

Outline of the lecture

This lecture introduces us to the topic of **supervised learning**. Here the data consists of **input-output** pairs. Inputs are also often referred to as **covariates**, **predictors** and **features**; while outputs are known as **variates** and **labels**. The goal of the lecture is for you to:

- Understand the supervised learning setting.
- Understand linear regression (aka **least squares**)
- Understand how to apply linear regression models to make predictions.
- Learn to derive the least squares estimate by optimization.

Linear supervised learning

- ❑ Many real processes can be **approximated** with linear models.
- ❑ Linear regression often appears as a **module** of larger systems.
- ❑ Linear problems can be solved **analytically**.
- ❑ Linear prediction provides an introduction to many of the **core concepts** of machine learning.

We are given a training dataset of n instances of input-output pairs $\{\mathbf{x}_{1:n}, \mathbf{y}_{1:n}\}$. Each input $\mathbf{x}_i \in \mathbb{R}^{1 \times d}$ is a vector with d attributes. The inputs are also known as predictors or covariates. The output, often referred to as the target, will be assumed to be univariate, $\mathbf{y}_i \in \mathbb{R}$, for now.

$$\mathbf{x}_{1:n} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} \quad \hat{=}$$



A typical dataset with $n = 4$ instances and 2 attributes would look like the following table:

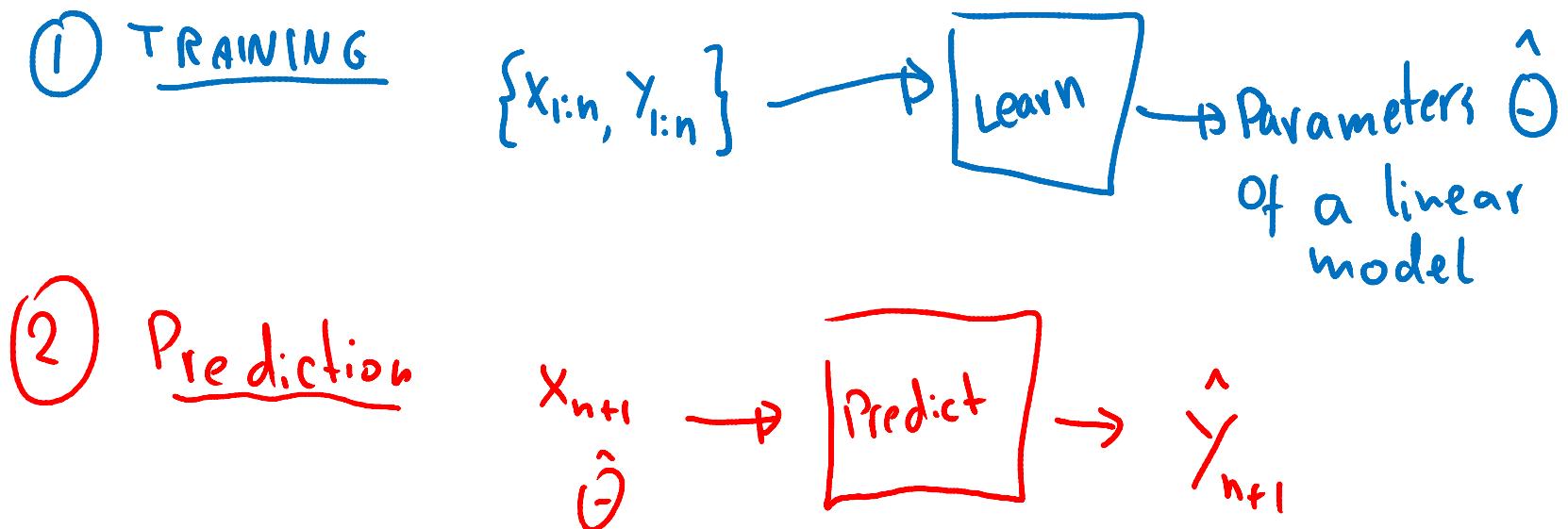
	Wind speed	People inside building	Energy requirement
1	100	2	5
2	50	42	25
3	45	31	22
$n=4$	60	35	18

$$\mathbf{x}_1 = [100 \quad 2] \quad \mathbf{y}_1 = [5]$$

Energy demand prediction



Given the training set $\{\mathbf{x}_{1:n}, \mathbf{y}_{1:n}\}$, we would like to learn a model of how the inputs affect the outputs. Given this model and a new value of the input \mathbf{x}_{n+1} , we can use the model to make a prediction $\hat{y}(\mathbf{x}_{n+1})$.



