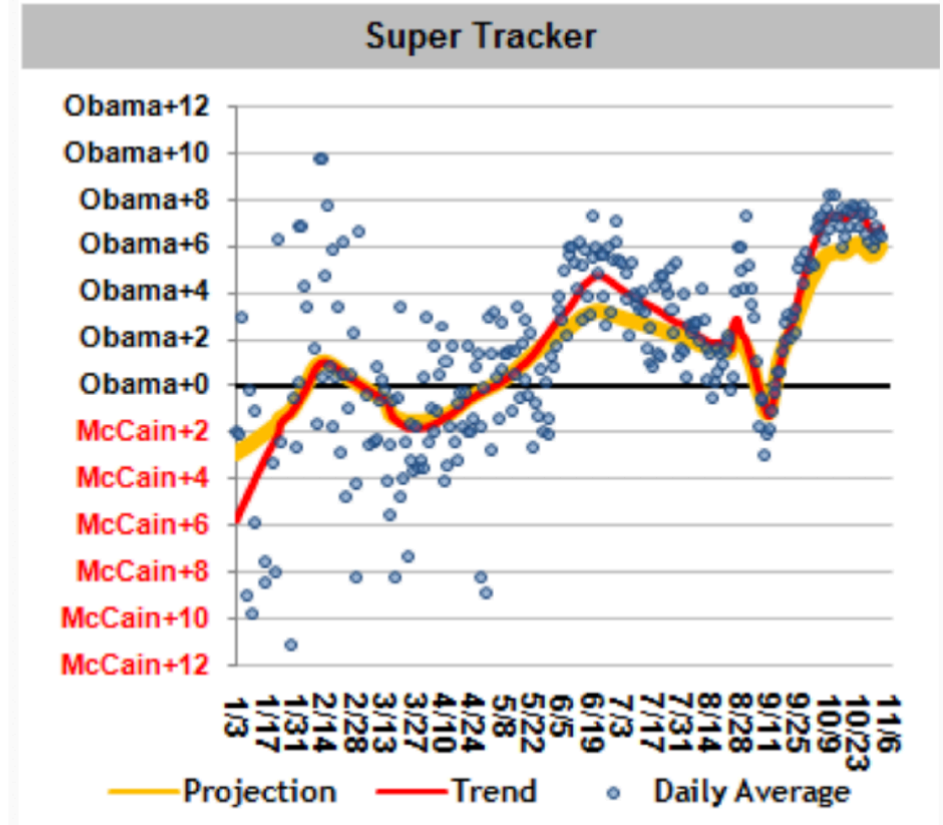
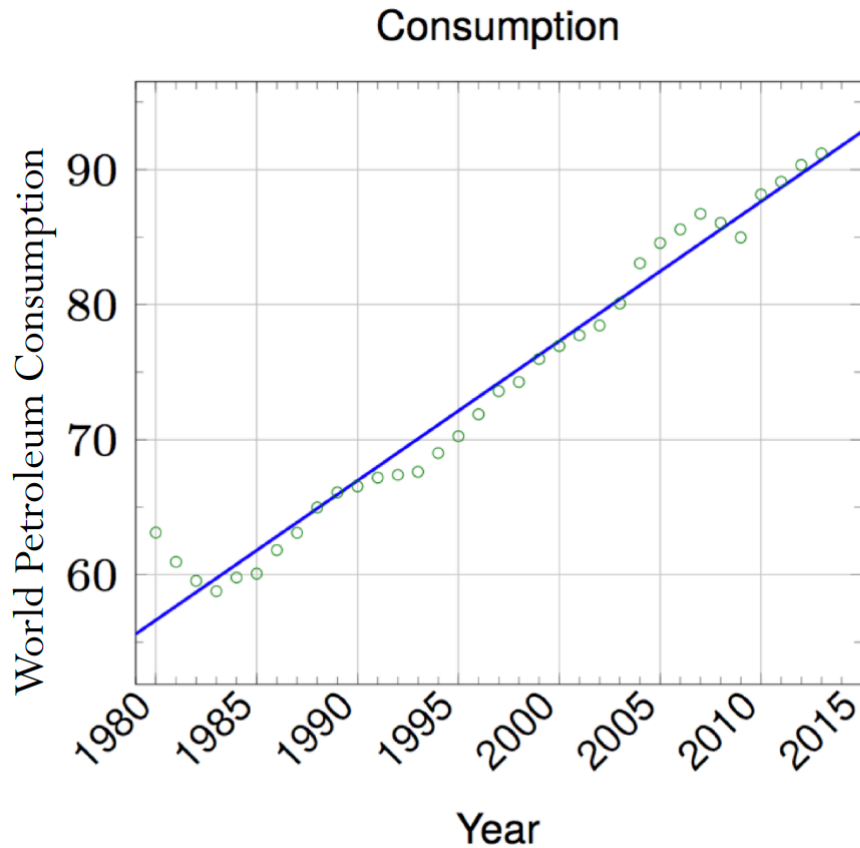


