# Linear regression

**Setup:** Given $\mathcal{D} := \{(x_1, y_1), \ldots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y \in \mathbb{R}$

**Workflow:**

i) Model definition: parametrization, likehood, prior
   $\hookrightarrow \theta \in \mathbb{R}^p$

ii) Training to identify $\theta^* \underset{=}{} \underbrace{(p(\theta|\mathcal{D}))}_{\text{posterior}}$

iii) Perform predictions $f_\theta(x^*)$? $\underbrace{p(f(x^*)|x^*, \mathcal{D})}_{\substack{\text{predictive} \\ \text{posterior}}}$

**1) Model:** $y = \underbrace{f_\theta(x)}_{\text{model form}} + \underbrace{\varepsilon}_{\text{noise}}$ $\left\{ \begin{array}{l} \cdot \ f_\theta(x) = w^T x, \quad w \in \mathbb{R}^d \\ \cdot \ \varepsilon \sim \mathcal{N}(0, \sigma_n^2) \end{array} \right\}$ ✓

**Parameters:** $\theta := \{w_1, \ldots, w_d, \sigma_n^2\}$

**Likehood:** $y_i \overset{\text{i.i.d.}}{\sim} p(y_i | x_i, \theta) = \mathcal{N}(y_i | w^T x_i, \sigma_n^2)$

$$y_1, \ldots, y_n \sim \mathcal{N}(y | Xw, \sigma_n^2 I)$$

$\underset{n \times 1}{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$, $\underset{n \times d}{X} = \begin{bmatrix} x_{11} & \cdots & x_{1d} \\ \vdots & & \\ x_{n1} & \cdots & x_{nd} \end{bmatrix}$, $\underset{d \times 1}{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix}$

**2.) Training:** Goal $\to$ estimate $\theta^*$ that "best" explain $\mathcal{D}$.

$$\theta_{MLE} = \arg\max_\theta p(y | X, w, \sigma^2)$$

**Likelihood:** $p(\mathcal{D}|\theta) = p(y_1, \ldots, y_n | x_1, \ldots, x_n, w_1, \ldots, w_d, \sigma_n^2)$

$$\overset{i.i.d.}{=} \prod_{i=1}^{n} \mathcal{N}(y_i \mid w^T x_i, \sigma_n^2)$$

$$\mathcal{L}(w, \sigma_n^2) := -\log p(y \mid X, w, \sigma^2) \qquad \sum_{i=1}^{n}(y_i - \vec{w}x_i)^2$$

$$= \frac{n}{2} \log(2\pi\sigma_n^2) + \frac{1}{2\sigma_n^2} \overbrace{(y - Xw)^T(y - xw)}$$

Assume that $\sigma_n^2$ is known, hence we will only estimate $W$.

- optimization $\begin{cases} \text{GD:} \quad w^{n+1} = w^n - \eta \nabla_w \mathcal{L}(w^n) \\ \text{Newton:} \quad w^{n+1} = w^n - \eta \, H_n^{-1} \nabla_w \mathcal{L}(w^n) \end{cases}$

$\left. \begin{aligned} \nabla_w \mathcal{L}(w) &:= -X^T y + X^T X w \\ \nabla_w^2 \mathcal{L}(w) &= X^T X \end{aligned} \right\} \Rightarrow$

$\text{GD:} \quad w^{n+1} = w^n - \eta[X^T X w^n - X^T y$

$\text{Newton:} \quad w^{n+1} = w^n - \eta(X^T X)^{-2}[X^T X w^n$

Solve for $W_{MLE}$ analytically:

$$W_{MLE} = \underset{w}{\arg\min} \; -\log p(y \mid X, w) := \mathcal{L}(w)$$

$$\mathcal{L}(w) = \frac{n}{2} \log(2\pi\sigma_n^2) + \frac{1}{2\sigma_n^2} \underbrace{(y - Xw)^T(y - Xw)}_{\text{quadratic form}} \leftarrow$$

Identify critical points:

$$\nabla_w \mathcal{L}(w) = 0$$

- $\frac{1}{2}(y - xw)^T(y - xw) = \frac{1}{2}\left( y^T y - y^T Xw - (xw)^T y + (xw)^T xw \right)$

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

$$\nabla_w \mathcal{L}(w) = -X^T y + X^T X w$$

Condition satisfied by critical points:

$$\nabla_w \mathcal{L}(w) = 0 \implies \boxed{W_{MLE} = \underbrace{(X^T X)^{-1}}_{\substack{\text{needs to be} \\ \text{invertible.}}} X^T y} \to \begin{array}{l}\text{Solution} \\ \text{a} \\ \text{least squ} \\ \text{regressio} \\ \text{proble}\end{array}$$

Observe that $\quad H := \nabla_w^2 \mathcal{L}(w) = \underbrace{X^T X}_{} \to$ symmetric positive-definite

$\implies \quad \mathcal{L}(w)$ is strictly convex in $W$ and $W_{MLE}$ is a unique
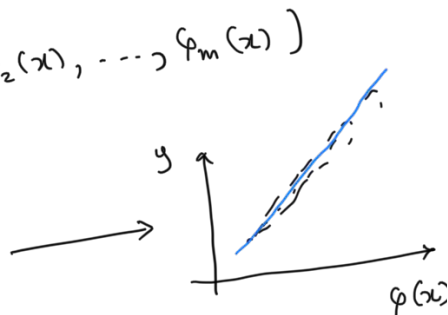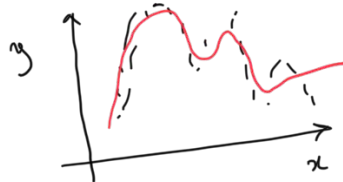
global minimizer.

