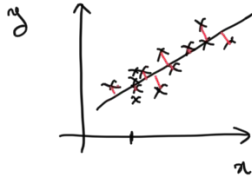
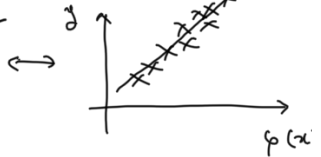
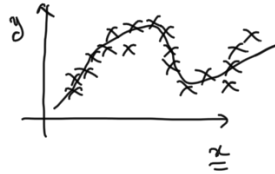


Linear Regression



$$\phi: \mathcal{X} \rightarrow \mathcal{X}'$$



Setup: Given n data-points $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$
 $i = 1, \dots, n$

Goal: Choose/find/"learn" $f: \mathbb{R}^d \rightarrow \mathbb{R}$ to predict the value of y
 for any new given input x .

e.g. $f(x) = w^T x$, $w \in \mathbb{R}^d$

$$= \sum_{i=1}^d w_i x_i$$

↓

ϕ : identity map, i.e.
 $\phi(x) = x$

e.g. $f(x) = w^T \phi(x)$

$$= \sum_{j=1}^m w_j \phi_j(x)$$

$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$
 ↪ features' space with m -features

$$\phi(x) = (\underbrace{\phi_1(x), \phi_2(x), \dots, \phi_m(x)}_{\text{features / basis functions}})$$

⊛ Linear means that the model depends linearly in its parameters.

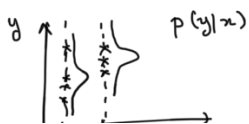
Choosing a model:

Assume an observation rule, typically taking the form:

$$y = f(x) + \epsilon$$


Since y may be corrupted by noise, it is natural to model it as a random variable. Then our goal to characterize its

distribution: $p(y|x)$



We'll have to assume a family $p_\theta(y|x)$

parametrized by θ , and then we'll try to estimate


 x & given the observed data \mathcal{D} .

- Which family to use for $p_\theta(y|x)$? ... our first would be a Gaussian

i.e. $p_\theta(y|x) = \mathcal{N}(y | \mu(x), \sigma^2(x))$, $\theta := \{\mu, \sigma^2\}$

 \downarrow a Gaussian random variable
 \swarrow mean
 \searrow variance

 \hookrightarrow functions in \mathbb{R}^d .

- What $\mu(x)$ and $\sigma^2(x)$ to use?

The simplest choice is to assume that $\sigma^2(x)$ are constants (i.e. independent of x).

$p_\theta(y|x) = \mathcal{N}(y | \mu(x), \sigma^2)$, $\mu(x) = W^T x$, $\sigma^2(x) = \sigma^2$
 unknown parameters $\theta := \{W_1, \dots, W_d, \sigma^2\}$

Recap:

Our starting point was that: $y = f(x) + \varepsilon$

If we assume $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ and $f(x) = W^T x$, then this

implies that: $p(y|x) = \mathcal{N}(y | W^T x, \sigma^2) \iff \boxed{y = \underline{W^T x} + \underline{\varepsilon}, \varepsilon \sim \mathcal{N}(0, \sigma^2)}$
 quantifies the "likelihood" of the observed data

according to the observation model we have assumed

- Modeling choices we've made so far: $\left\{ \begin{array}{l} \text{mean } \mu(x) = W^T x \\ \text{latent function } f(x) \end{array} \right\}$, $\left\{ \varepsilon \sim \mathcal{N}(0, \sigma^2) \right\}$

⊛ How to estimate the unknown model parameters θ ?

... via Maximum Likelihood Estimation (MLE)

Setup: Given $\mathcal{D} := \{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$, $\theta = \{W, \sigma^2\}$

Model: $\boxed{y_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(y_i | \underline{W^T x_i}, \underline{\sigma^2})}$, i.e. our observed outputs are
 i.i.d. = independent and distributed according to a Gaussian

identically distributed

likelihood.

$$\underline{\text{Then}} : \theta_{MLE} = \arg \max_{\theta} p(\mathcal{D}|\theta),$$

The likelihood is : $p(\mathcal{D}|\theta) = p(y|x, \theta)$

$$= p(y_1, y_2, \dots, y_n | x_1, x_2, \dots, x_n, w_1, w_2, \dots, w_d, \sigma^2)$$

joint distribution of all observed outputs, given the corresponding inputs and model parameters!

$$\text{i.i.d.} \Rightarrow \prod_{i=1}^n p(y_i | x_i, \theta) = \prod_{i=1}^n \underbrace{\mathcal{N}(y_i | w^T x_i, \sigma^2)}_{\text{Gaussian pdf}}$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (y_i - w^T x_i)^2 \right]$$

can be re-written in vectorized form :

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp \left[-\frac{1}{2\sigma^2} \underbrace{\sum_{i=1}^n (y_i - w^T x_i)^2}_{\text{sum of squares}} \right] = (y - Xw)^T (y - Xw) \quad \textcircled{1}$$

$$\text{where : } y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}, \quad X = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix}_{n \times d}, \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}_{d \times 1}$$

Recall, we want to compute $\theta_{MLE} = \arg \max_{\theta} p(\mathcal{D}|\theta)$



$$\theta_{MLE} = \arg \min -\log p(\mathcal{D}|\theta)$$

$$\textcircled{1} \Rightarrow -\log p(\mathcal{D}|\theta) = \frac{n}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) := \mathcal{L}(\theta)$$

$\theta := \{w, \sigma\}$

We need to compute the critical points of $\mathcal{L}(w)$

i.e. the points for which $\nabla_w \mathcal{L}(w) = 0$.

Let's focus on estimating w_{MLE} (i.e. for now assume that σ^2 is known)

First notice : $\frac{1}{2}(y - Xw)^T(y - Xw) = \frac{1}{2}(y^T y - y^T Xw - (Xw)^T y + (Xw)^T(Xw))$

$$= \frac{1}{2} y^T y - \underbrace{y^T Xw}_{w^T X^T y} + w^T X^T X w$$

Now we can solve for :

$$\nabla_w L(w) = 0 \Rightarrow \boxed{w_{MLE} = (X^T X)^{-1} X^T y}$$

⊕ $(X^T X)^{-1}$ need to be invertible.
dxd.

$X^T X$ is invertible if the columns of X are linearly independent.

To check whether w_{MLE} is actually a minimizer we can examine the second derivative of our loss $L(w)$:

$$\nabla_w^2 L(w) = X^T X \rightarrow \text{this is a symmetric positive-definite matrix, which implies that indeed}$$

w_{MLE} is a (global) minimizer

Remark :

In the case of linear regression with basis functions :

$$w_{MLE} = (\Phi^T \Phi)^{-1} \Phi^T y$$

, $\varphi : \mathbb{R}^d \rightarrow \mathbb{R}^m$
input space feature space

where

$$\Phi_{n \times m} = \begin{bmatrix} \varphi_1(x_1) & \dots & \varphi_m(x_1) \\ \vdots & & \vdots \\ \varphi_1(x_n) & \dots & \varphi_m(x_n) \end{bmatrix}$$