

7 Statistical learning

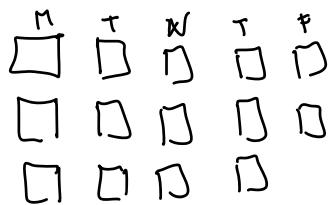
• Introduction

- * Useful for generalization / prediction
- * Concepts:
 - Parameter estimator
 - Bias
 - Variance
- * Point estimation:
 - single best prediction of quantity of interest
 - Example: vector θ — we are interested in estimate the output
- * Parameter γ
- * Point estimate of parameter $\hat{\gamma}$
- * Point estimate or statistic $\hat{\gamma}_m = g(x^{(1)}, \dots, x^{(n)})$ $\{x^{(k)}\}$ i.i.d. data points
- * Good estimator: $\hat{\gamma}_m$ close to γ
- * Frequentist perspective:
 - γ is fixed but unknown
 - $\hat{\gamma} = f(x^{(k)})$. Since $x^{(k)}$ is random $\hat{\gamma}$ is random
- * Function estimation (or approximation):
 - $y = f(x) + \epsilon$
 - \hat{f} is a point estimator
 - Example:
 - Point estimation: $\min \|X^{(\text{test})} \theta - y^{(\text{test})}\|^2$
 - Function estimation: $\min \|f(X^{(\text{test})}) - y^{(\text{test})}\|$ where $f = X^{(\text{test})} \theta$
- * Bias

$$\text{bias}(\hat{\gamma}_m) = \mathbb{E}_{x \sim P(x, \theta)} (\hat{\gamma}_m(x)) - \gamma;$$

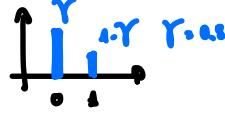


 - $\mathbb{E}(\cdot)$ over the data (samples from a random variable)
 - γ true underlying value of γ used to define data-generating distribution
- * Unbiased: $\text{bias}(\hat{\gamma}_m) = 0 \Rightarrow \mathbb{E}(\hat{\gamma}_m) = \gamma$
- * Asymptotically unbiased: $\lim_{m \rightarrow \infty} \mathbb{E}(\hat{\gamma}_m) = \gamma$



* Example 1: Bernoulli Distribution $\{x^{(1)}, \dots, x^{(m)}\}$ iid with mean γ

$$\varphi(x^{(i)}, \gamma) = \gamma^{x^{(i)}} (1-\gamma)^{1-x^{(i)}}$$



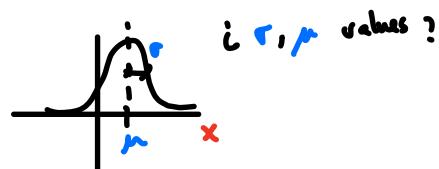
How much is γ : $\hat{\gamma}_m = \frac{1}{m} \sum_i^m x^{(i)}$ (empirical mean)

$$\begin{aligned} \text{bias}(\hat{\gamma}_m) &= E(\hat{\gamma}_m) - \gamma = E\left(\frac{1}{m} \sum_i^m x^{(i)}\right) - \gamma = \frac{1}{m} \sum_i^m E(x^{(i)}) - \gamma = \\ &= \frac{1}{m} \sum_i^m \sum_{x^{(i)}} \varphi(x^{(i)}, \gamma) x^{(i)} - \gamma = \frac{1}{m} \sum_i^m \sum_{x^{(i)}} \gamma^{x^{(i)}} (1-\gamma)^{1-x^{(i)}} x^{(i)} - \gamma = \\ &= \frac{1}{m} \sum_i^m \gamma - \gamma = 0 \end{aligned}$$

* Example 2: Gaussian distribution Estimation of the mean

: $\{x^{(1)}, \dots, x^{(m)}\}$ iid from $\varphi(x^{(i)}) = N(x^{(i)}; \mu, \sigma^2)$

$$\varphi(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}}$$



How much is μ : $\hat{\mu}_m = \frac{1}{m} \sum_i^m x^{(i)}$

$$\text{bias}(\hat{\mu}_m) = E(\hat{\mu}_m) - \mu = E\left(\frac{1}{m} \sum_i^m x^{(i)}\right) - \mu = \frac{1}{m} \sum_i^m E(x^{(i)}) - \mu$$

$$= \frac{1}{m} \sum_i^m \int_{x^{(i)}} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x^{(i)} - \mu)^2}{\sigma^2}} dx^{(i)} - \mu = \frac{1}{m} \sum_i^m \mu - \mu = 0$$

