

## PROBABILITY

OFTEN THE GOAL OF MACHINE LEARNING IS TO DETERMINE THE PROBABILITY OF AN EVENT

## EXAMPLES

- $P_X(\text{YEAR THAT POLAR ICE CAP MELTS} \leq 2030)$
- $P_X(\text{A NEW EMAIL IS SPAM})$
- $P_X(\text{A PERSON IS AT RISK FOR A DISEASE})$

□



## STATISTICS

THERE ARE TWO MAIN PARADIGMS IN STATISTICS:

1. FREQUENTISTS
2. BAYESIAN

### 1. FREQUENTISTS

— PROBABILITIES ARE LONG RUN FREQUENCIES.

EXAMPLE: FLIP A COIN MILLION TIMES TO DETERMINE IF ITS FAIR

— PROBABILITY OF AN EVENT IS DEFINED AS THE FREQUENCY THAT SPECIFIC EVENT OCCURRED IN A LONG LIST OF REPEATED TRIALS.

### 2. BAYESIAN

— PROBABILITIES QUANTIFY OUR UNCERTAINTY IN EVENTS DESIGNED TO GET THE CLOSEST TO THE TRUTH GIVEN A SPECIFIC SET OF DATA.

— WE HAVE TO WORK WITH THE GIVEN SET OF DATA POINTS. WE DO NOT HAVE THE LUXURY OF OBSERVING MULTIPLE TRIALS. WE THEREFORE REQUIRE THE ABILITY TO GENERALIZE FROM A SMALL NUMBER OF EVENTS.



## BAYESIAN STATISTICS

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

WE WANT TO ESTIMATE  $\theta$

THAT IS, WE WANT TO FIND MOST LIKELY VALUE OF  $\theta$  GIVEN  $\mathcal{D}$ .

THIS IS  $\underset{\theta}{\text{ARG MAX}} P(\theta|\mathcal{D})$

AS PER BAYES' RULE

$$P(\theta|\mathcal{D}) = \frac{P(\mathcal{D}|\theta)P(\theta)}{P(\mathcal{D})} ; P(\mathcal{D}) > 0$$

THE TERMS HAVE THE FOLLOWING MEANING.

- $P(\theta|\mathcal{D})$  = POSTERIOR DISTRIBUTION OF  $\theta$  GIVEN  $\mathcal{D}$
- $P(\mathcal{D}|\theta)$  = LIKELIHOOD OF OBSERVATIONS, GIVEN THE MODEL
- $P(\theta)$  = PRIOR DISTRIBUTION OR BELIEF
- $P(\mathcal{D})$  = EVIDENCE, OR MARGINAL LIKELIHOOD OF  $\mathcal{D}$ .  
IT IS A CONSTANT WITH RESPECT TO  $\theta$ .



-  $P(\theta|x)$  = POSTERIOR DISTRIBUTION OF  $\theta$  GIVEN  $x$  AND CALCULATING (OR APPROXIMATING) IS THE MAIN GOAL OF BAYESIAN INFERENCE

IT EXPRESSES UNCERTAINTY ABOUT  $\theta$ , AFTER WE HAVE SEEN ("LEARNED FROM") THE DATA. THERE IS STILL UNCERTAINTY BECAUSE WHILE THE DATA TOLD US SOMETHING ABOUT  $\theta$ , WE STILL DO NOT KNOW EVERYTHING ABOUT IT.

-  $P(x|\theta)$  = LIKELIHOOD. IT IS A FUNCTION OF  $\theta$ .

THIS IS A WAY OF DESCRIBING THE PROBABILITY OF THE DATA AS A FUNCTION OF THESE UNKNOWN PARAMETERS.

-  $P(\theta)$  = PRIOR DISTRIBUTION. IT DESCRIBES OUR BELIEF ABOUT THE QUANTITY OF INTEREST BEFORE WE SEE DATA. OFTEN TIMES WE USE WHAT ARE CALLED CONJUGATE PRIORS TO THE LIKELIHOOD TO DESCRIBE OUR A PRIORI BELIEFS.

-  $P(x)$  = MARGINAL LIKELIHOOD OF  $x$ , AND IT IS CONSTANT WITH RESPECT TO  $\theta$ .

THEREFORE WE WRITE  $P(\theta|x) \propto P(x|\theta)P(\theta)$   
↑  
PROPORTIONALITY SYMBOL

$$P(x) = \int_{\theta} P(x|\theta)P(\theta)d\theta$$

THIS IS CALLED THE MARGINAL LIKELIHOOD, BECAUSE WE ARE MARGINALISING (OR AVERAGING) OVER  $\theta \in \Theta$

LOWER-CASE  $\theta$  → UPPER-CASE  $\Theta$