Prostate cancer example

□ **Goal:** Predict a prostate-specific antigen (log of `lpsa`) from a number of clinical ~~measures~~ in men who are about to receive a radical prostatectomy.



□ The **inputs** are:

- Log cancer volume (`lcavol`)
- Log prostate weight (`lweight`)
- ~~Age~~
- Log of the amount of benign prostatic hyperplasia (`lbph`)
- Seminal vesicle invasion (`svi`) - *binary*
- Log of capsular penetration (`lcp`)
- Gleason score (`gleason`) – *ordered categorical*
- Percent of Gleason scores 4 or 5 (`pgg45`)

Which inputs are more important?

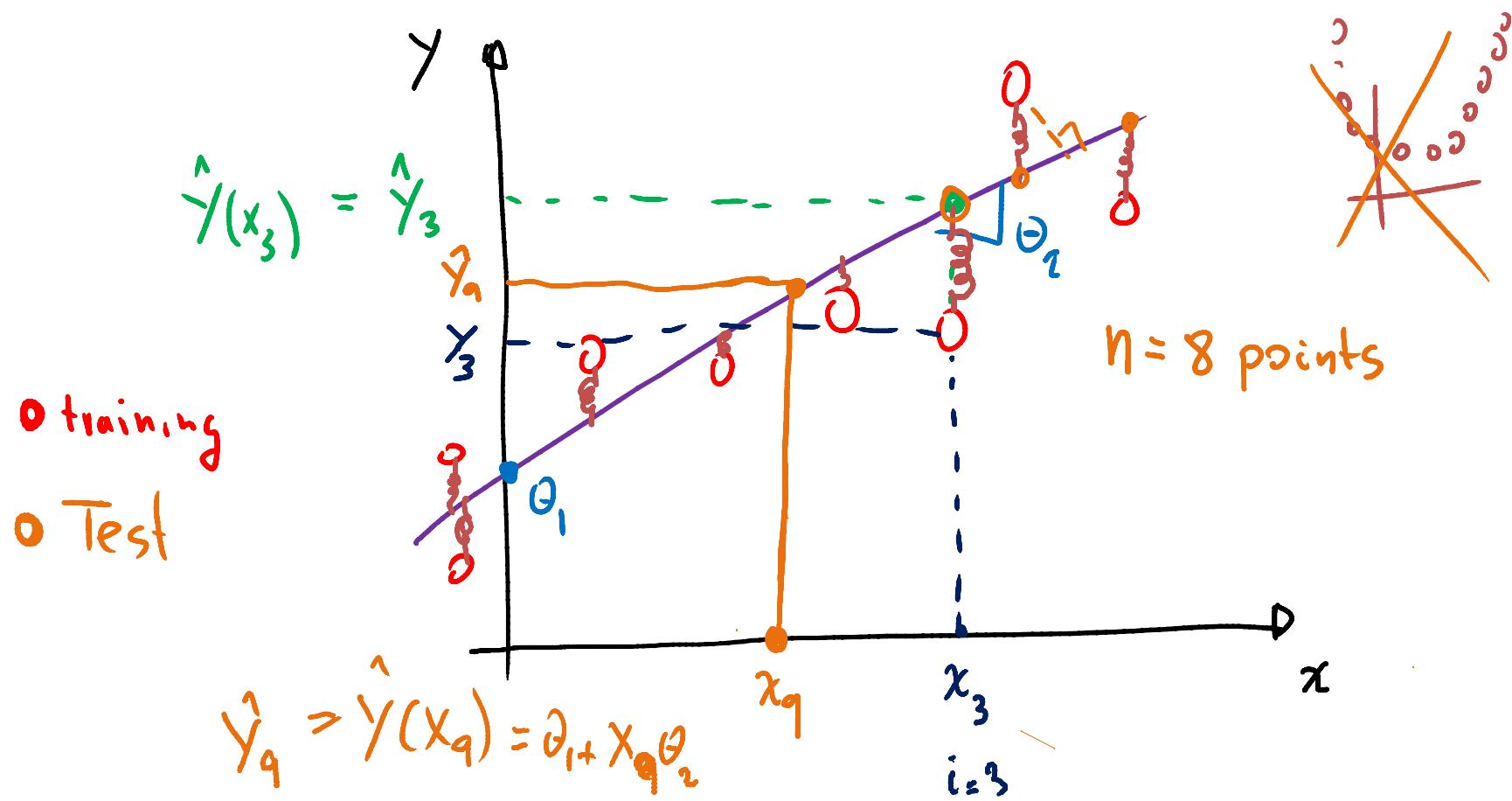
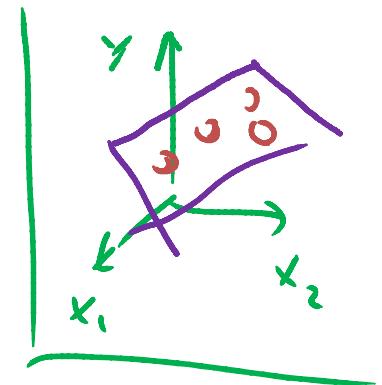


