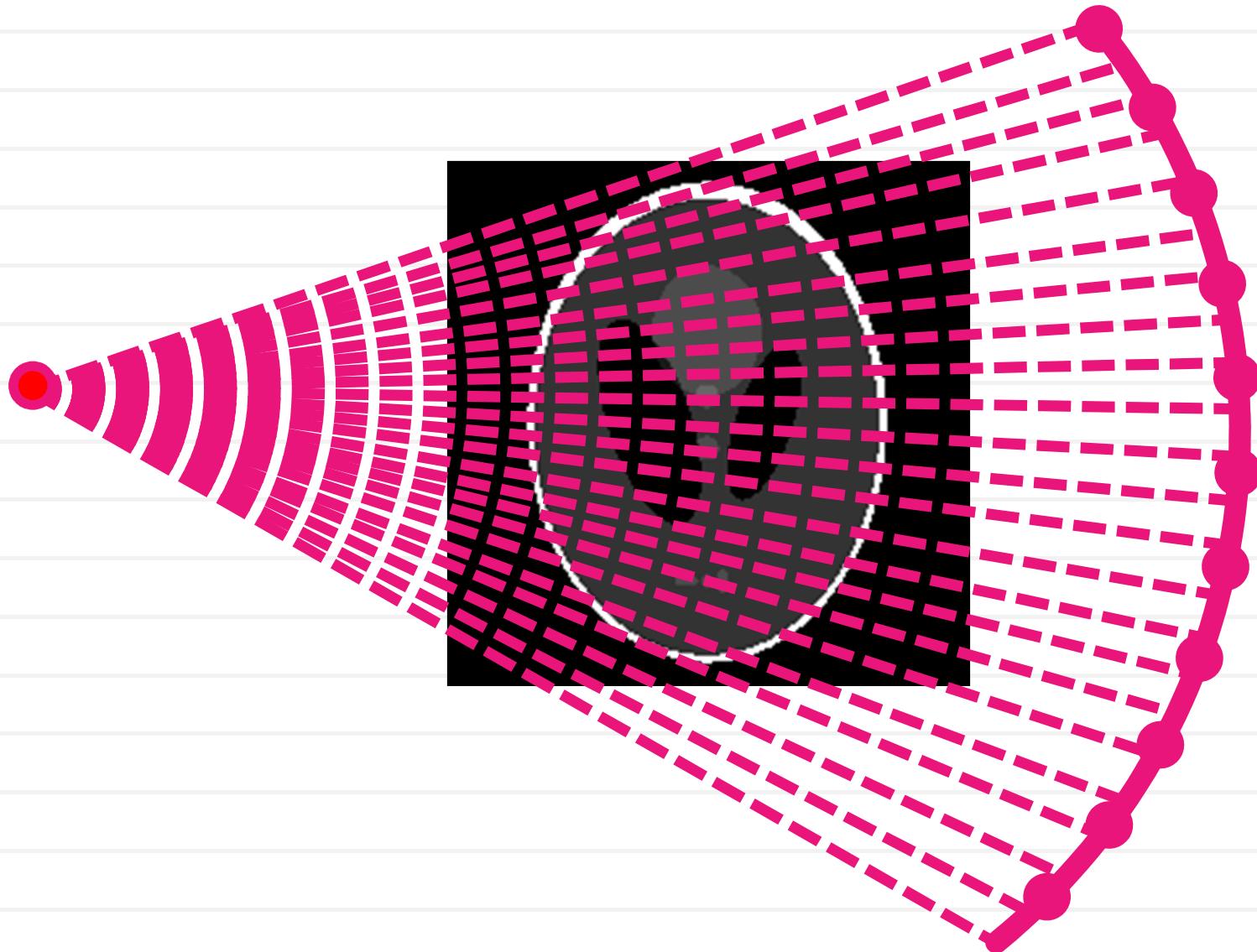
# Computer tomography



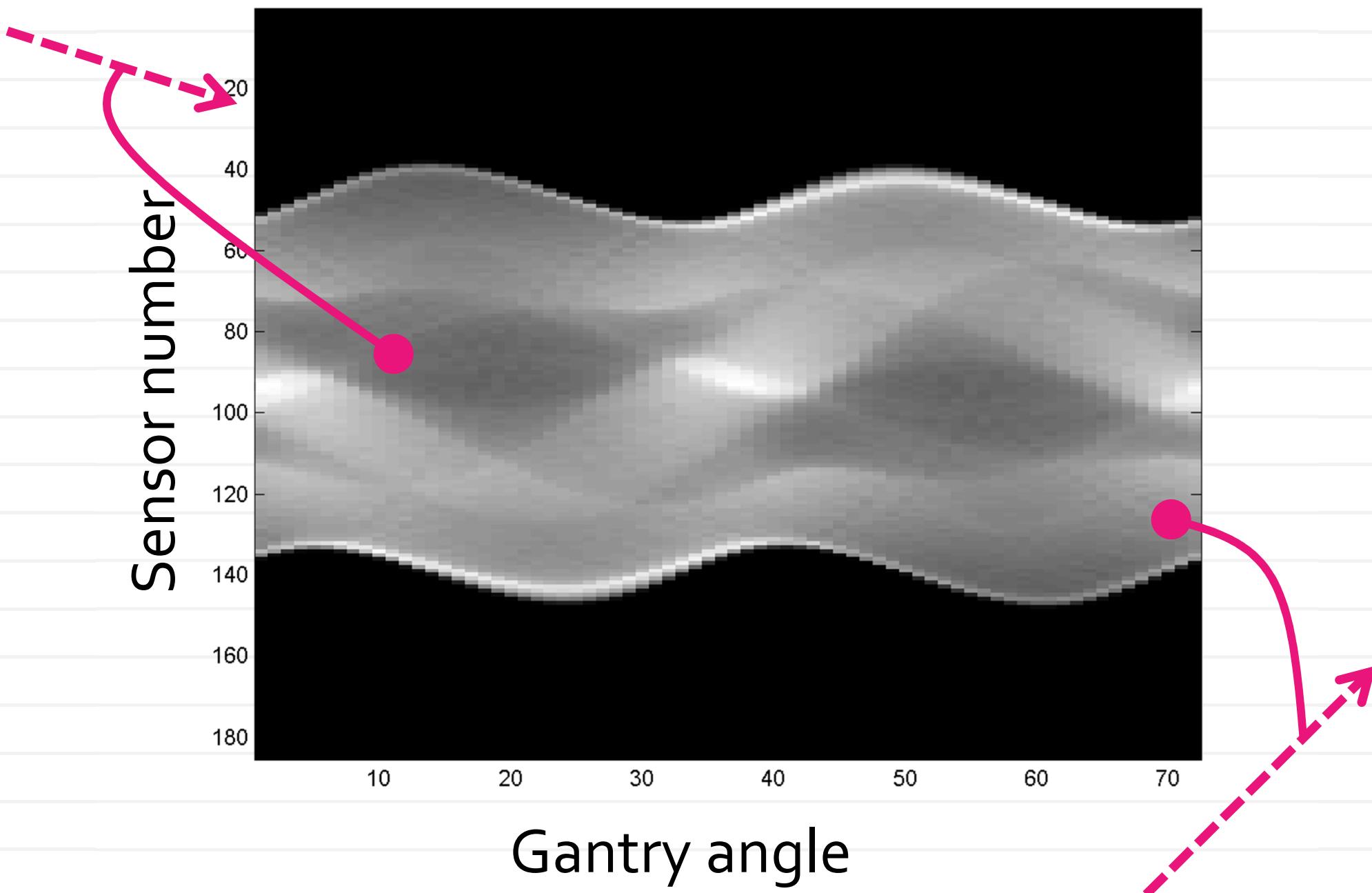
# Computer tomography



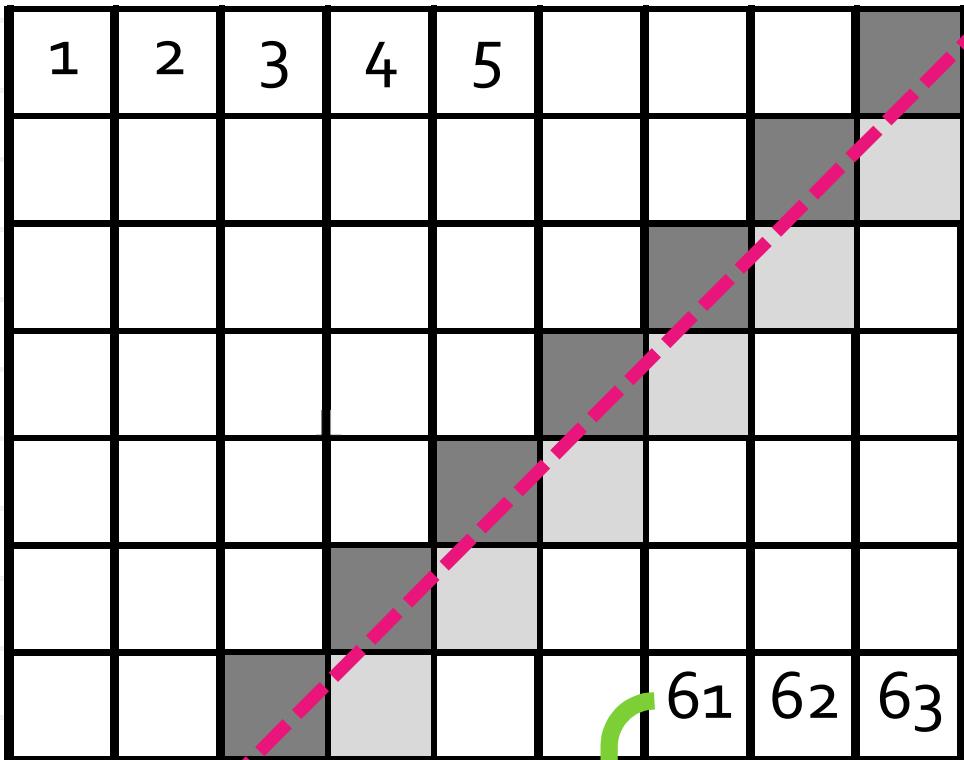
# Computer tomography



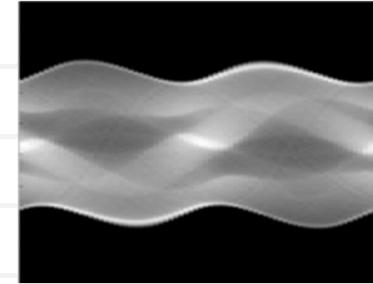
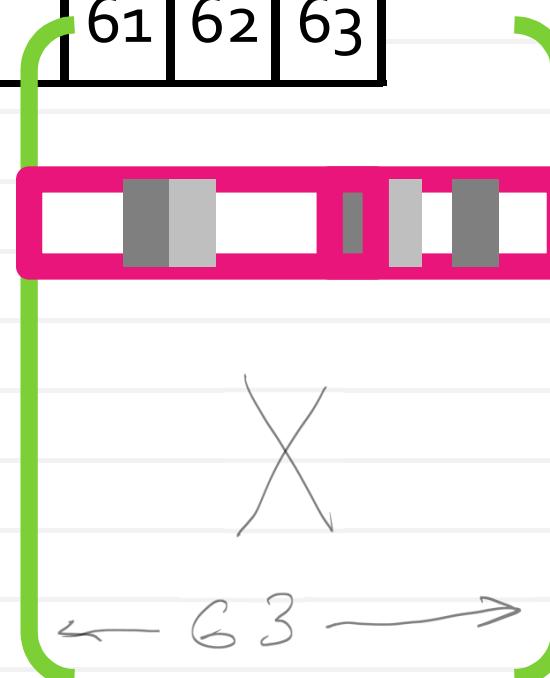
# A sinogram



# The least squares model



each row  
corresponds  
to one ray



*y*

*w*

$$y = w + \text{NOISE}$$

CT IMAGE

NOISE

# Finding the minimum

$$\min_w \frac{1}{2} \|Xw - y\|_2^2$$

Differentiating w.r.t.  $w$  and setting the gradient to zero, we get:

$$(Xw - y)^T X = 0$$

After algebra, we get *normal equations*:

$$X^T X w = X^T y$$

If normal equations are *full rank*, then:

$$w = (X^T X)^{-1} X^T y$$

# Adding non-linearity

We can adapt linear least squares to fitting non-linear functions, by transforming the input  $x$  non-linearly:

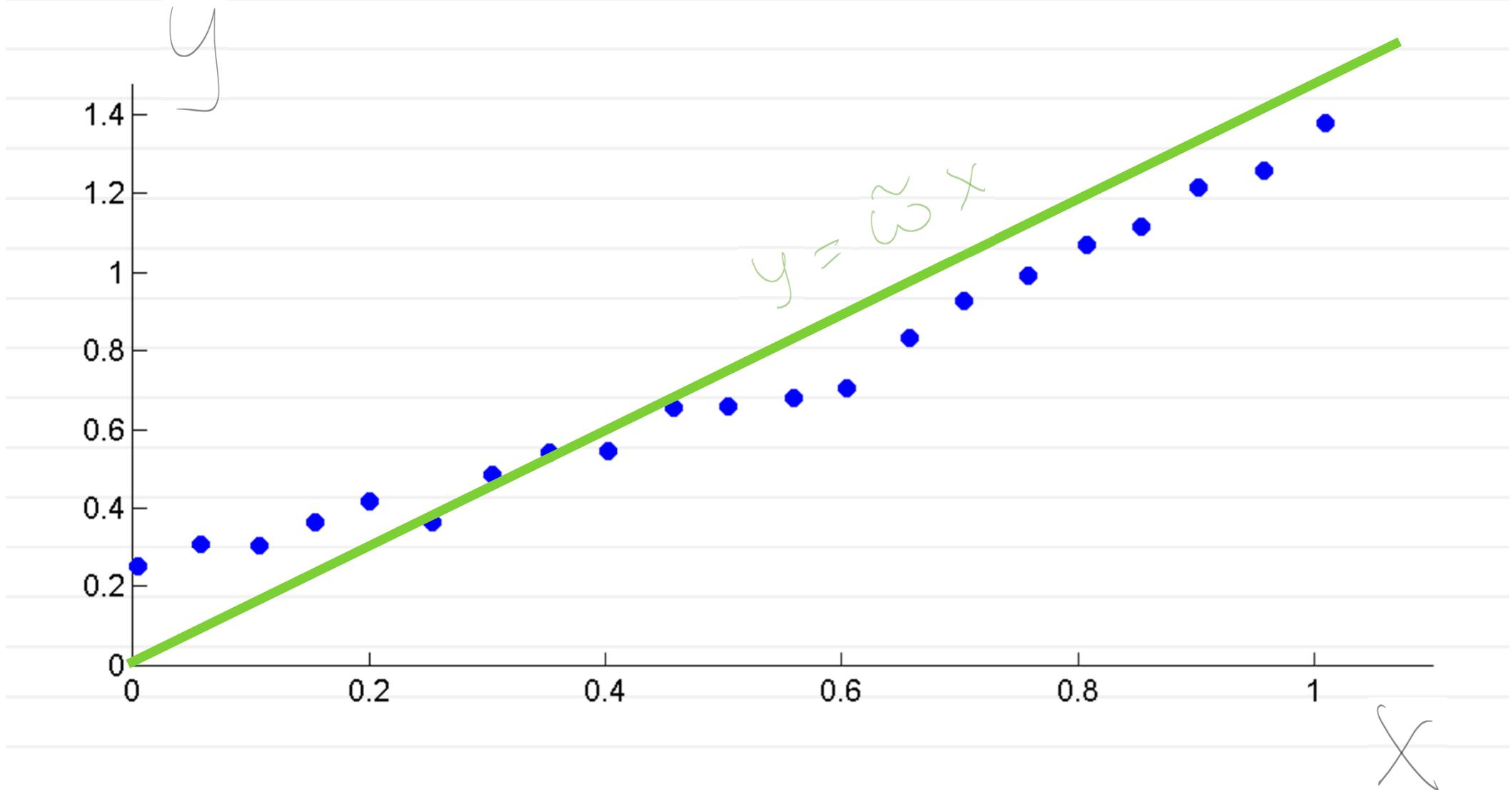
$$\mathbb{R}^m \ni x \longrightarrow \phi(x) \in \mathbb{R}^m$$

Our data matrix then becomes:

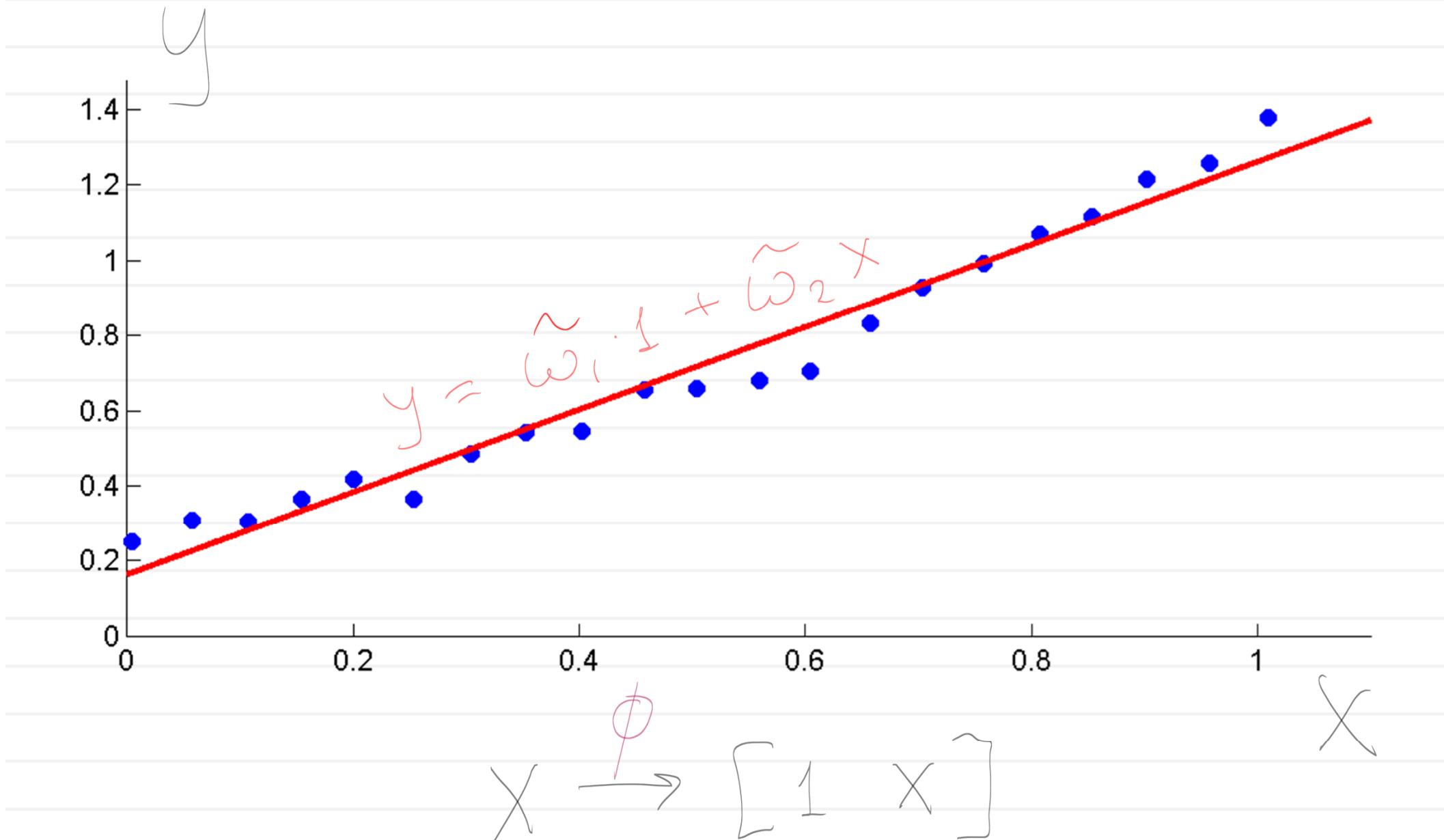
$$X = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_m(x_1) \\ \vdots & \ddots & \ddots & \vdots \\ \phi_1(x_n) & \phi_2(x_n) & \dots & \phi_m(x_n) \end{bmatrix}$$

Then, linear least squares proceed as before.