# Least Squares and Data Fitting

# Data fitting

How do we best fit a set of data points?



# Linear Least Squares

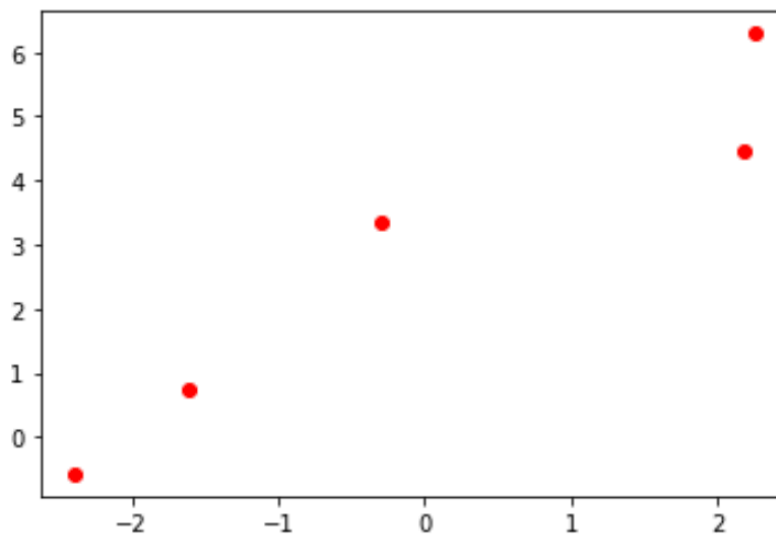
## 1) Fitting with a line

Given  $m$  data points  $\{(t_1, y_1), \dots, (t_m, y_m)\}$ , we want to find the function

$$y = x_0 + x_1 t$$

that best fit the data (or better, we want to find the coefficients  $x_0, x_1$ ).

Thinking geometrically, we can think “what is the line that most nearly passes through all the points?”



Given  $m$  data points  $\{\{t_1, y_1\}, \dots, \{t_m, y_m\}\}$ , we want to find  $x_0$  and  $x_1$  such that

$$y_i = x_0 + x_1 t_i \quad \forall i \in [1, m]$$

Given  $m$  data points  $\{\{t_1, y_1\}, \dots, \{t_m, y_m\}\}$ , we want to find  $x_0$  and  $x_1$  such that

$$y_i = x_0 + x_1 t_i \quad \forall i \in [1, m]$$

or in matrix form:

$$\begin{bmatrix} 1 & t_1 \\ \vdots & \vdots \\ 1 & t_m \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \quad \mathbf{A} \mathbf{x} = \mathbf{b}$$