How much is σ : $\hat{\sigma}_m = \sqrt{\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2}$

$$\text{bias}(\hat{\sigma}_m) = E(\hat{\sigma}_m) - \sigma = E\left(\sqrt{\frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2}\right) - \sigma =$$

↑
we will see it
later

* Variance

$$\text{Var}(\hat{\gamma}_m) = E[(\hat{\gamma}_m - E(\hat{\gamma}_m))^2] ; \quad \text{SE}(\hat{\gamma}_m) = \sqrt{\text{Var}(\hat{\gamma}_m)}$$

In ML we typically measure:

$$(\mu_m - 1.96 \text{SE}(\mu_m), \mu_m + 1.96 \text{SE}(\mu_m))$$

It corresponds to the 95% of the interval

* MSE as bias/variance

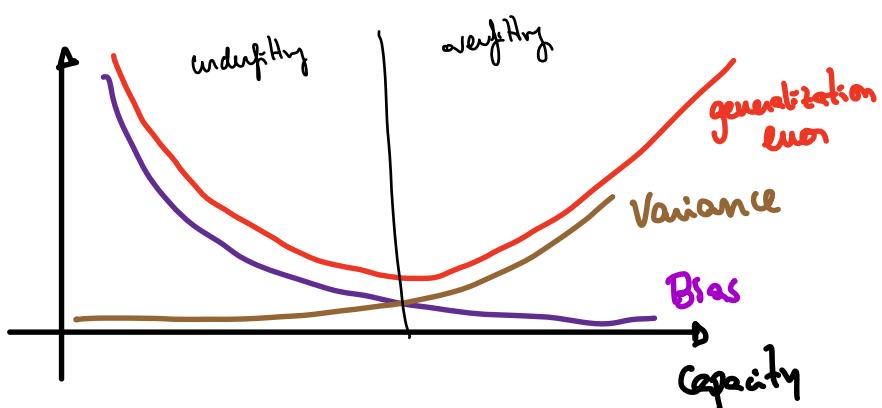
$$\mathbb{E}[(Y_m - Y)^2] = \mathbb{E}(Y_m^2) - 2\mathbb{E}(Y_m)Y + Y^2$$

$$\text{Bias}^2(Y_m, Y) = (\mathbb{E}(Y_m) - Y)^2 = \mathbb{E}(Y_m^2) - 2\mathbb{E}(Y_m)Y + Y^2$$

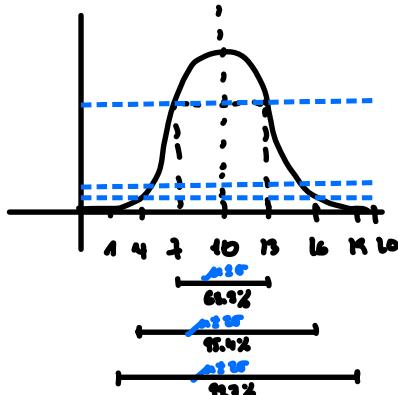
$$\text{Var}(Y_m) = \mathbb{E}(Y_m^2) - \mathbb{E}(Y_m)^2$$

$$\begin{aligned} \text{Bias}^2(Y_m, Y) + \text{Var}(Y_m) &= \mathbb{E}(Y_m^2) - 2\mathbb{E}(Y_m)Y + Y^2 + \mathbb{E}(Y_m^2) - \cancel{\mathbb{E}(Y_m)^2} \\ &= \mathbb{E}(Y_m^2) - 2\mathbb{E}(Y_m)Y + Y^2 = \mathbb{E}[(Y_m - Y)^2] \end{aligned}$$

$$\boxed{\mathbb{E}[(Y_m - Y)^2] = \text{Bias}^2(Y_m, Y) + \text{Var}(Y_m)}$$



* Example $p(x) = N(x|\mu, \sigma^2)$ with $\mu = 10; \sigma^2 = 9$



Likelihood:

$$q(\vec{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_c - \mu)^2}$$

iid

Statistically independent:

$$q(\vec{x}; \mu, \sigma^2) = q(x_1, x_2, \dots, x_n | \mu, \sigma^2)$$

$$= q(x_1; \mu, \sigma^2) q(x_2; \mu, \sigma^2) \dots q(x_n; \mu, \sigma^2)$$

$$= \prod_{i=1}^n q(x_i; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2} = \frac{1}{(2\pi\sigma^2)^n} e^{-\sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}}$$