## Linear regression with basis functions :

Model setup: $\quad y = f_\theta(x) + \varepsilon \quad \begin{cases} \cdot\ f_\theta(x) = w^T \varphi(x) \\ \cdot\ \varepsilon \sim \mathcal{N}(0, \sigma_n^2) \end{cases}$

basis
functions
↗
$\underset{\substack{\downarrow \\ \text{feature} \\ \text{mapping}}}{\varphi} : \underset{\substack{\downarrow \\ \text{input} \\ \text{space}}}{\mathbb{R}^d} \longrightarrow \underset{\substack{\downarrow \\ \text{feature} \\ \text{space}}}{\mathbb{R}^m} \quad , \quad \varphi(x) = (\varphi_1(x), \varphi_2(x), \cdots, \varphi_m(x))$



$$\underset{n \times 1}{y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}} \quad , \quad \underset{n \times m}{\Phi} = \begin{bmatrix} \varphi_1(x_1) & \cdots & \varphi_m(x_1) \\ \vdots & & \vdots \\ \varphi_1(x_n) & \cdots & \varphi_m(x_n) \end{bmatrix} \quad , \quad \underset{m \times 1}{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

2.) MLE for $W$ : $\qquad \underset{m \times 1}{W_{MLE}} = \underbrace{\underset{\substack{m \times n \ \ n \times m}{\underbrace{(\Phi^T \Phi)^{-1}}_{m \times m}}}\ \underset{\substack{m \times n \ \ n \times 1}{\underbrace{\Phi^T y}_{m \times 1}}}$

## Maximum a-posteriori estimation (MAP) :

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

**Model** :  $y = f_\theta(x) + \varepsilon \iff p(y|x, \theta)$ : likelihood

assume a prior  $\theta \sim p(\theta)$ : prior

**Bayes rule** :  $\underbrace{p(\theta|D)}_{posterior} = \dfrac{\overbrace{p(D|\theta)}^{likelihood} \overbrace{p(\theta)}^{prior}}{\underbrace{p(D)}_{marginal\ likelihood\ |\ evid}}$  model / evid

$\hookrightarrow \int p(D|\theta)\, p(\theta)\, d\theta$

**Goal** :  $\theta_{MAP} = \underset{\theta}{\arg\max}\ p(\theta|D)$   $\left(vs\ \theta_{MLE} = \underset{\theta}{\arg\max}\ p(D|\theta)\right)$

**Recall , for linear regression** :

$$\underline{p(w|X,y)} \overset{Bayes}{\propto} \underbrace{p(y|X,w) \cdot p(w)}$$

( omitting $p(y)$ since it does not depend on $\theta$ )

$$\theta_{MAP} = \underset{\theta}{\arg\min}\ \underline{-\log p(w|X,y)}$$

( ✳ We need to assume a prior for $w \sim p(w)$. The simplest choice is to assume  $\boxed{p(w) = \mathcal{N}(0, \frac{-1}{\lambda})}$ )

$$-\log p(w|X,y) = -\underbrace{\log p(y|X,w)} - \log p(w)$$

$w^T w := \|w\|_2^2$ ✓

$$= \frac{n}{2}\log(2\pi\sigma_n^2) + \underbrace{\frac{1}{2\sigma_n^2}(y-Xw)^T(y-Xw)}_{likelihood} + \underbrace{\frac{\lambda}{2} w^T w}_{prior} \quad := \mathcal{L}(w)$$

**Critical points** :

$$\nabla_w \mathcal{L}(w) = 0 \implies W_{MAP} = (X^T X + \lambda I)^{-1} X^T y$$

vs

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

# Comments on MLE vs MAP:

## Pros of MAP:

- easy to compute and interpretable ( interpolates between the MLE and the prio

- It is more resilient against overfitting

- Tends to look like the MLE assymptotically ($n \to \infty$)

## Cons of MAP:

- It is just a point-estimate ( no quantification of uncertainty )

- Unlike the MLE, the MAP is not invariant to re-parametrization.

- Must assume an appropriate prior for $\theta$, possible choices :

$$\begin{cases} \|\theta\|_2 \longleftarrow p(\theta) \sim \mathcal{N}(0, b^{-1}) \to \text{promote "simple" models} \\ \|\theta\|_1 \longleftarrow p(\theta) \sim Lap(b) \to \text{promote "sparsity"} \end{cases}$$