$m \times n$     $n \times 1$     $m \times 1$

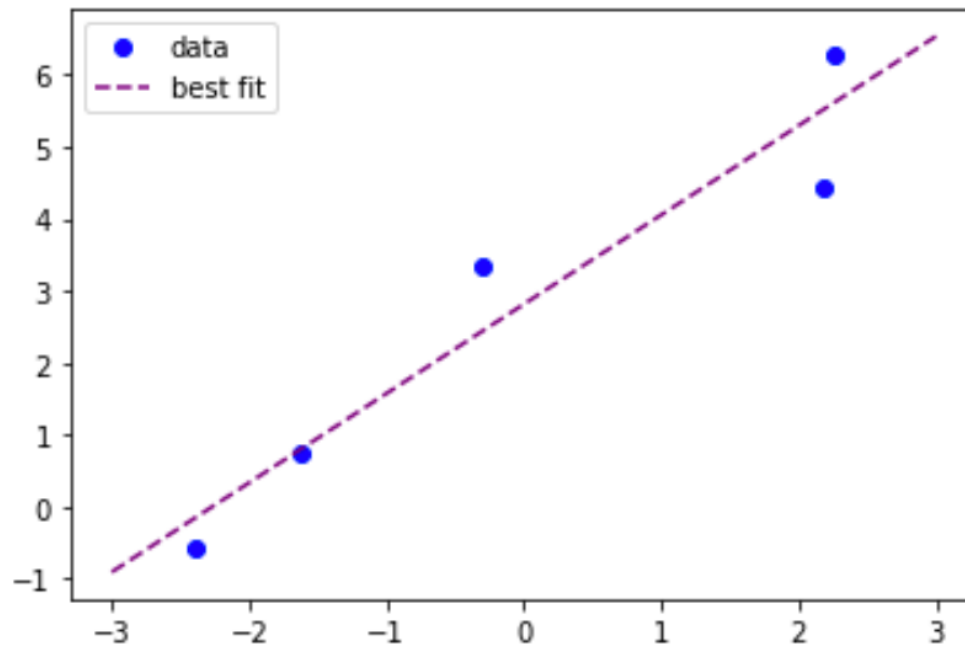
Note that this system of linear equations has more equations than unknowns –  
OVERDETERMINED  
SYSTEMS

We want to find the appropriate linear combination of the columns of  $\mathbf{A}$  that makes up the vector  $\mathbf{b}$ .

If a solution exists that satisfies  $\mathbf{A} \mathbf{x} = \mathbf{b}$  then  $\mathbf{b} \in \text{range}(\mathbf{A})$

# Linear Least Squares

- In most cases,  $\mathbf{b} \notin \text{range}(\mathbf{A})$  and  $\mathbf{A} \mathbf{x} = \mathbf{b}$  **does not have an exact solution!**



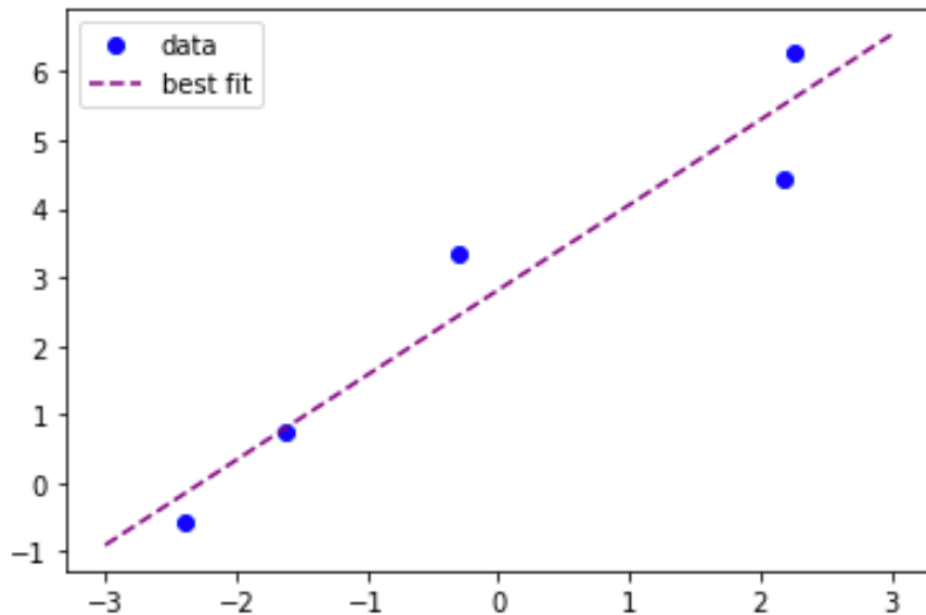
- Therefore, an overdetermined system is better expressed as

$$\mathbf{A} \mathbf{x} \cong \mathbf{b}$$

# Linear Least Squares

- **Least Squares:** find the solution  $\mathbf{x}$  that minimizes the residual