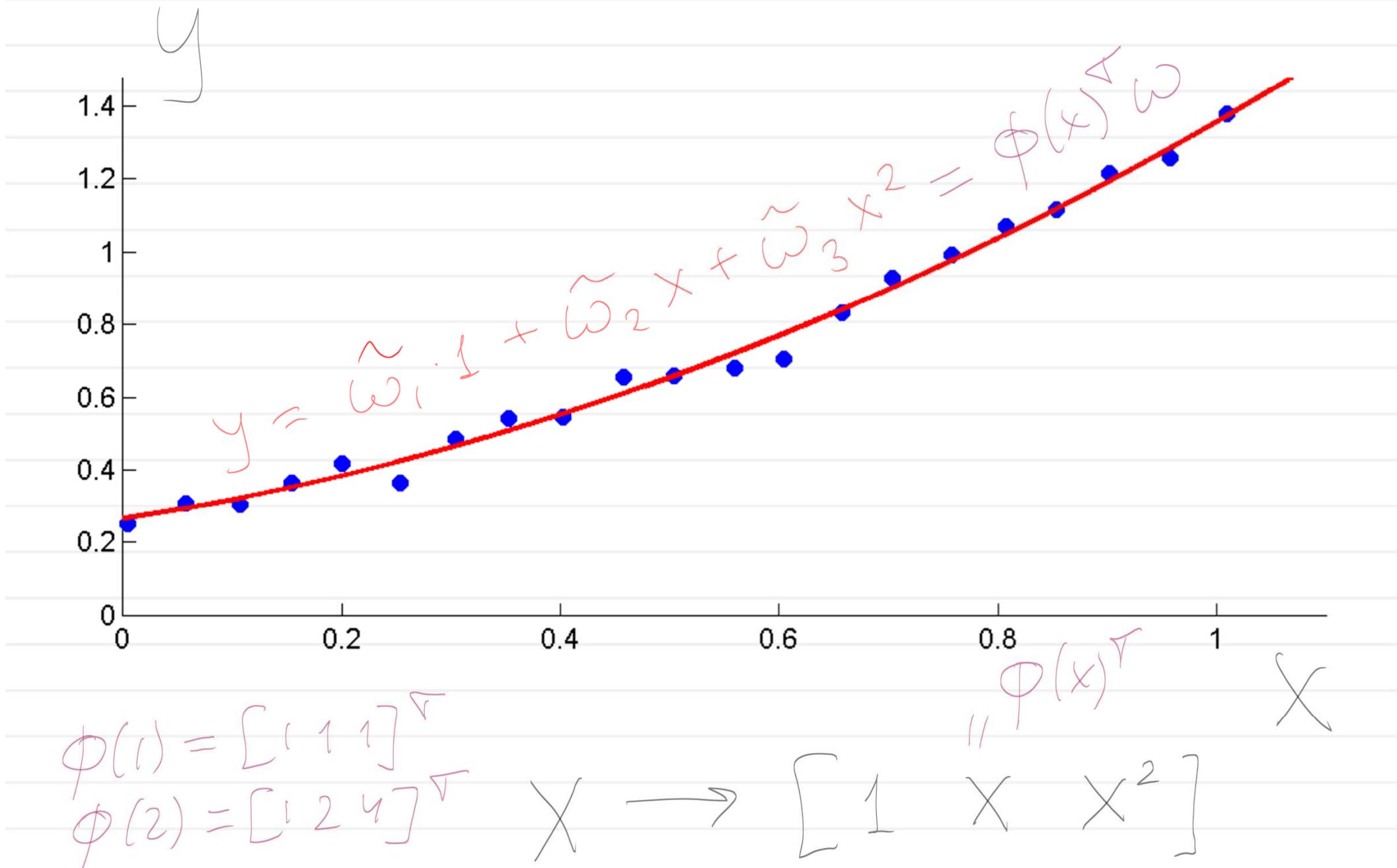
# Example fitting the data



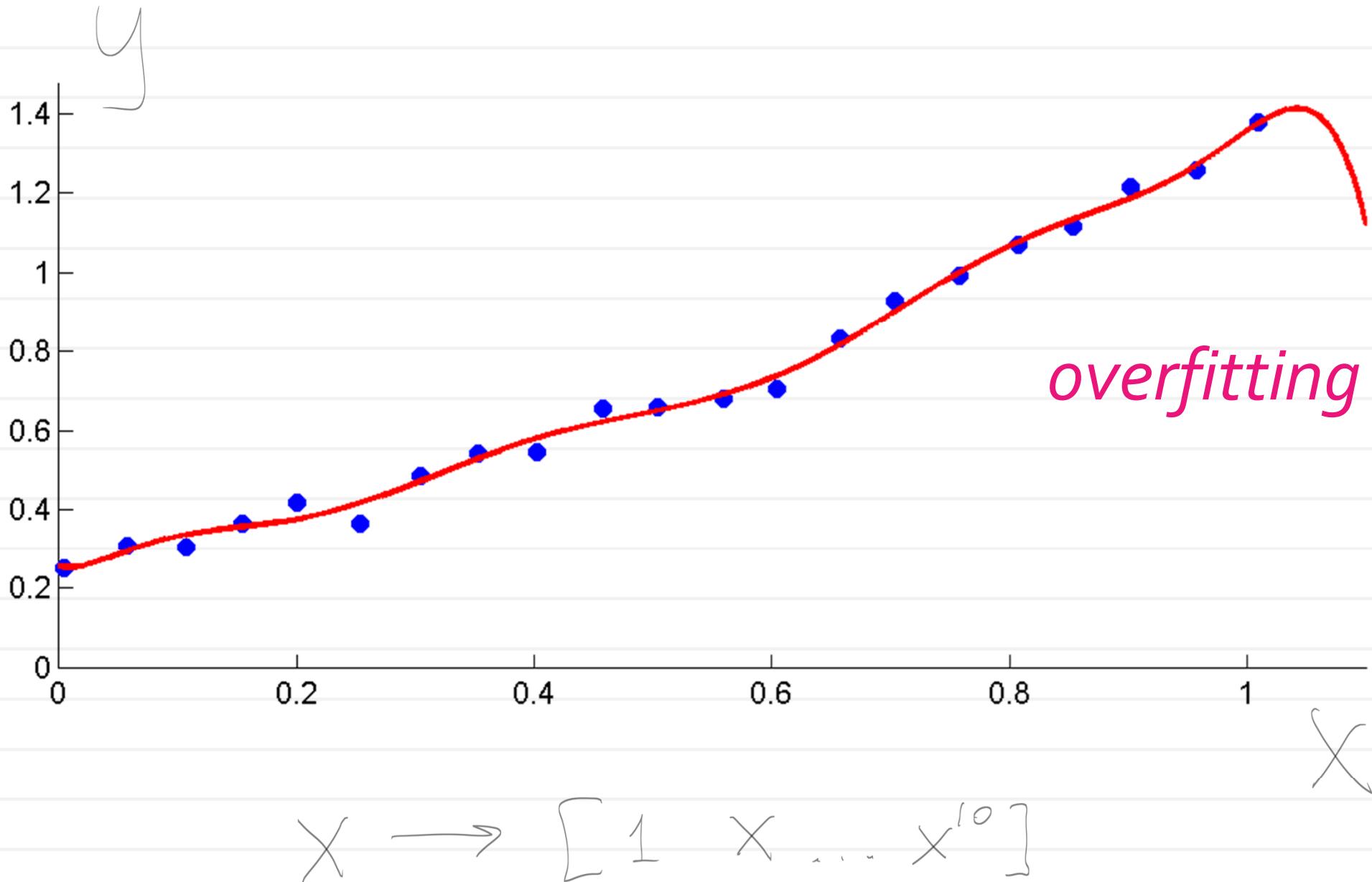
# Example fitting the data



# Example fitting the data



# Example fitting the data



# Choosing the right model

---

The most straight-forward and commonly used method is *validation*:

1. Split data into training and validation
2. Fit all models on training
3. Check how well do they perform on validation
4. Pick the best performer
5. Reestimate the model using full set

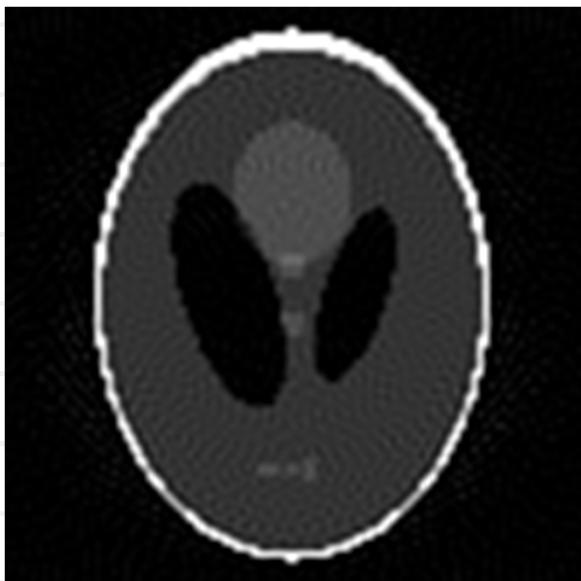
# Overfitting in computer tomography

Tomographic reconstruction  
suffers from overfitting as well:

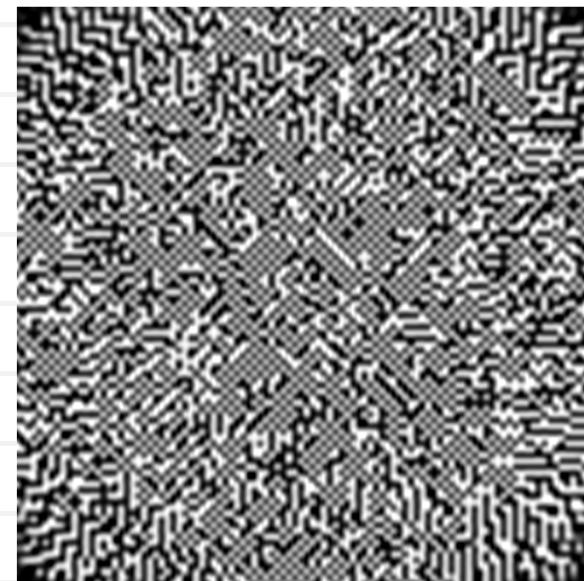
"True" phantom used to  
generate the data:



Reconstruction (no noise):



Reconstruction with noise:



# (Not fully) Bayesian approach

Our model (in the linear case):

$$X \omega = y + \text{"noise"}$$

Let us look at:

$$P(\omega | X, y)$$

Using Bayes rule:

$$P(\omega | X, y) = \frac{1}{P(X, y)} P(X, y | \omega) P(\omega)$$

*posterior*                    *likelihood*                    *prior*

We will now model the factors and then try to find  $\omega$  that maximizes the probability.

# (Not fully) Bayesian approach

$$\max_{\omega} P(\omega | X, Y)$$

$$\max_{\omega} \frac{1}{P(X, Y)} P(X, Y | \omega) P(\omega)$$

$$\min_{\omega} -\log P(X, Y | \omega) - \log(P(\omega))$$

*log-likelihood*                    *log-prior*

Different choices of log-likelihood and log-prior give rise to different regression methods with different statistical and computational performance

# Tikhonov-regularized least-squares

$$\cancel{X} \omega = y + \text{"noise"}$$

$$\min_{\omega} -\log P(x, y | \omega) - \log(P(\omega))$$

*log-likelihood*                           *log-prior*

We assume that the noise is independent between measurements and Gaussian. Then our likelihood is:

$$P(x, y | \omega) \propto \prod_i \exp\left(-\frac{\|y_i - x_i^\top \omega\|_2^2}{2}\right)$$

We assume that  $\omega$  is drawn from an isotropic zero-mean Gaussian prior:

$$\omega \sim \frac{1}{Z} \exp(-\lambda \|\omega\|_2^2)$$

Thus we get:

$$\min_{\omega} \frac{1}{2} \|X\omega - y\|_2^2 + \frac{M}{2} \|\omega\|_2^2$$

*M = 2λ*  $6^2$

# Finding the minimum

$$\min_w \frac{1}{2} \|Xw - y\|_2^2 + \frac{1}{2} \mu \|w\|_2^2$$

Differentiating w.r.t.  $w$  and setting the gradient to zero, we get:

$$(Xw - y)^T X + \mu I w = 0$$

$\begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix}$

After algebra, we get *normal equations*:

$$(X^T X + \mu I) w = X^T y$$

We can express  $w$  as:

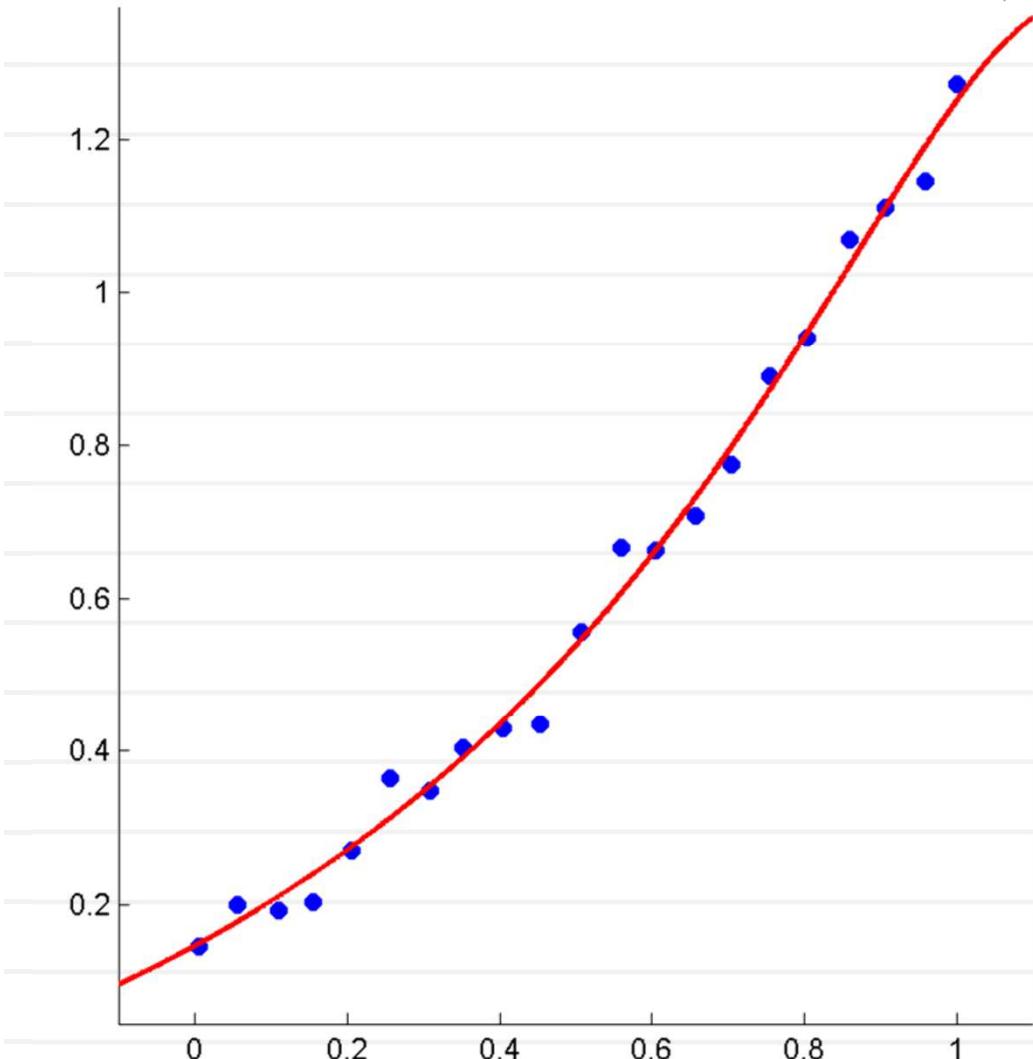
$$w = (X^T X + \mu I)^{-1} X^T y$$

# Fitting with regularization

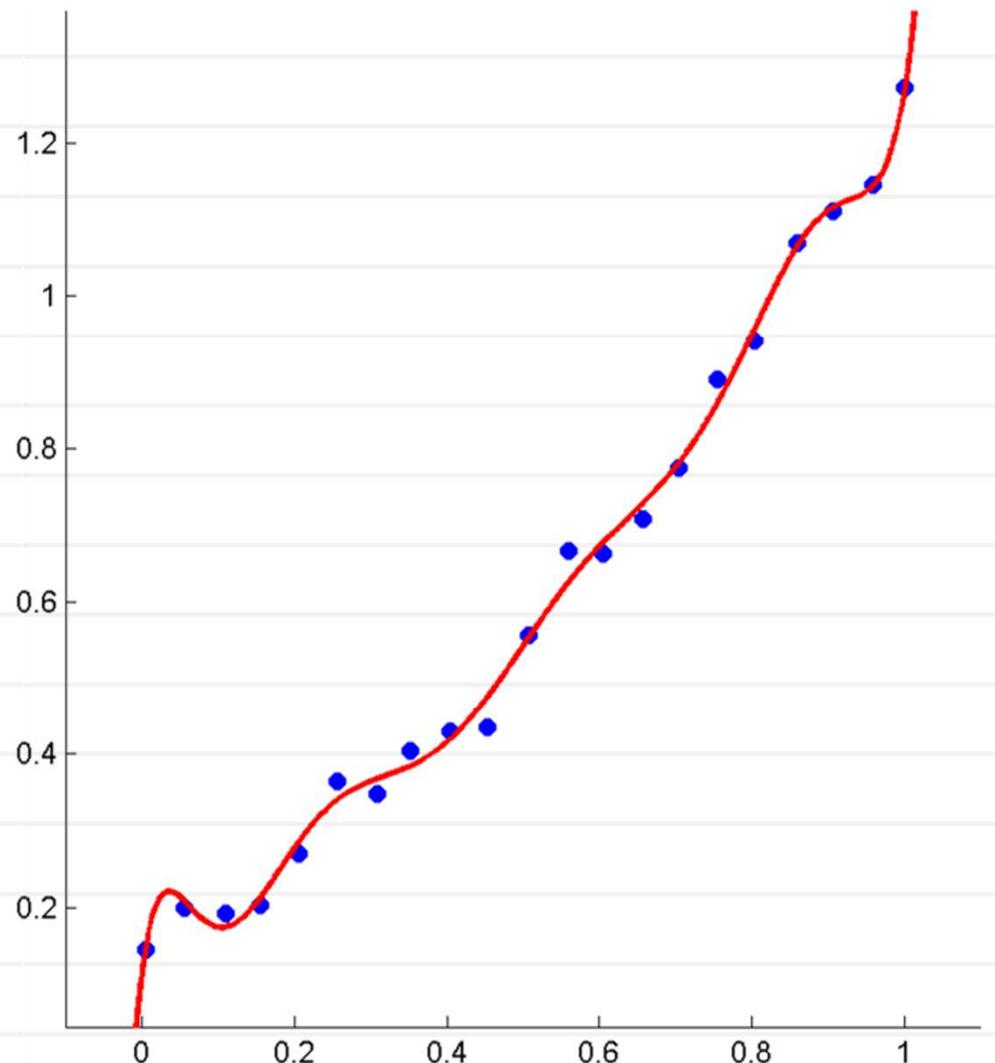
Fitting degree 10 polynomial:



$[1 \ X \ \dots \ X^{10}]$

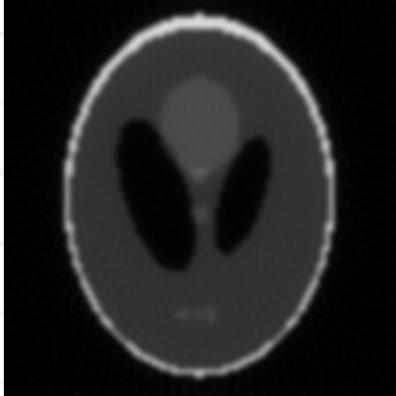
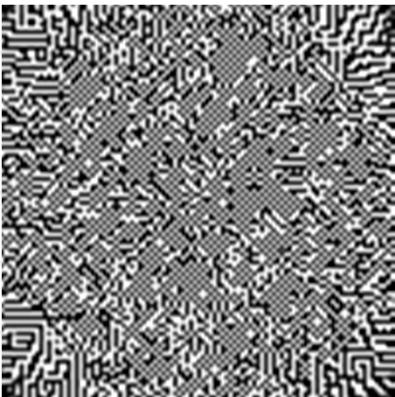
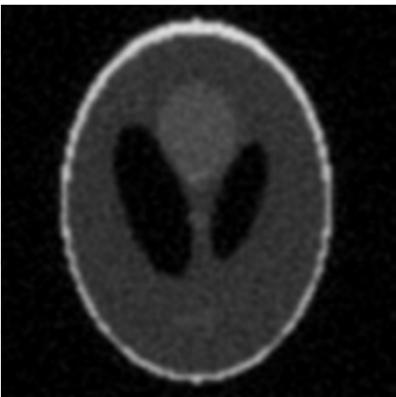


$$+ 0.01 \|\omega\|_2^2$$



No regularization

# Tomography with regularization

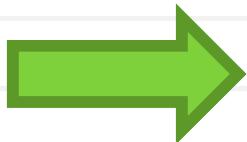
	No regularization	Strong regularization
Result from 90 angles		
No noise		
Strong noise		

