

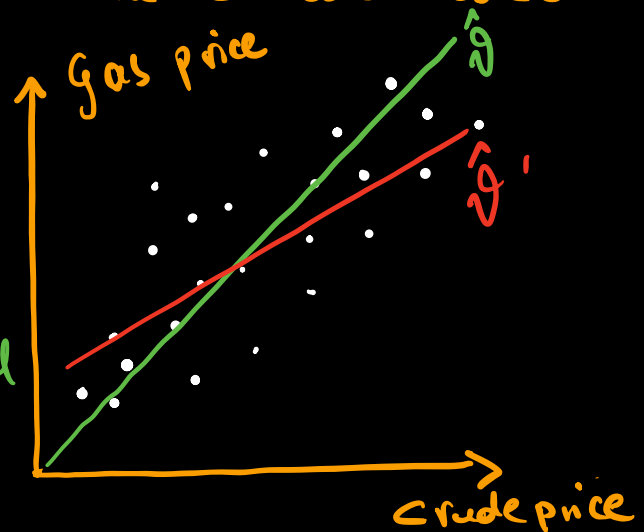
# Linear Models:

Example: We collect data  $(x_j, y_j)$  on crude oil price and gas price. We consider the statistical model

$$y_j | x_j \sim \theta x_j + \varepsilon_j$$

influence  
of crude price  
on gas

noise  
of external  
random  
factors



Goal: we want to find line so as to predict data. (line to fit the data)

Idea: Estimate  $\theta$  with the slope of the line minimizing the errors between the data.

$$\hat{g}(x, y) \text{ minimizes } \sum (y_j - \theta \cdot x_j)^2$$

Note: Throughout,  $x_j$  are considered fixed  
if  $x_j \sim P_X$  then  $x_j$  RV, results are conditional  
on  $(x_1, \dots, x_n)$ .

Example (offsets): Consider the affine model

$$y_j = \vartheta_0 + \vartheta_1 x_j + \varepsilon_j \quad \varepsilon_j \sim \mathcal{N}_0(0, \sigma)$$

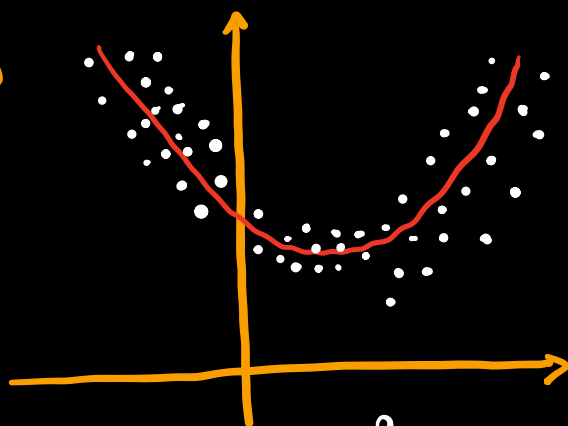
Then we can fit

$$\begin{aligned} \hat{\vartheta}(x, y) &= \arg \min \sum (y_j - (\vartheta_0 + \vartheta_1 x_j))^2 \\ &= \arg \min \sum \left( \underset{\text{observed}}{y_j} - \underbrace{\begin{pmatrix} \vartheta_0 \\ \vartheta_1 \end{pmatrix} \cdot \begin{pmatrix} 1 \\ x_j \end{pmatrix}}_{\text{predicted}} \right)^2 \end{aligned}$$

Example (Quadratic): Consider now

$$y_j = \vartheta_0 + \vartheta_1 x_j + \vartheta_2 x_j^2 + \varepsilon_j$$

We can still find a model fitting our data



$$\hat{\vartheta}(x, y) = \arg \min \sum \left( y_j - \begin{pmatrix} 1 \\ x_j \\ x_j^2 \end{pmatrix} \begin{pmatrix} \vartheta_0 \\ \vartheta_1 \\ \vartheta_2 \end{pmatrix} \right)^2$$

$\leadsto$  The model is still linear in  $\vartheta$ !