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \mathbf{x}$$



- Let's define the function  $\phi$  as the square of the 2-norm of the residual

$$\phi(\mathbf{x}) = \|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2$$

# Linear Least Squares

- **Least Squares:** find the solution  $\mathbf{x}$  that minimizes the residual

$$\mathbf{r} = \mathbf{b} - \mathbf{A} \mathbf{x}$$

- Let's define the function  $\phi$  as the square of the 2-norm of the residual

$$\phi(\mathbf{x}) = \|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2$$

- Then the least squares problem becomes

$$\min_{\mathbf{x}} \phi(\mathbf{x})$$

- Suppose  $\phi: \mathcal{R}^m \rightarrow \mathcal{R}$  is a smooth function, then  $\phi(\mathbf{x})$  reaches a (local) maximum or minimum at a point  $\mathbf{x}^* \in \mathcal{R}^m$  only if

$$\nabla \phi(\mathbf{x}^*) = 0$$



# How to find the minimizer?

- To minimize the 2-norm of the residual vector

$$\min_{\mathbf{x}} \phi(\mathbf{x}) = \|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2 = (\mathbf{b} - \mathbf{A} \mathbf{x})^T (\mathbf{b} - \mathbf{A} \mathbf{x})$$

# Linear Least Squares (another approach)

- Find  $\mathbf{y} = \mathbf{A} \mathbf{x}$  which is closest to the vector  $\mathbf{b}$
- What is the vector  $\mathbf{y} = \mathbf{A} \mathbf{x} \in \text{range}(\mathbf{A})$  that is closest to vector  $\mathbf{y}$  in the Euclidean norm?

# Summary:

- $\mathbf{A}$  is a  $m \times n$  matrix, where  $m > n$ .
- $m$  is the number of data pair points.  $n$  is the number of parameters of the “best fit” function.
- Linear Least Squares problem  $\mathbf{A} \mathbf{x} \cong \mathbf{b}$  *always* has solution.
- The Linear Least Squares solution  $\mathbf{x}$  minimizes the square of the 2-norm of the residual:

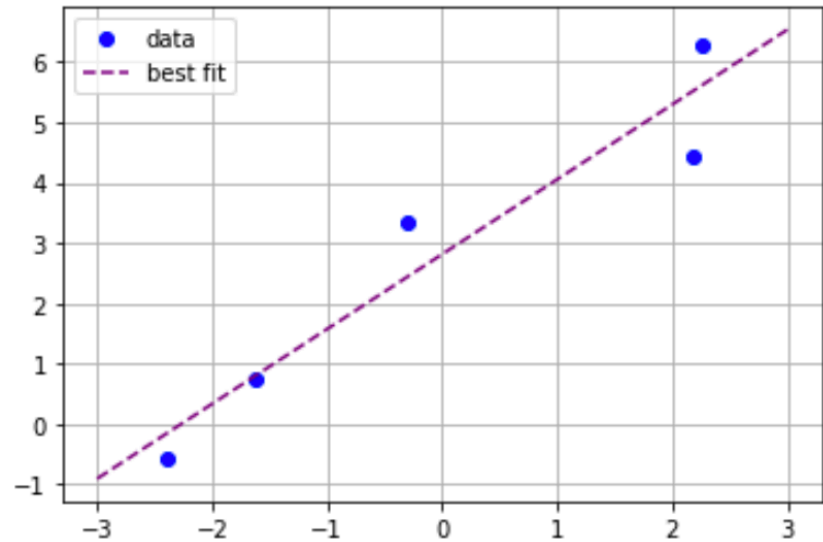
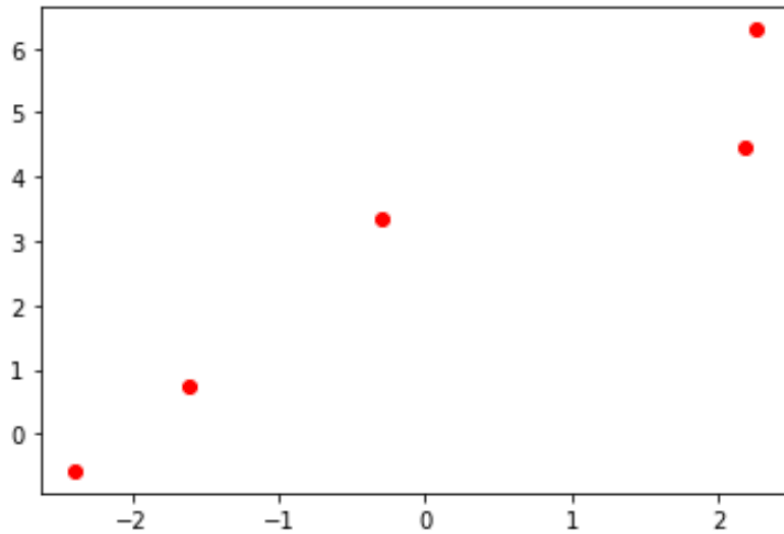
$$\min_{\mathbf{x}} \|\mathbf{b} - \mathbf{A} \mathbf{x}\|_2^2$$