[Hastie, Tibshirani & Friedman book]

$$\hat{y}(\mathbf{x}_i) = \theta_1 + x_i \theta_2 = \hat{x}_i$$

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \theta_1 - x_i \theta_2)^2$$

\nwarrow objective, cost, loss, energy, error function.



Linear prediction

$$\hat{y}_i = \underline{x_{i1}}\theta_1 + x_{i2}\theta_2 \\ = \underline{x_{i1}}\theta_1 + \underline{x_{i2}}\theta_2$$

In general, the linear model is expressed as follows:

$$\hat{y}_i = \sum_{j=1}^d x_{ij}\theta_j,$$

$$i = 1, 2, \dots, n \\ j = 1, 2, \dots, d$$

where we have assumed that $x_{i1} = 1$ so that θ_1 corresponds to the intercept of the line with the vertical axis. θ_1 is known as the bias or offset.

In matrix form, the expression for the linear model is:

$$\hat{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta}, \quad \hat{y}_i = \cancel{x_{i1}}\theta_1 + x_{i2}\theta_2 + \dots + \cancel{x_{id}}\theta_d$$

with $\hat{\mathbf{y}} \in \mathbb{R}^{n \times 1}$, $\mathbf{X} \in \mathbb{R}^{n \times d}$ and $\boldsymbol{\theta} \in \mathbb{R}^{d \times 1}$. That is,

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \boldsymbol{\theta}$$

Wind speed	People inside building	Energy requirement
100	2	5
50	42	25
45	31	22
60	35	18

For our energy prediction example, we would form the following matrices with $n = 4$ and $d = 3$:

$$\mathbf{y} = \begin{bmatrix} 5 \\ 25 \\ 22 \\ 18 \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & 100 & 2 \\ 1 & 50 & 42 \\ 1 & 45 & 31 \\ 1 & 60 & 35 \end{bmatrix}, \quad \boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}.$$

Suppose that $\boldsymbol{\theta} = [1 \ 0 \ 0.5]^T$. Then, by multiplying \mathbf{X} times $\boldsymbol{\theta}$, we would get the following predictions on the training set:

$$\hat{\mathbf{y}} = \begin{bmatrix} 2 \\ 22 \\ 16.5 \\ 18.5 \end{bmatrix} = \begin{bmatrix} 1 & 100 & 2 \\ 1 & 50 & 42 \\ 1 & 45 & 31 \\ 1 & 60 & 35 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix}. \quad \text{Predictions on the training set}$$

