# MLE, MAP, & NB

MLE = MAXIMUM LIKELIHOOD ESTIMATION

MAP = MAXIMUM A POSTERIORI

NB = NAIVE BAYES'

$\mathcal{D}$ = SET OF DATA GENERATED FROM SOME DISTRIBUTION PARAMETRIZED BY $\theta$.

WE WANT TO ESTIMATE $\theta$

MLE & MAP ARE PARAMETER ESTIMATION METHODS.

# MAXIMUM LIKELIHOOD ESTIMATION (MLE)

IN THE MLE METHOD, WE WANT TO FIND $\theta$ THAT BEST EXPLAINS THE DATA.

$\Rightarrow$ WE MAXIMIZE $P(\mathcal{X}|\theta)$

DENOTE SUCH VALUE BY $\hat{\theta}_{ML}$

$$\hat{\theta}_{ML} = \underset{\theta}{\text{ARG MAX}}\ P(\mathcal{X}|\theta)$$

IF THE SET OF DATA POINTS $\mathcal{X} = (x_1, x_2, \ldots, x_n)$, THEN

$$\hat{\theta}_{ML} = \underset{\theta}{\text{ARG MAX}}\ P(x_1, x_2, \ldots, x_n|\theta)$$

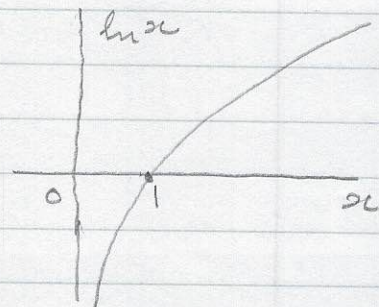IF OBSERVATIONS ARE <u>INDEPENDENT</u>: $P(x_1, x_2, \ldots, x_n) = \prod$

$$P(x_1, x_2, \ldots, x_n|\theta) = \prod_{i=1}^{m} P(x_i|\theta)$$

AS LOGARITHM IS A MONOTONICALLY INCREASING FUNCTION.

$$\hat{\theta}_{ML} = \underset{\theta}{\text{ARG MAX}}\ \ln \prod_{i=1}^{m} P(x_i|\theta)$$

$$= \underset{\theta}{\text{ARG MAX}}\ \sum_{i=1}^{m} \ln P(x_i|\theta)$$

IF WE KNOW $P(\cdot|\cdot)$ WE CAN SOLVE THIS BY TAKING DERIVATIVES WITH RESPECT TO $\theta$ AND MAKING IT EQUAL TO ZERO.

# MAXIMUM A POSTERIORI (MAP)

$\theta$ = ASSUMPTIONS (MODE, HYPOTHESIS)

$\alpha$ = OBSERVATIONS (DATA)

$P(\theta|\alpha)$ = POSTERIOR DISTRIBUTION
= POSTERIORI BELIEF GIVEN OBSERVATIONS

$P(\alpha|\theta)$ = LIKELIHOOD OF OBSERVATIONS, GIVEN THE MODEL

$P(\theta)$ = PRIOR DISTRIBUTION
= PRIOR BELIEF

$P(\alpha)$ = MARGINAL LIKELIHOOD
= EVIDENCE : USED AS A NORMALIZATION FACTOR

$$P(\theta|\alpha) = \frac{P(\alpha|\theta)\,P(\theta)}{P(\alpha)} \quad ; \quad P(\alpha) > 0$$

LIKELIHOOD OF DATA, GIVEN MODEL

PRIOR BELIEF

$$P(\text{MODEL}|\text{DATA}) = \frac{P(\text{DATA}|\text{MODEL})\;P(\text{MODEL})}{P(\text{DATA})}$$

EVIDENCE

POSTERIOR BELIEF
GIVEN DATA

WE WANT TO FIND MOST LIKELY VALUE OF $\theta$ GIVEN $\alpha$.

THIS IS $\arg\max\limits_{\theta} P(\theta|\alpha)$

MAXIMUM A-POSTERIORI (MAP) ESTIMATE IS DEFINED AS

$$\widehat{\theta}_{\text{MAP}} = \arg\max\limits_{\theta} P(\theta|\alpha)$$

MODEL       DATA

## NAIVE BAYES' CLASSIFIER

$\Omega = \{\omega_1, \omega_2, \ldots, \omega_m\}$ = SET OF LABELS FOR DIFFERENT
    CLASSES OF DATA

$P(\omega_i)$ = PROBABILITY THAT A DATA POINT BELONGS TO CLASS $\omega_i \in \Omega$
    = A PRIORI PROBABILITY

$R$ = FEATURE SPACE = SET OF ALL DATA POINTS

$\chi$ = SET OF DATA POINTS, $\chi \subseteq R$

$\tilde{X}$ = RANDOM DATA VECTOR

$x \in \tilde{X}$ IS AN INSTANCE OF RANDOM DATA VECTOR

$p(x)$ = PROBABILITY DENSITY FUNCTION OF THE RANDOM
    DATA VECTOR $\tilde{X}$

$$p(\omega_i | x) = \frac{p(x | \omega_i) P(\omega_i)}{p(x)} \; ; \; p(x) > 0 \; ; \; 1 \leq i \leq m$$

$$p(x) = \sum_{i=1}^{m} p(x | \omega_i) P(\omega_i)$$

TEST DATA POINT $x \in R \backslash \chi$ IS ASSIGNED TO CLASS $\omega_i$

IF $P(\omega_i | x) > P(\omega_j | x) \quad \forall j \neq i \; ; \; 1 \leq j \leq m$

THIS SCHEME MINIMIZES CLASSIFICATION ERROR
PROBABILITY.

THIS IS CALLED NB CLASSIFIER, BECAUSE IF $x_a \in R$ IS
$(x_{a_1}, x_{a_2}, \ldots, x_{a_t})$ IT IS ASSUMED THAT

$$p(x_a | \omega_j) = \prod_{k=1}^{t} p(x_{a_k} | \omega_j) \; ; \; 1 \leq j \leq m$$

## DISCUSSION

1. MLE, MAP, AND NB ARE ALL CONNECTED
   - MLE & MAP ARE PARAMETER ESTIMATION METHODS.
   - NB IS A CLASSIFIER THAT PREDICTS THE PROBABILITY OF THE CLASS THAT THE EXAMPLE BELONGS TO.

2. MLE DOES NOT ALLOW US TO 'INJECT' OUR BELIEFS ABOUT THE LIKELY VALUES FOR THE PARAMETER (PRIOR) IN THE ESTIMATION PROCESS.

3. MAP ALLOWS THE FACT THAT THE PARAMETER CAN TAKE VALUES FROM A PRIORI (NON-UNIFORM) DISTRIBUTION THAT EXPRESS OUR PRIOR BELIEFS REGARDING THE PARAMETERS.

4. MAP RETURNS PARAMETER VALUE, WHERE THE PROBABILITY IS HIGHEST FOR GIVEN DATA.

5. MLE AND MAP, EACH RETURN A SINGLE SPECIFIC VALUE FOR THE PARAMETER.

   NB COMPUTES THE FULL POSTERIOR DISTRIBUTION $P(\theta | \mathcal{D})$.

## OPINION

- MLE CAN BE REGARDED AS TRADITIONAL "FREQUENTIST" THINKING
- MAP CAN BE REGARDED AS BAYESIAN (AS IT IS A DIRECT APPLICATION OF BAYES' THEOREM).