- One method to solve the minimization problem is to solve the system of **Normal Equations**

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$$

- Let's see some examples and discuss the limitations of this method.

# Example:



t

```
array([-1.61477467, -2.3970584 , -0.30372944,  2.26304537,  2.188127  ])
```

b

```
array([ 0.74112251, -0.57768693,  3.33523097,  6.29377547,  4.44786481])
```

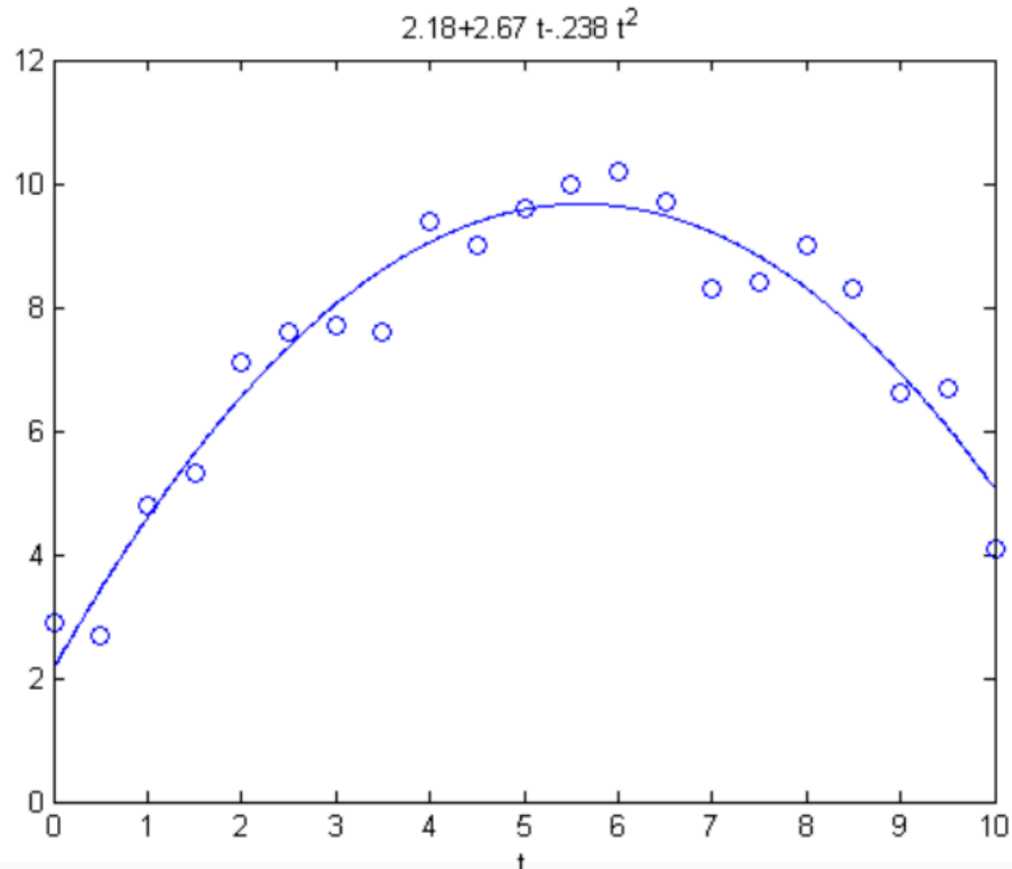
Solve:  $A^T A x = A^T b$

x

```
array([2.81441707, 1.24048133])
```

# Data fitting - not always a line fit!

- Does not need to be a line! For example, here we are fitting the data using a quadratic curve.