Definition: A linear model is a statistical model where:

$$y_j = \sum_{\ell=1}^L \vartheta_\ell \bar{\phi}(x_j)_\ell + \sigma \cdot z_j \quad j \in \{1, \dots, n\}$$

$$\mathbb{E}[z_j] = 0$$

Unknown Parameters:  $\vartheta_1, \dots, \vartheta_p$

Note: We can represent this in vector notation

$$y = \Phi \vartheta + \epsilon \cdot z$$

where •  $\Phi_{i \ell} = \Phi(x_i)_\ell$  is the design matrix

•  $\vartheta = (\vartheta_1, \dots, \vartheta_p)$  and  $\epsilon^2$  are the parameters

We say that the linear model is underparametrized if  $\text{rank}(\Phi) = p$

Definition: The least square estimators

$\hat{\vartheta}$  of  $\vartheta$  is given by:

$$\hat{\vartheta} = \arg \min_{\vartheta \in \mathbb{R}^p} \|y - \Phi \vartheta\|_2^2$$

Lemma: In the underparametrized setting  $\hat{\vartheta}$  is unique and can be written as

$$\hat{\vartheta}(y) = (\Phi^T \Phi)^{-1} \Phi^T y$$

Proof:  $\|y - \Phi \vartheta\|_2^2 = (y - \Phi \vartheta)^T (y - \Phi \vartheta)$

Solve  $\nabla_{\vartheta} (y - \Phi \vartheta)^T (y - \Phi \vartheta) = 0$

$\Rightarrow \Phi^T \Phi \vartheta = \Phi^T y$

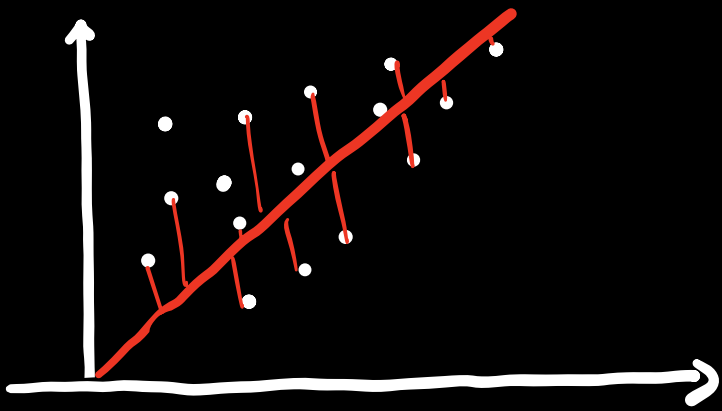
Since  $\Phi^T \Phi \in \mathbb{R}^{p \times p}$  is full rank, it is invertible so

$\hat{\vartheta}(y) = (\Phi^T \Phi)^{-1} \Phi^T y \quad \square$

Note:  $\hat{\vartheta}(y)$  is linear in  $y$

Geometric interpretation: (least squares)