# Non-linear least squares

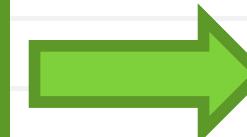
Typical usage scenario:

$$f(x; \omega)$$

$X$

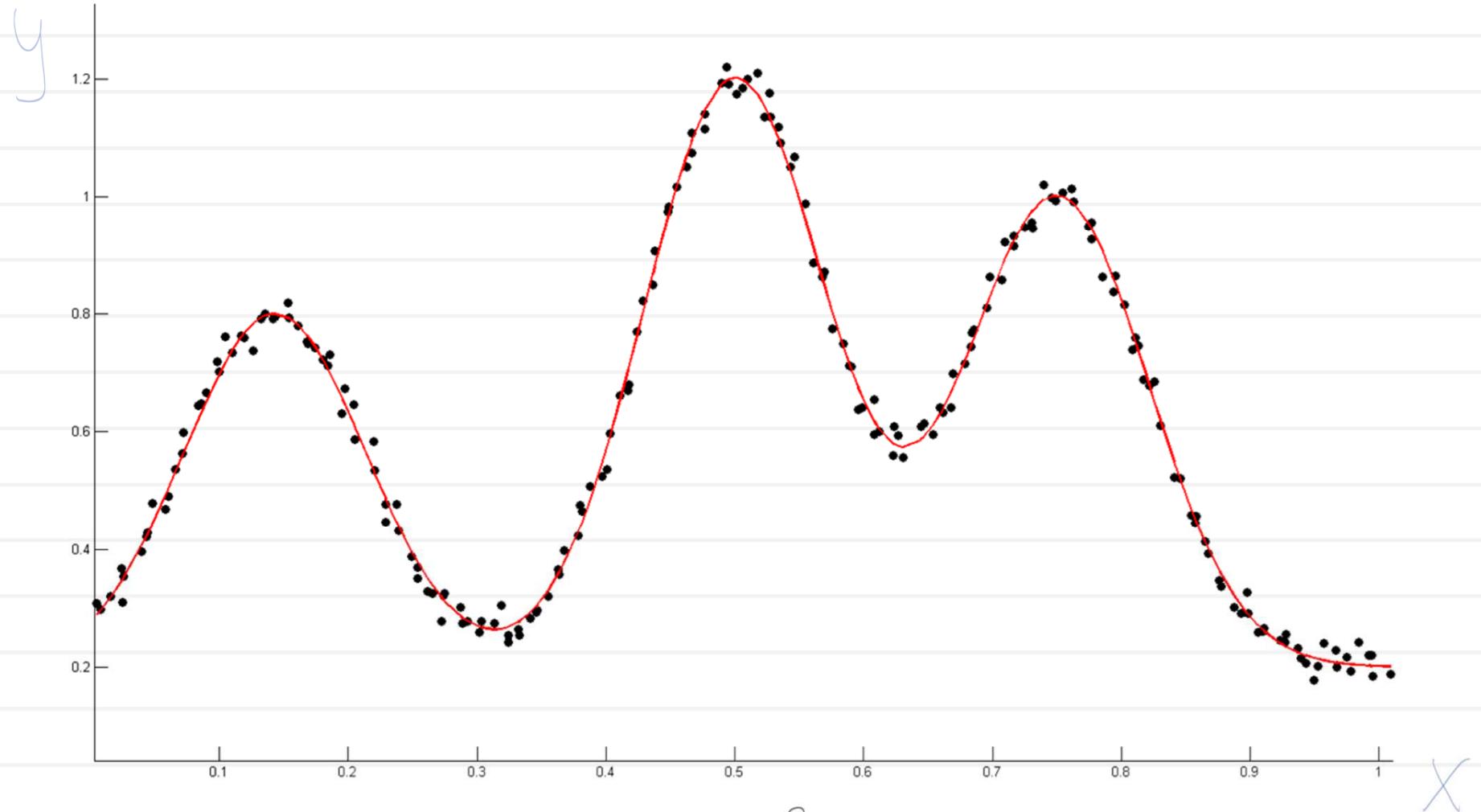


Complex model (e.g.  
physical simulation) that  
depends in a  
complicated non-linear  
way on  $w$



$y$

# Nonlinear curve fitting



$$f(x; \omega) = \omega_1 e^{-100(x - \omega_2)^2} + \omega_3 e^{-100(x - \omega_4)^2} + \omega_5 e^{-100(x - \omega_6)^2}$$

# Nonlinear curve fitting

$$f(x; \omega) = \omega_1 e^{-100(x - \omega_2)^2} + \omega_3 e^{-100(x - \omega_4)^2} + \omega_5 e^{-100(x - \omega_6)^2}$$

Assume that the noise is Gaussian, then to recover  $w$  we need to solve:

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i; \omega))^2 = F(\omega)$$

$$\min_{\omega} \frac{1}{2} \| r(\omega) \|^2 \quad r(\omega) = \begin{bmatrix} f(x_1; \omega) - y_1 \\ \vdots \\ f(x_n; \omega) - y_n \end{bmatrix}$$

# Non-linear least squares

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i; \omega))^2$$

The gradient of the objective function is:

$$\frac{d F(\omega)}{d \omega} = \nabla F(\omega) = J^T r(\omega)$$

, where the Jacobian  $J$  is:

$$J = \begin{bmatrix} \frac{\partial f(x_1; \omega)}{\partial \omega_1} & \dots & \frac{\partial f(x_1; \omega)}{\partial \omega_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(x_n; \omega)}{\partial \omega_1} & \dots & \frac{\partial f(x_n; \omega)}{\partial \omega_m} \end{bmatrix}$$

# Non-linear LS : gradient descent

We can try a standard gradient descent. We use a first order approximation in the vicinity of the current estimate  $w^t$ :

$$F(w_t + \Delta w) = F(w_t) + \nabla F(w_t)^T \Delta w$$

We then do steps along the gradient according some schedule  $s_t$  (or using line search):

$$w_{t+1} = w_t - s_t \gamma(w_t)^T r(w_t)$$

Reminder:  $\nabla F(w) = \gamma(w)^T r(w)$

# Non-linear LS: second-order method

Gradient descent is slow. Instead, we can get much faster convergence by looking at a second-order approximation:

$$F(\omega_t + \Delta\omega) = F(\omega_t) + \nabla F(\omega_t)^\top \Delta\omega + \frac{1}{2} \Delta\omega^\top \nabla^2 F(\omega_t) \Delta\omega$$