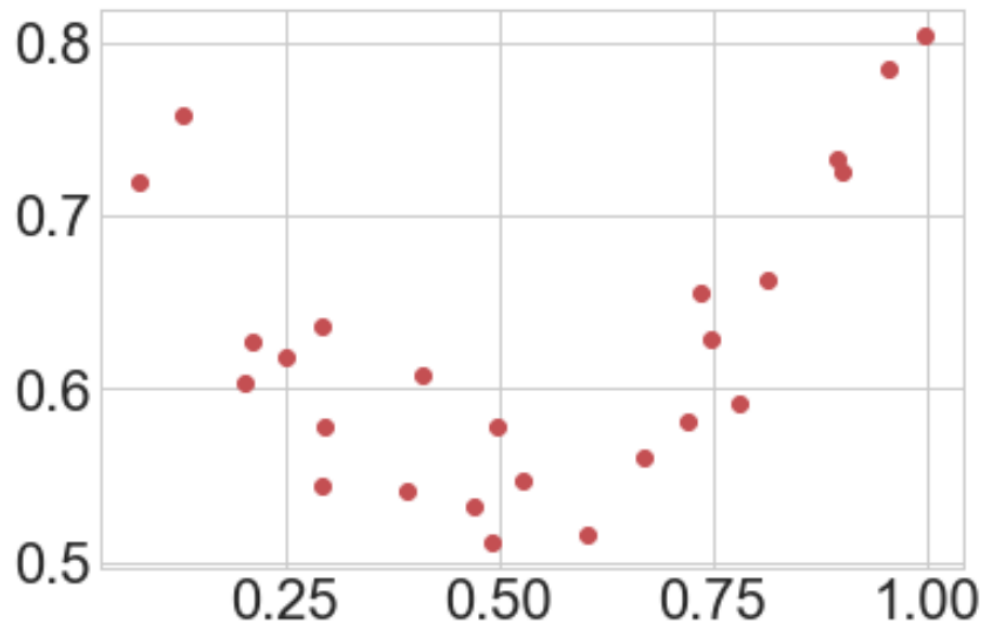


Linear Least Squares: The problem is **linear in its coefficients!**

# Another example

We want to find the coefficients of the quadratic function that best fits the data points:

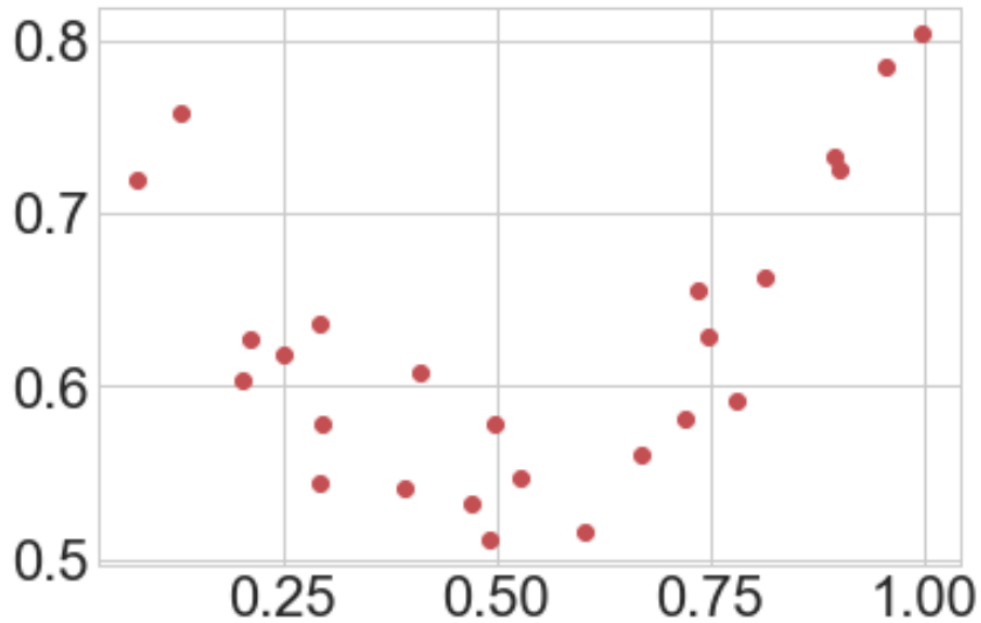
$$y = x_0 + x_1 t + x_2 t^2$$



We would not want our “fit” curve to pass through the data points exactly as we are looking to model the general trend and not capture the noise.