• Estimators

$$\left\{ \begin{array}{l} \text{Empirical mean: } \hat{\mu}_m = \frac{1}{n} \sum x_i \\ \text{Empirical variance: } \hat{\sigma}_m^2 = \frac{1}{n} \sum (x_i - \hat{\mu}_m)^2 \\ \text{Unbiased empirical variance: } \bar{\sigma}_m^2 = \frac{1}{n-1} \sum (x_i - \hat{\mu}_m)^2 \end{array} \right.$$

Date:

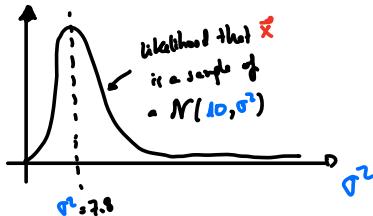
$$x_i = \{10, 12, 7, 5, 11\} \Rightarrow \left\{ \begin{array}{l} \hat{\mu}_m = 9 \\ \hat{\sigma}_m^2 = 6.8 \\ \bar{\sigma}_m^2 = 9.5 \end{array} \right.$$

• We would like:

$$\left\{ \begin{array}{l} E(\hat{\mu}_m) = \mu \Rightarrow \text{small bias}(\hat{\mu}_m) = E(\hat{\mu}_m) - \mu; \text{ small } \text{Var}(\hat{\mu}_m) = E[(\mu - \hat{\mu}_m)^2] \\ E(\hat{\sigma}_m^2) = \sigma^2 \Rightarrow \text{small bias}(\hat{\sigma}_m^2) = E(\hat{\sigma}_m^2) - \sigma^2; \text{ small } \text{Var}(\hat{\sigma}_m^2) = E[(\sigma^2 - \hat{\sigma}_m^2)^2] \\ E(\bar{\sigma}_m^2) = \sigma^2 \Rightarrow \text{small bias}(\bar{\sigma}_m^2) = E(\bar{\sigma}_m^2) - \sigma^2; \text{ small } \text{Var}(\bar{\sigma}_m^2) = E[(\sigma^2 - \bar{\sigma}_m^2)^2] \end{array} \right.$$

Maximum likelihood estimation (MLE)

* Motivation



Suppose fixed \bar{x} and $\mu = 10$; $q(\bar{x}; \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2\sigma^2}(\bar{x}-\mu)^2}$

$$\text{In fact: } \hat{\sigma}_m^2 = \frac{1}{N} \sum (x_i - \bar{x})^2 = 7.8$$

Recall:

$$\begin{cases} \mu = E(x) \\ \sigma^2 = \text{Var}(x) = E[(x-\mu)^2] = E(x^2) - 2E[\underbrace{(x)}_{\mu}] + \mu^2 = E(x^2) - \mu^2 \\ \text{Var}(x) = E(x^2) - E(x)^2 \end{cases}$$

• Suppose μ known and call $s = \sigma^2$

$$\max_s q(\bar{x}; \mu, s) \Rightarrow \frac{\partial q}{\partial s} = 0$$

$$\max_s \log[q(\bar{x}; \mu, s)] \Leftrightarrow \max_s \log[q(\bar{x}; \mu, s)] = \log\left[\frac{1}{(2\pi s)^{1/2}}\right] - \frac{1}{2s} \sum (x_i - \bar{x})^2$$

$$-\frac{1}{2}\left(\frac{1}{s}\right) \ln s + \left(\frac{1}{2\pi s}\right)^{-1/2} + \frac{1}{2s} \sum (x_i - \bar{x})^2 = 0 \Rightarrow \left(\frac{1}{s}\right) \frac{1}{(2\pi s)^{1/2}} + \frac{1}{2s} \sum (x_i - \bar{x})^2 = 0 \Rightarrow \boxed{s = \frac{1}{N} \sum (x_i - \bar{x})^2}$$

$$E(s) = E(\hat{\sigma}_m^2) = \frac{1}{m} \sum E[(x_i - \bar{x})^2] = \frac{1}{m} \sum [E(x_i^2) - 2\mu E(x_i) + \mu^2] = \frac{1}{m} \sum [E(x_i^2) - \mu^2] = \sigma^2 + \mu^2 - \mu^2 = \sigma^2$$