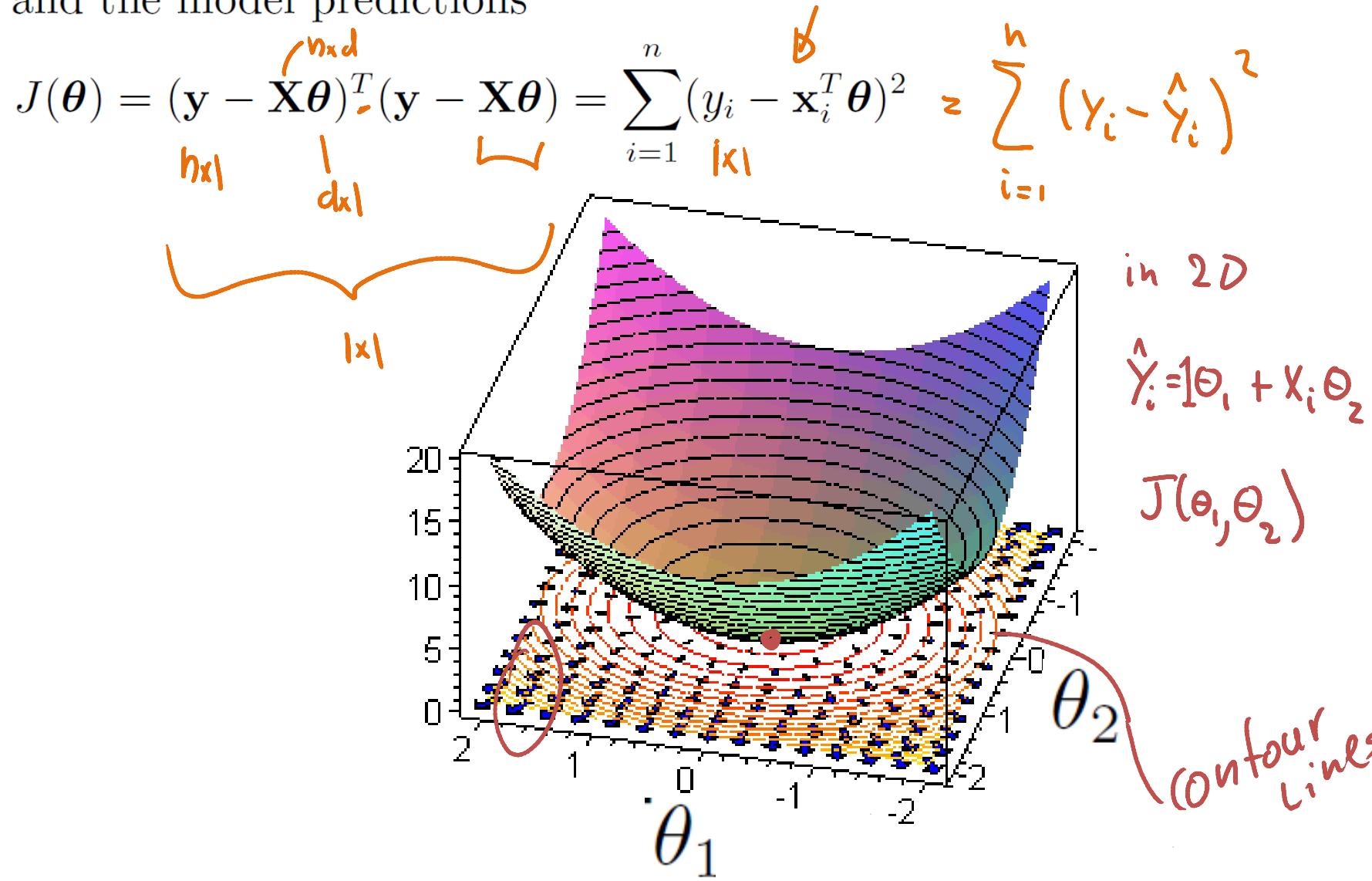
Linear prediction

Likewise, for a point that we have never seen before, say $x = [50 \ 20]$, we generate the following prediction:

$$\hat{y}(x) = [1 \ 50 \ 20] \begin{bmatrix} 1 \\ 0 \\ 0.5 \end{bmatrix} = 1 + 0 + 10 = 11.$$