# Data fitting

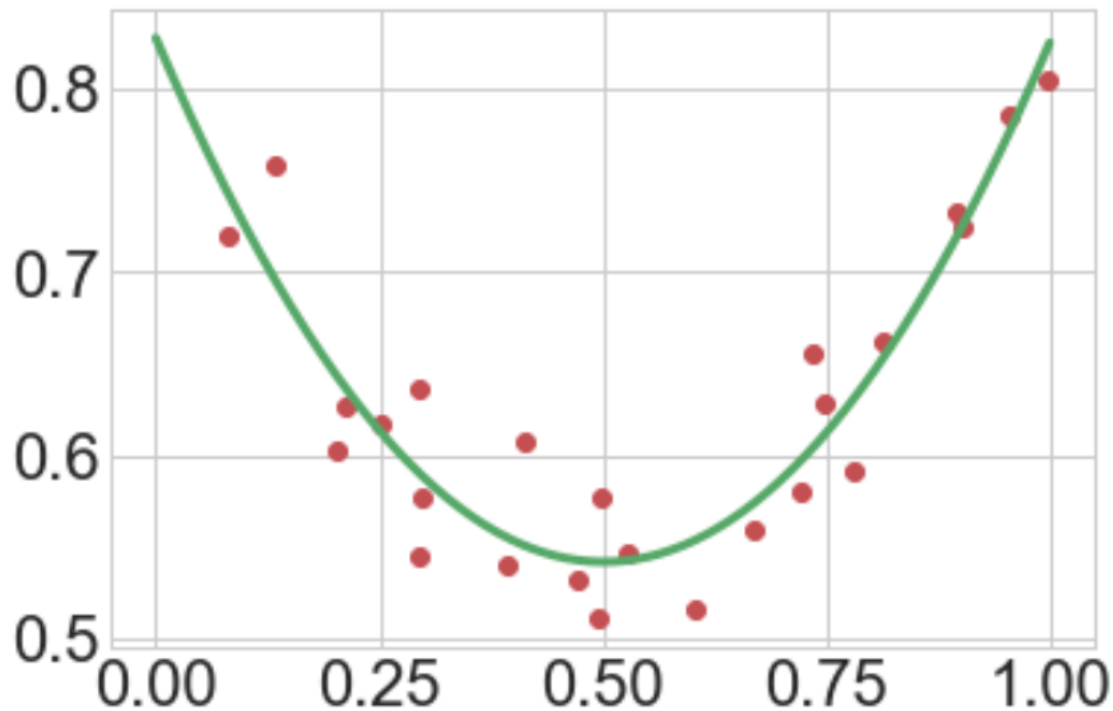
$$y = x_0 + x_1 t + x_2 t^2$$



# Data fitting

$$\begin{bmatrix} 1 & t_1 & t_1^2 \\ \vdots & \vdots & \vdots \\ 1 & t_m & t_m^2 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Solve:  $\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$





Which function is not suitable for linear least squares?