$$E(\hat{\sigma}_m^2) = \sigma^2$$

• Suppose μ, σ^2 unknown

$$\arg \max_{\mu, s} q(\bar{x}; \mu, s) \Leftrightarrow \max_{\mu, s} \log[q(\bar{x}; \mu, s)] \Rightarrow \begin{cases} \frac{\partial \log q(\bar{x}; \mu, s)}{\partial \mu} = 0 \\ \frac{\partial \log q(\bar{x}; \mu, s)}{\partial s} = 0 \end{cases} \Rightarrow (\mu_m, s_m)$$

$$\frac{\partial \log q(\bar{x}; \mu, s)}{\partial \mu} = \frac{1}{s} \log\left[\frac{1}{(2\pi s)^{1/2}} e^{-\frac{1}{2s}(\bar{x}-\mu)^2}\right] = \frac{1}{s} [-\ln[(2\pi s)^2] - \frac{1}{2s}(\bar{x}-\mu)^2] = -\frac{1}{s} \bar{x} + (\bar{x}-\mu) = 0 \Rightarrow \mu_m = \frac{1}{s} \bar{x}$$

$$E(\mu_m) = \frac{1}{m} \sum_i E(x_i) = \frac{1}{m} m \mu = \mu \Rightarrow E(\mu_m) = \mu$$

$$\frac{\partial \log q(\bar{x}; \mu, s)}{\partial s} = \frac{1}{s} [-\ln[(2\pi s)^2] - \frac{1}{2s}(\bar{x}-\mu)^2] = -\frac{1}{(2\pi s)^{1/2}} \sum_i (x_i - \bar{x})^2 - \frac{1}{s} (\bar{x} - \mu)^2 = -\frac{m}{(2\pi s)} + \frac{1}{s} (\bar{x} - \mu)^2$$

$$= -\frac{m}{2s} + \frac{1}{s} (\bar{x} - \mu)^2 = \frac{m}{1s} [\frac{1}{m} (\bar{x} - \mu)^2]^{1-1} = 0 \Rightarrow s_m = \frac{1}{m} (\bar{x} - \mu)^2$$

$$s_m = \frac{1}{m} (\bar{x} - \mu)^2 = \frac{1}{m} [\sum x_i^2 - 2\bar{x} \sum x_i + \frac{m}{m} \mu^2] = \frac{1}{m} [\sum x_i^2 - 2\bar{x} \sum x_i + \frac{m}{m} (\bar{x})^2]$$

$$= \frac{1}{m} [\sum x_i^2 - \underbrace{\frac{m}{m} (\bar{x})^2}_{m(\bar{x})^2}] + \frac{1}{m} [\sum (\bar{x})^2] = \frac{1}{m} [\sum x_i^2 - \frac{1}{m} (\bar{x})^2] = \frac{1}{m} [\sum x_i^2 - \frac{1}{m} (\bar{x})^2]; \quad \boxed{G_m^{-1} = \frac{1}{m} [\sum x_i^2 - \frac{1}{m} (\bar{x})^2]}$$

$$E[(\bar{x} - \mu)^2] = E[\bar{x}^2 + \bar{x}^2 - 2\bar{x} \bar{x}] = m E(x^2) + 2 \sum_i E(x_i) E(\bar{x}) - m(\bar{x}^2) = m(\sigma^2 + \mu^2) + m(m-1)\mu^2 = m[\sigma^2 + m\mu^2]$$

$$E(G_m^{-1}) = \frac{1}{m} E[\sum x_i^2 - \frac{1}{m} (\bar{x})^2] = \frac{1}{m} \left[E(\sum x_i^2) - \frac{1}{m} E[(\bar{x})^2] \right] = \frac{1}{m} [m(\sigma^2 + \mu^2) - \frac{1}{m} m(\sigma^2 + m\mu^2)] = \frac{1}{m} [(m-1)\sigma^2 + (m-n)\mu^2] = \frac{m-1}{m} \sigma^2$$

$$\boxed{E(G_m^{-1}) = \frac{m-1}{m} \sigma^2}$$

MLE: Maximum likelihood estimator

- Idea: Pick θ that assign the maximum probability to the training data D
- Assumption: iid assumption: $p(D|\theta) = \prod_{n=1}^N p(y_n|x_n; \theta)$
- Definition: $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} p(D|\theta) = \arg \max_{\theta} \log(p(D|\theta))$
 Thus, $\hat{\theta}_{\text{MLE}} = \arg \max_{\theta} \sum_{i=1}^n \log(p(y_i|x_i; \theta)) = \arg \min_{\theta} -\sum_{i=1}^n \log(p(y_i|x_i; \theta))$
 (conditional) Negative log likelihood (NLL)