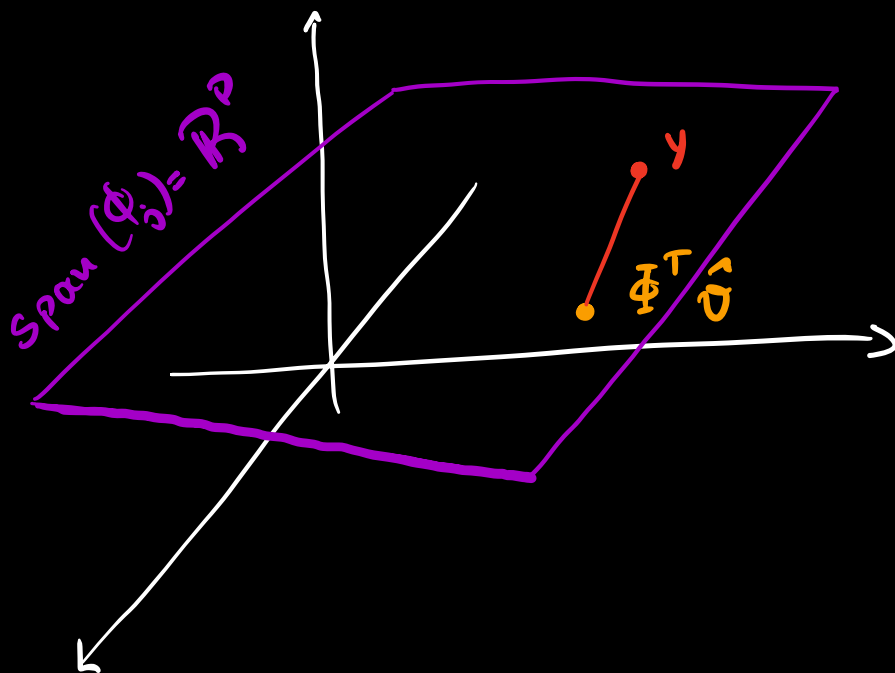
- $\arg \min \sum (y_i - \Phi \vartheta_i)^2$ :  $\vartheta$  minimizes the total vertical distance between  $n$  points in  $(p+1)$  dimensions. We are fitting a  $p$ -dimensional hyperplane. If we include an intercept, we fit  $(p-1)$ -dimensional plane offset from the origin.



- $\arg \min \|\mathbf{y} - \Phi \boldsymbol{\theta}\|_2^2 = \arg \min \|\mathbf{y} - \sum_{j=1}^p \bar{\Phi}_j^T \theta_j\|_2^2$

$\hat{\boldsymbol{\theta}}$  contains the coefficients of the column vector of  $\Phi$ .  $\Phi \hat{\boldsymbol{\theta}}$  is the orthogonal projection of  $\mathbf{y}$  onto the column space of  $\Phi$  denoted by  $\text{span}(\bar{\Phi}_j)$

we are minimizing the perpendicular distance with  $\mathbf{y}$



This means that

$$\bar{\Phi} \hat{\theta} = \bar{\Phi} (\bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T y = \Pi_{\bar{\Phi}} y$$

where  $\Pi_{\bar{\Phi}}$  is a projection on  $\text{Span}_{\mathbb{R}} \bar{\Phi}$

Definition: We say that the linear model

is Gaussian if  $z \sim \mathcal{N}_0(0, 1)$

Then  $y | x \sim \mathcal{N}_0(\theta \phi(x), \sigma^2)$

with density:

$$f(y; \theta, \sigma^2) = (2\pi)^{-n/2} \sigma^{-n} \exp\left(-\frac{1}{2\sigma^2} \|y - \bar{\Phi} \theta\|^2\right)$$

$$= (2\pi)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \|y\|^2 + \frac{1}{\sigma^2} \langle y, \bar{\Phi} \theta \rangle\right)$$

$$- \frac{1}{2\sigma^2} (\bar{\Phi} \theta)^2 - 2n \log \sigma$$

$$= \underbrace{(2\pi)^{-n/2}}_{h(x)} \exp\left(-\frac{1}{2\sigma^2} \|y\|^2 + \left\langle \frac{\bar{\Phi} \theta}{\sigma^2}, \Pi_{\bar{\Phi}} y \right\rangle - \underbrace{\frac{1}{2\sigma^2} \|\bar{\Phi} \theta\|^2 - 2n \log \sigma}_{d(\theta)}\right)$$

$\Rightarrow$  This is a 2-dim exponential family  
with statistics:

$$T_1(Y) = \Pi_\Phi Y$$

$$C_1(\vartheta) = \frac{\Phi \vartheta}{\sigma^2}$$

$$T_2(Y) = \|Y\|^2$$

$$C_2(\vartheta) = -\frac{1}{2\sigma^2}$$

⚡ if  $\Phi$  has full rank then  $T_1, T_2$   
 $T(Y) = (T_1(Y), T_2(Y))$  is sufficient  
 and complete.

