# Artificial Intelligence

A Course About Foundations



# Probabilistic Reasoning: Supervised Learning

Marco Piastra

Artificial Intelligence 2024–2025 Supervised Learning [1]

# Machine Learning

Artificial Intelligence 2024–2025 Supervised Learning [2]

# Types of machine learning problems

Consider a number of observations (i.e. a dataset) made by an agent  $\{d^{(1)}, d^{(2)}, ..., d^{(N)}\}$ 

#### Supervised learning

Learning form <u>complete</u> observations: each of the observations  $\{d^{(1)}, d^{(2)}, ..., d^{(N)}\}$  include values for <u>all</u> the random variables in the model

The objective is learning a <u>distribution P</u>

#### Unsupervised learning

Learning form <u>incomplete</u> observations: observations  $\{d^{(1)}, d^{(2)}, ..., d^{(N)}\}$  do <u>not</u> necessarily include values for all the random variables in the model The objective is learning a *distribution P* 

#### ■ Reinforcement learning

The observations  $\{d^{(1)}, d^{(2)}, ..., d^{(N)}\}$  are states o situations, at each state  $S^{(i)}$  the agent must perform an **action**  $a_i$  that produces a **result**  $r_i$ .

The objective is learning a distribution  $\pi$  over possible actions in each state which describes the policy that the agent will follow

Such policy should maximize the expected value of a reward function  $v(< r_1, r_2, ..., r_n >)$  of the sequence of results

### Observations and Independence

Each observation could be the outcome of an experiment or a test

The outcome of a particular experiment can be represented by a set of *random variables* 

For example, if the model makes use of the two random variables  $\{X, Y\}$ , the N outcomes of the experiments are  $d^{(1)} = (X^{(1)}, Y^{(1)}), \dots, d^{(N)} = (X^{(N)}, Y^{(N)})$ 

That is, a dataset

a dataset 
$$D:=\{d^{(i)}\}_{i=1}^N\;,\qquad d^{(i)}=\{X_1^{(i)},\ldots,X_n^{(i)}\}^{\textit{may involve several variables}}$$

#### Independent observations, same probability distribution

*Independent and Identically Distributed* (IID) random variables

Definition

Random variables in a set  $\{X^{(1)}, X^{(2)}, \dots, X^{(N)}\}$  are *Independent and Identically Distributed* (IID) iff:

$$\langle X^{(i)} \perp X^{(j)} \rangle, \ \forall i \neq j$$
 (independence)

$$P(X^{(i)} = x) = P(X^{(j)} = x), \ \forall i \neq j, \ \forall x$$
 (identical distribution)

CAUTION: Being IID is not an obvious property of observations

e.g. different measurements on different patients <u>may</u> be IID, but different measurements over time on the same patient are <u>not</u> IID

### ML = Representation + Evaluation + Optimization

#### Representation

The objective is learning a specific distribution

$$P({X_r};\theta)$$

where  $\{X_r\}$  are all the random variables of interest and  $\theta$  is a *set* of parameters

Which kind of distribution (i.e. the *model* or also the *learner*) do we select?

Example: assume we select the anti-spam filter (i.e. Naïve Bayesian Classifier) as the model the parameters in such case are the numerical probabilities in the CPTs

#### Evaluation

Given a dataset D, how well does a specific set of parameter values  $\hat{\theta}$  make the distribution P fit the dataset?

An estimator, i.e. a scoring function of some sort, must be selected

#### Optimization

How can we find the optimal set of parameter values  $\theta^*$  with respect to the *estimator* of choice?

In general, this is an optimization problem

# Maximum Likelihood Estimator (MLE)

Artificial Intelligence 2024-2025

### Likelihood

A probabilistic model P(X), with parameters  $\theta$ 

 $\theta$  is a set of values that characterizes P(X) completely: once  $\theta$  is defined, P(X) is also defined.

A set of IID observations (data items)  $D = \{d^{(1)}, d^{(2)}, ..., d^{(N)}\}$ 

#### Likelihood function (or conditional probability)

A function, or a conditional probability, derived from the model P(X)

$$L(\theta \mid D) = P(D \mid \theta) = P(d^{(1)}, \dots, d^{(N)} \mid \theta)$$

*Note the 'trick':* 

where  $P(D \mid \theta)$  is the conditional probability that the parameter  $\theta$ , considered as a random variables, could <u>generate</u> the observations D

likelihood of the dataset given the parameters

When the observations  $\{D^{(1)}, \dots, D^{(N)}\}$  are IID:

$$P(D \mid \theta) = P(d^{(1)} \mid \theta) \dots P(d^{(N)} \mid \theta) = \prod_{m} P(d^{(m)} \mid \theta)$$

Artificial Intelligence 2024–2025 Supervised Learning [7]

### Maximum Likelihood Estimator (MLE)

A probabilistic model P(X), with parameters  $\theta$ 

 $\theta$  is a set of values that characterizes P(X) completely: once  $\theta$  is defined, P(X) is also defined.