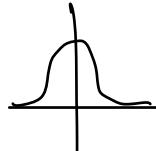
Example: MLE for linear regression

$$p(y|\theta) = N(y|\theta^T x, \sigma^2). \text{ Assume } \sigma^2 \text{ fixed}$$

$$\begin{aligned} \arg \min \text{NLL}(\theta) &= -\sum_{n=1}^N \log \left[\frac{1}{(2\pi\sigma^2)^{1/2}} e^{-\frac{1}{2\sigma^2}(y_n - \theta^T x)^2} \right] \\ &= \sum_{n=1}^N \log \cancel{(2\pi\sigma^2)^{1/2}} + \frac{1}{2\sigma^2} (y_n - \theta^T x)^2 \end{aligned}$$

$$\arg \min \text{RSS}(\theta) = \sum_{n=1}^N (y_n - \theta^T x)^2 \text{ residual sum of squares} =$$

$$\begin{aligned} \arg \min \text{MSE}(\theta) &= \frac{1}{N} \text{RSS}(\theta) = \frac{1}{N} \sum_{n=1}^N (y_n - \theta^T x)^2 && \text{Ordinary least square} \\ \Rightarrow \frac{1}{N} (\mathbf{x}_w - \mathbf{y})^T (\mathbf{x}_w - \mathbf{y}) &\Rightarrow \theta^* = (\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y} \end{aligned}$$



Example Logistic regression model

$$\text{Given } y \in \{0, 1\}, \text{ Bernoulli: } \Pr(y|X; \theta) = g(h_\theta(x))^y (1 - g(h_\theta(x)))^{1-y}$$

$$\begin{aligned} \text{likelihood (iid)}: L(\theta|y, x) &= \Pr(y|X; \theta) = \prod_i \Pr(y_i|X_i; \theta) \\ &= \prod_i g(h_\theta(x_i))^y (1 - g(h_\theta(x_i)))^{1-y} \end{aligned}$$

$$\arg \max_{\theta} \log \left[\prod_i g(h_\theta(x_i))^y (1 - g(h_\theta(x_i)))^{1-y} \right] =$$

$$\arg \min_{\theta} -\sum_{i=1}^n (1-y_i) \cdot \log(1 - g(h_\theta(x_i))) + y_i \cdot \log(g(h_\theta(x_i)))$$

Maximum a posteriori probability estimation (MAP)

* Bayesian version of MLE + a prior probability density on the parameter θ

* $\begin{cases} \theta : \text{vector to be estimated} \\ D = \{X, Y\} : \text{observations} \end{cases}$ are random variables with joint probability density $p(D, \theta)$

* Statistical estimation vs Bayesian

$\begin{cases} \theta \text{ is an unknown parameter} \\ \hat{\theta} \text{ estimate is random} \\ \text{since } D \text{ is seen as random} \end{cases}$ $\begin{cases} \theta \text{ is a random variable} \\ D \text{ is observed (not random)} \end{cases}$

* Prior density: $\begin{cases} p(\theta) = \int p(D, \theta) dD \\ p(D) = \int p(D, \theta) d\theta \end{cases}$ { Prior info about vector θ before we observe D

{ Prior info about the measurements or observations D .

* Conditional density of D given θ : $\begin{cases} p(D|\theta) = \frac{p(D, \theta)}{p(\theta)} \\ \text{Plays the role of probability in the MLE} \end{cases}$

* Conditional density of θ given D : $\underbrace{p(\theta|D)}_{\text{Posterior density of } \theta} = \frac{p(D, \theta)}{p(D)} = \frac{p(D|\theta) \cdot p(\theta)}{p(D)}$ = Bayes rule

Bayes theorem

$$\text{joint probability} = p(H, E) \quad E = \text{evidence} \\ H = \text{hypothesis}$$

Error: $p(H)$ = probability an hypothesis is true (before any evidence)