Theorem (Gauss-Markov): BLUE

The estimator  $(\Phi^T \Phi)^{-1} \Phi^T Y$  is an unbiased estimator of  $\vartheta$ , optimal among linear unbiased estimators and:

$$\| \Phi (\Phi^T \Phi)^{-1} \Phi^T Y - Y \|^2$$

is an unbiased estimator of  $\sigma^2$ . In the Gaussian case, these estimators are optimal among all unbiased estimators of  $\vartheta$ ,  $(n-k) \sigma^2$

Proof: Let  $VY$  be a linear estimator of  $\vartheta$ .

$$E_{\vartheta, \sigma^2}(VY) = E_{\vartheta, \sigma^2}(V(\Phi \vartheta + \sigma Z)) = V\Phi \vartheta + V\sigma E(Z)$$

$$\text{So } V\Phi \vartheta = \vartheta \Rightarrow V\Phi = I_p \quad \forall \vartheta$$

so any unbiased linear estimator satisfy

$$V\Phi = I_p = \mathbb{1}_p$$

But for the chosen  $V = (\Phi^T \Phi)^{-1} \Phi$

$$\Rightarrow V\Phi = (\Phi^T \Phi)^{-1} \Phi^T \Phi = \mathbb{1}_p$$

$\Rightarrow (\Phi^T \Phi)^{-1} \Phi^T y$  is unbiased. Now we minimize the risk (MSE).

$$\begin{aligned} E_{\theta, \sigma^2}(\|Vy - \theta\|^2) &= E(\|V\Phi\theta + \epsilon z - \theta\|^2) = \sigma^2 E_{\epsilon^2, \theta}(\|Vz\|^2) \\ &= \sigma^2 E\left(\bar{z}_i \left(\sum_l v_{il} z_l\right)^2\right) = \sigma^2 \sum_i E\left(\sum_{j,l} v_{ij} v_{il} z_j z_l\right) \\ &= \sigma^2 \sum_i v_{ij} v_{il} = \sigma^2 \sum_{i,l} v_{il}^2 \end{aligned}$$

Now we have the constrained optimization problem  $E(\|Vy - \theta\|^2)$   
s.t.  $V\Phi = \mathbb{1}_p$

$$\begin{aligned} E(\|Vy - \theta\|^2) - \lambda(V\Phi - \mathbb{1}_p) &= \sigma^2 \|V\|^2 - \lambda(V\Phi - \mathbb{1}_p) \\ \Rightarrow \frac{\partial}{\partial V} & 2\sigma^2 \cdot V = \lambda \Phi^T \Rightarrow V = \frac{\lambda}{2\sigma^2} \Phi^T \end{aligned}$$

but since  $V\Phi = \frac{\lambda}{2\sigma^2} \Phi^T \Phi = \mathbb{1}_p$



$$\Rightarrow \lambda = 2\sigma^2 (\Phi^T \Phi)^{-1}$$

$$\Rightarrow \hat{\gamma} = \frac{2\sigma^2}{2\sigma^2} (\Phi^T \Phi)^{-1} \cdot \Phi^T$$

Consider now the estimator

$$\|\Pi_\Phi Y - Y\|^2 = \|\Pi_\Phi (\Phi \theta + \epsilon Z) - (\Phi \theta + \epsilon Z)\|^2$$

$$= \|(\Pi_\Phi \Phi - \Phi) \theta + \epsilon (\Pi_\Phi Z - Z)\|^2$$

$$= \epsilon^2 \|\Pi_\Phi \cdot Z - Z\|^2$$

if  $\Pi_\Phi = \Pi_{1 \dots k} = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 & \\ & & & & 0 \dots 0 \end{pmatrix}$

$$\mathbb{E} (\|\Pi_\Phi Z - Z\|^2) = \sum_{j=k+1}^n \mathbb{E} (Z_j^2) = n - k$$