Optimization approach

Our aim is to minimise the quadratic cost between the output labels and the model predictions



Optimization approach

$$J(\theta) = (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta)^2$$



[Prove]

Optimization: Finding the minimum

$$J(\theta) = (\mathbf{y} - \mathbf{X}\theta)^T(\mathbf{y} - \mathbf{X}\theta) = \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta)^2$$

Suppose $n=3, d=2$

$$\begin{aligned} J(\theta) &= (y_1 - \theta_0 - x_1\theta_1)^2 + (y_2 - \theta_0 - x_2\theta_1)^2 + (y_3 - \theta_0 - x_3\theta_1)^2 \\ \frac{\partial J(\theta)}{\partial \theta_1} &= \sum_{i=1}^3 (y_i - \theta_0 - x_i\theta_1)^2 \\ &\quad \sum_{i=1}^3 2(y_i - \theta_0 - x_i\theta_1)(-1) = -2 \sum_{i=1}^3 (y_i - \theta_0 - x_i\theta_1) \end{aligned}$$

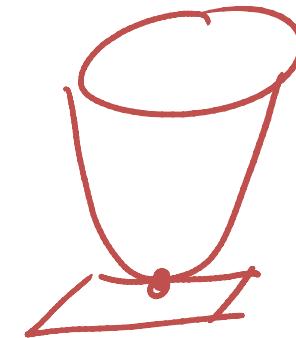
Optimization

$$J(\theta) = \underbrace{(\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta)}_{\text{nx1 } n \times d \text{ d}x \text{l}}$$

We will need the following results from matrix differentiation:

$$\left. \frac{\partial \mathbf{A}\theta}{\partial \theta} \right\} = \mathbf{A}^T \text{ and } \frac{\partial \theta^T \mathbf{A}\theta}{\partial \theta} = 2\mathbf{A}^T \theta$$

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left[\mathbf{y}^T \mathbf{y} + \theta^T \underbrace{\mathbf{X}^T \mathbf{X} \theta}_{\mathbf{A}} - 2 \underbrace{\mathbf{y}^T \mathbf{X} \theta}_{\mathbf{A}' \theta} \right] \\ &= 0 + 2 \mathbf{X}^T \theta - 2 \mathbf{X}^T \mathbf{y} \end{aligned}$$



Least squares estimates

$$2X^T X \theta = 2X^T y \quad \text{Normal eq.}$$

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

L. S.
estimate

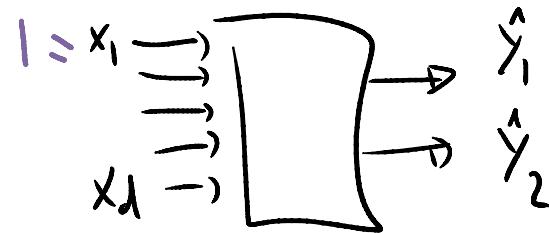
$$\hat{y} = \cancel{X} \hat{\theta} = \underbrace{X (X^T X)^{-1} X^T}_H y = Hy$$

H_{AT}

Multiple outputs

If we have several outputs $\mathbf{y}_i \in \mathbb{R}^c$, our linear regression expression becomes:

e.g. $c=2$



$$\begin{bmatrix} \hat{y}_{11} & \hat{y}_{12} \\ \vdots & \vdots \\ \hat{y}_{n1} & \hat{y}_{n2} \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \Theta_{11} & \Theta_{12} \\ \vdots & \vdots \\ \Theta_{d1} & \Theta_{d2} \end{bmatrix}$$