$p(E)$ = probability of seeing the evidence $p(D|E)$

Likelihood: $p(E|H)$ = probability of seeing the evidence if Hypothesis is true

Posterior: $p(H|E)$ = probability an hypothesis is true given some evidence

Bayes theorem:

$$p(H|E) = p(E|H) \frac{p(H)}{p(E)}$$

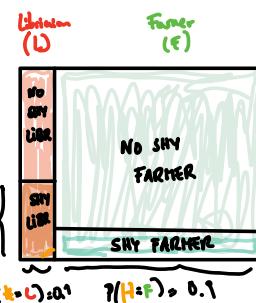
* Farmers/librarian example

		$H=L$	$H=F$
$E=0$	$H=L$	$p(E=0, H=L) = 0.1$	$p(E=0, H=F) = 0.9$
	$H=F$	$(0.1)(0.1)$	$(0.9)(0.9)$
$E=1$	$H=L$	$p(E=1, H=L) = 0.4$	$p(E=1, H=F) = 0.6$
	$H=F$	$(0.4)(0.5)$	$(0.6)(0.5)$

joint probability

$$p(E=1, H=L) = (0.4)(0.1) = 0.04$$

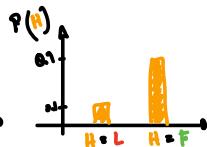
$$p(H=L) = 0.1 \quad p(H=F) = 0.9$$



Marginals

$$P(E) = 0.52$$

$$P(H) = 0.5$$



+ Likelihood (Conditional probability)

a) Probability of being shy $E=1$ when being a liberal $H=L$

$$p(E=1|H=L) = \frac{p(E=1, H=L)}{p(H=L)} = 0.4$$

b) Probability of being shy $E=1$ when being a farmer $H=F$

$$p(E=1|H=F) = \frac{p(E=1, H=F)}{p(H=F)} = 0.6$$

* Posterior

a) Probability of being a liberal $H=L$ when being shy $E=1$

$$p(H=L|E=1) = \frac{p(E=1, H=L)}{p(E=1)} = \frac{(0.4)(0.1)}{0.13} = 0.307$$

b) Probability of being a farmer $H=F$ when being shy $E=1$

$$p(H=F|E=1) = \frac{p(E=1, H=F)}{p(E=1)} = \frac{(0.6)(0.1)}{0.13} = 0.462$$

Remark

Typically the interest is to know the posterior $p(H|E=1)$ by knowing $p(E=1|H)$ and $p(H)$ and $p(E=1)$

- considering that we know the probability of being shy when is a liberal $p(E=1|H=L)$, knowing the probability of being a liberal $p(H=L)$ and the probability of being shy $p(E=1)$, the

? what is the probability to be a liberal if the person is shy? i.e. $p(H=L|E=1)$?

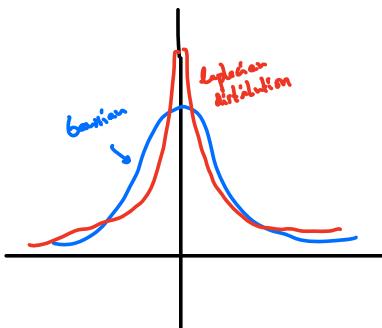
* Our estimate of θ , given the observation D :
$$\hat{\theta}_{\text{map}} = \arg \max_{\theta} p(\theta | D; \theta, D)$$

We take as estimate of θ , the value that maximizes the conditional density of θ given the observed value of D

* The difference with MLE is $p_x(\theta)$ is now considered (we consider the prior info of x)
 { If the prior is uniform MAP \equiv MLE }

* Taking logarithms :
$$\hat{\theta}_{\text{map}} = \arg \max_{\theta} \underbrace{\log(p(D|\theta; \theta, D))}_{\text{Same term as the log-likelihood}} + \underbrace{\log(p_x(\theta))}_{\substack{\text{Penalization of choices of } \theta \\ \text{that are very unlikely to happen}}}$$

$$= \arg \min_{\theta} -\log(p(D|\theta; \theta, D)) - \underbrace{\log(p_x(\theta))}_{\text{penalty}}$$



Examples :

Gaussian:	$p_x(\theta) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\theta^2}{2\sigma^2}\right)$
with zero mean	$-\log[p_x(\theta)] \approx \frac{\theta^2}{2\sigma^2}$ L^2 regularization
Laplacean:	$p_x(\theta) = \frac{1}{(2a)} \exp\left(\frac{ \theta }{a}\right)$
	$-\log[p_x(\theta)] \approx \frac{ \theta }{a}$ L^1 regularization