A)  $y = a + b x + c x^2 + d x^3$

B)  $y = x(a + b x + c x^2 + d x^3)$

C)  $y = a \sin(x) + b / \cos(x)$

D)  $y = a \sin(x) + x / \cos(bx)$

E)  $y = a e^{-2x} + b e^{2x}$

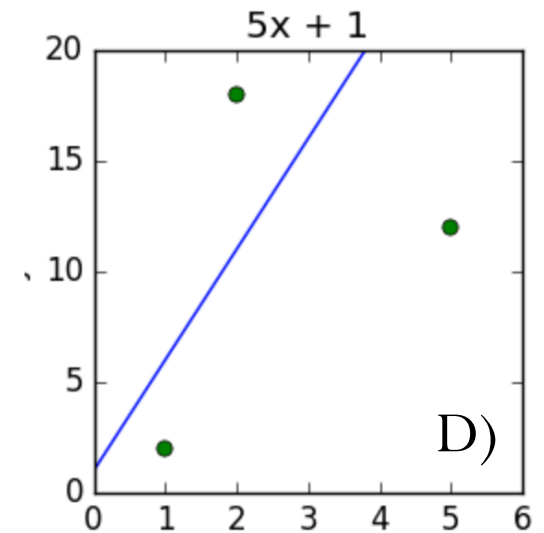
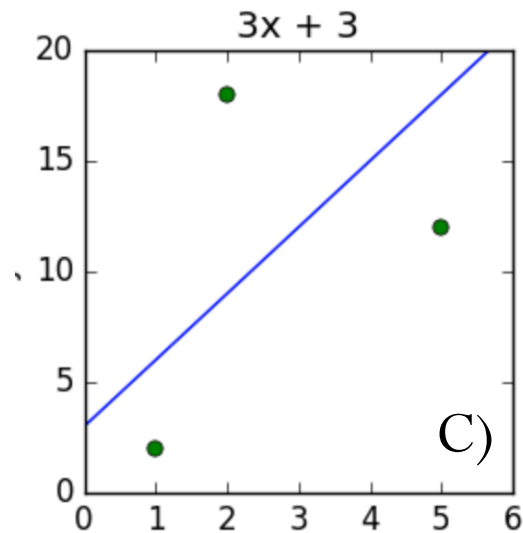
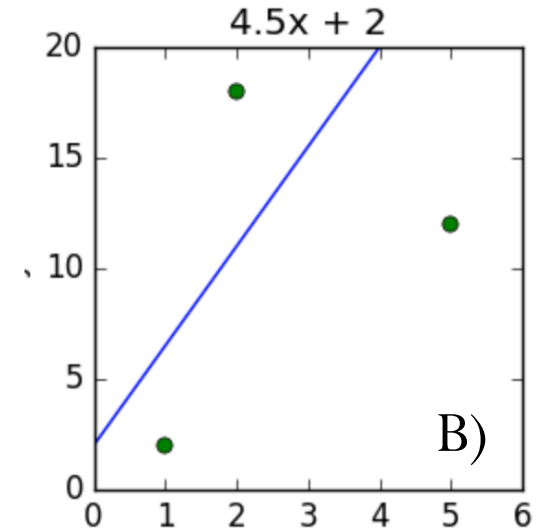
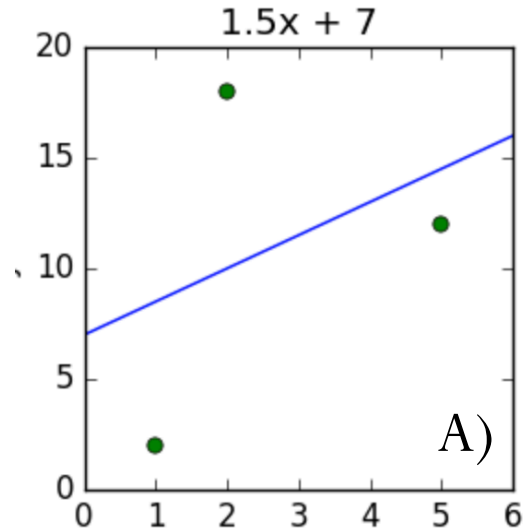
# Computational Cost

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{b}$$

# Short questions

Given the data in the table below, which of the plots shows the line of best fit in terms of least squares?

$x$	1	2	5
$y$	2	18	12



# Short questions

Given the data in the table below, and the least squares model

$$y = c_1 + c_2 \sin(t\pi) + c_3 \sin(t\pi/2) + c_4 \sin(t\pi/4)$$

written in matrix form as

$$A \begin{bmatrix} c_1 \\ c_2 \\ c_3 \\ c_4 \end{bmatrix} \cong \mathbf{y}$$

determine the entry  $A_{23}$  of the matrix  $\mathbf{A}$ .

Note that indices start with 1.

A)  $-1.0$

B)  $1.0$

C)  $-0.7$

D)  $0.7$

E)  $0.0$

$t_i$	$y_i$
0.5	0.72
1.0	0.79
1.5	0.72
2.0	0.97
2.5	1.03
3.0	0.96
3.5	1.00