**Reminder:**  $\nabla F(\omega) = \gamma(\omega)^\top r(\omega)$

Let us estimate the matrix of second order derivatives:

$$\nabla(\nabla F(\omega)) = \underbrace{\frac{d \gamma(\omega)^\top}{d\omega} \cdot r(\omega)}_{\text{pink}} + \gamma(\omega)^\top \gamma(\omega)$$

Or, in more details:

$$\begin{aligned} [\nabla^2 F(\omega)]_{i,j} &= \sum_k \underbrace{\frac{\partial^2 f(x_k; \omega_t)}{\partial \omega_i \cdot \partial \omega_j} r_k(\omega_t)}_{\text{pink}} + \\ &\quad \sum_k \underbrace{\frac{\partial f(x_k; \omega_t)}{\partial \omega_i} \cdot \frac{\partial f(x_k; \omega_t)}{\partial \omega_j}}_{\text{green}} \end{aligned}$$

# Simplifying the second-order term

$$\left[ \nabla^2 F(\omega) \right]_{i,j} = \sum_k \frac{\partial^2 f(x_k; \omega_t)}{\partial \omega_i \cdot \partial \omega_j} r_k(\omega_t) + \sum_k \frac{\partial f(x_k; \omega_t)}{\partial \omega_i} \cdot \frac{\partial f(x_k; \omega_t)}{\partial \omega_j}$$

- Computing green is much easier than red ( $n$  first order derivatives vs  $n^2$  second-order derivatives per each square)
- The second-order term is approximated with the green term only (red term is omitted).
- This approximation is better for almost linear functions (second derivatives are small) or close to a good minimum (the residual vector is small).

# Properties of the approximate Hessian

$$\nabla^2 F(\omega) \approx \cancel{\frac{d}{d\omega} J(\omega)^T \cdot r(\omega)} + J(\omega)^T J(\omega)$$

- Easily computable
- Easily invertible (when Jacobian is sparse)
- Positive (semi)-definite
- We are replacing one approximation with another

$$\nabla^2 F(\omega) \cdot \Delta \omega = -\nabla F(\omega)$$

# Non-linear LS: Gauss-Newton method

At each step, the reduced second-order approximation is:

$$F(\omega_t + \Delta\omega) = F(\omega_t) + \nabla F(\omega_t) \cdot \Delta\omega + \frac{1}{2} \Delta\omega^T \mathcal{J}(\omega_t)^T \mathcal{J}(\omega_t) \Delta\omega$$

Equate the gradient to zero:

$$\mathcal{J}(\omega_t)^T r(\omega_t) + \mathcal{J}(\omega_t)^T \mathcal{J}(\omega_t) \cdot \Delta\omega = 0$$

A non-linear version of **normal equations** (consider linear normal equations for a linear function  $f(x, w) = x^T w$  at the origin)

Reminder: normal equations in the linear case

$$X^T y + X^T X \Delta w = 0$$

# Non-linear LS: Gauss-Newton method

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i; \omega))^2$$

Iterative algorithm:

1. Solve

$$J(\omega_t)^\top r(\omega_t) + J(\omega_t)^\top J(\omega_t) \cdot \Delta \omega = 0$$

2. Update the current position:

$$\omega_{t+1} = \omega_t + \Delta \omega$$

**Problem:** after two approximations some steps are too large and may lead to the increase/very slow decrease/oscillations.

# Levenberg-Marquardt method

$$\min_{\omega} \frac{1}{2} \sum_{i=1}^n (y_i - f(x_i; \omega))^2$$

On each step of Gauss-Newton we optimize:

$$F(\omega_t + \Delta\omega) = F(\omega_t) + \nabla F(\omega_t) \cdot \Delta\omega + \frac{1}{2} \Delta\omega^T Y(\omega_t)^T Y(\omega_t) \Delta\omega$$

The approximation makes sense in the vicinity.

Prevent the steps that are too large by adding an extra term (introducing the "**trust region**"):

$$\begin{aligned} \min_{\Delta\omega} & F(\omega_t) + \nabla F(\omega_t) \cdot \Delta\omega + \frac{1}{2} \Delta\omega^T Y(\omega_t)^T Y(\omega_t) \Delta\omega \\ & + \frac{1}{2} \lambda \|\Delta\omega\|^2 \end{aligned}$$

# Levenberg-Marquardt method

$$\min_{\Delta \omega} F(\omega_t) + \nabla F(\omega_t) \cdot \Delta \omega + \frac{1}{2} \Delta \omega^T J(\omega_t)^T J(\omega_t) \Delta \omega + \frac{1}{2} \lambda \|\Delta \omega\|^2$$

Equate the gradient to zero:

$$J(\omega_t)^T r(\omega_t) + J(\omega_t)^T J(\omega_t) \cdot \Delta \omega + \lambda I \Delta \omega = 0$$

We arrive at the **Levenberg-Marquardt** algorithm:

1. Solve

$$J(\omega_t)^T r(\omega_t) + J(\omega_t)^T J(\omega_t) \cdot \Delta \omega + \lambda I \Delta \omega = 0$$

2. Update the current position:  $\omega_{t+1} = \omega_t + \Delta \omega$

3. If the value is decreased “enough”, decrease  $\lambda$ .  
Otherwise, step back and increase it.

# Interpretation 1

$$\nabla^2 F(\omega) = \cancel{\frac{d \mathcal{J}(\omega)^T}{d \omega} r(\omega)} + \mathcal{J}(\omega)^T \mathcal{J}(\omega) + \lambda I$$

- Easily computable
- Easily invertible (when Jacobian is sparse)

$$\nabla^2 F(\omega) \cdot \Delta \omega = -\nabla F(\omega)$$

- **Positive definite, well conditioned**

# Interpretation 2

$$\mathcal{Y}(\omega_f)^\top r(\omega_f) + \mathcal{Y}(\omega_f)^\top \mathcal{Y}(\omega_f) \cdot \Delta\omega + \lambda I \Delta\omega = 0$$

$$\Delta\omega = - (\mathcal{Y}(\omega_f)^\top \mathcal{Y}(\omega_f) + \lambda I)^{-1} \mathcal{Y}^\top r(\omega_f)$$

How does the size of the extra term affects the algorithm?

$\lambda^t$  is small

$$\Delta\omega \approx - (\mathcal{Y}(\omega_f)^\top \mathcal{Y}(\omega_f))^{-1} \mathcal{Y}(\omega_f)^\top r(\omega_f)$$

**This is Gauss-Newton!**

$\lambda^t$  is large

$$\Delta\omega \approx -\frac{1}{\lambda} \mathcal{Y}(\omega_f)^\top r(\omega_f) = -\frac{1}{\lambda} \nabla F(\omega_f)$$

**This is gradient descent!**

**Depending on the size of the extra term, Levenberg-Marquardt switches between gradient descent like behaviour and Newton method-like behaviour.**