

# Machine Learning and Data Mining - Notes

Philip Warton

September 29, 2021

## 1 9/29

### 1.1 Lecture 1.2: Statistical Learning - MLE / MAP

Jhomas.

#### 1.1.1 Probability

**Definition 1.1.** A sample space  $\Omega$  is a set of all possible outcomes.

**Definition 1.2.** An event  $A$  is a subset of  $\Omega$ . That is,  $A \subset \Omega$ .

A probability must be non-negative for any event. Must be 1 for the entire sample space, 0 for the empty set, and must not be double-counting.

Marginalization:

$$P(A) = \sum_{b \in \text{Val}(B)} P(A, B = b) \quad (\text{discrete}) \quad (1)$$

$$P(A) = \int_{b \in \text{Val}(B)} P(A, B = b) \quad (\text{continuous}) \quad (2)$$

Conditional Distribution:

$$P(A | B) = \frac{P(A, B)}{P(B)} \quad (3)$$

Chain Rule:

$$P(A, B) = P(A | B)P(B) = P(B | A)P(A) \quad (4)$$

Bayes Rule:

$$P(A | B) = \frac{P(B | A)P(A)}{P(B)} \quad (5)$$

**Definition 1.3.** A random variable  $X$  is a mapping between events in  $\Omega$  to numbers. They can be discrete or continuous.

A probability density describes the mapping from values of a random variable  $X$  to probabilities. Some common discrete distributions are the following:

Bernoulli:  $p_X(x) = \theta^x(1 - \theta)^{(1-x)}$  (6)

Categorical:  $p_X(x) = \theta_x$  (7)

Common continuous distributions:

Gaussian:  $f_X(x) = \frac{1}{\sigma\sqrt{2\pi}}e^{-\frac{(x-\mu)^2}{2\sigma^2}}$  (8)

**Definition 1.4.** The expectation of a random variable is given by

$$E_X[g(x)] = \int_{x \in \text{Val}(X)} f_x(x)g(x)dx$$

### 1.1.2 MLE Algorithm

There are two steps to maximum likelihood estimation:

- (i) Assume a probabilistic model of how the data was generated
- (ii) Find  $\hat{\theta}_{MLE}$  that maximizes the probability (or likelihood) of generating the training data under the probabilistic model.