A set of IID observations (data items)  $D = \{d^{(1)}, ..., d^{(N)}\}$ 

#### Maximum Likelihood Estimation

$$\theta_{ML}^* := \operatorname{argmax}_{\theta} L(\theta|D)$$

Since the observations are IID, using *log-Likelihood* could ease computations:

$$\ell(\theta \mid D) = \log L(\theta \mid D) = \log \prod_{m} P(d^{(m)} \mid \theta) = \sum_{m} \log P(d^{(m)} \mid \theta)$$

$$\theta_{ML}^* = \operatorname{argmax}_{\theta} \ell(\theta \mid D)$$
true because  $\log$  is monotonically increasing

# MLE as optimization (Analytical Way)

# Example: coin tossing (Bernoulli Trials)

**Experiment**: tossing a coin X, not necessarily fair (X = 1 head, X = 0 tail)

**Parameters**: 
$$\theta := \{ \pi \} \Leftrightarrow P(X=1) = \pi, P(X=0) = 1 - \pi \}$$

**Observations**: a sequence of experimental outcomes

$$D = \{d^{(1)} = \{X^{(1)} = x^{(1)}\}, \ d^{(2)} = \{X^{(2)} = x^{(2)}\}, \dots, \ d^{(N)} = \{X^{(N)} = x^{(N)}\}\}$$

Binomial distribution

$$P(D|\theta) = \binom{N}{N_{X=1}} \prod_{i} P(X^{(i)}|\theta) = \binom{N}{N_{X=1}} P(X = 0|\theta)^{N_{X=0}} P(X = 0|\theta)^{N_{X=0}}$$

$$P(D|\theta) = \binom{N}{N_{X=1}} \prod_{i} P(X^{(i)}|\theta) = \binom{N}{N_{X=1}} P(X = 1|\theta)^{N_{X=1}} P(X = 0|\theta)^{N_{X=0}}$$

$$N_{X=1} \text{ is the number of } X=1 \text{ (i.e. heads) in a sequence of } N \text{ trials}$$

$$= \binom{N}{N_{X=1}} \pi^{N_{X=1}} (1-\pi)^{N_{X=0}}$$

It is the probability of obtaining  $N_{X=1}$  times 'head' in a sequence of N trials In this case, it is assumed to be the likelihood of  $\{d^{(1)}, \ldots, d^{(N)}\}$  given the parameters  $\theta$ 

Artificial Intelligence 2024–2025 Supervised Learning [10]

### Example: coin tossing (Bernoulli Trials)

(Log-)Likelihood Function

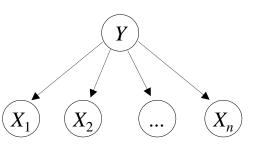
$$\begin{split} &\ell(\theta|D) = \log P(D|\theta) = \log P(\{X^{(i)}\}|\theta) = \log \binom{N}{N_{X=1}} \prod_{i} P(X^{(i)}|\theta) = \log \binom{N}{N_{X=1}} + \sum_{i} \log P(X^{(i)}|\theta) \\ &\text{Rewrite } P(X|\theta) \text{ as:} \\ &P(X\mid\theta) = \pi^{[X=1]} (1-\pi)^{[X=0]} \quad \text{where:} \quad [X^{(i)} = v] := \left\{ \begin{array}{cc} 1 & \text{if} \quad X^{(i)} = v \\ 0 & \text{if} \quad X^{(i)} \neq v \end{array} \right. & \text{Also called} \\ &0 & \text{if} \quad X^{(i)} \neq v \end{array} \\ &\ell(\theta\mid D) = \log \binom{N}{N_{X=1}} + \sum_{i} \log \left( \pi^{[X^{(i)}=1]} \; (1-\pi)^{[X^{(i)}=0]} \; \right) = \\ &= \log \binom{N}{N_{X=1}} + \log \pi \sum_{i} [X^{(i)} = 1] \; + \log (1-\pi) \sum_{i} [X^{(i)} = 0] \\ &= \log \binom{N}{N_{X=1}} + N_{X=1} \log \pi \; + N_{X=0} \log (1-\pi) \end{split}$$

Maximum Likelihood Estimation

$$\frac{\partial \ell}{\partial \theta} = \frac{\partial \ell}{\partial \pi} = \frac{N_{X=1}}{\pi} - \frac{N_{X=0}}{(1-\pi)} \qquad \qquad \frac{\partial \ell}{\partial \theta} = 0 \quad \Rightarrow \quad \theta_{ML}^* = \frac{N_{X=1}}{N_{X=1} + N_{X=0}} = \frac{N_{X=1}}{N}$$

Artificial Intelligence 2024–2025 Supervised Learning [11]

$$P(Y, X_1, ..., X_n) = P(Y) \prod_{i=1}^{n} P(X_i \mid Y)$$



Parameters: the conditional probability tables in the graphical model

$$\theta := \{ \pi_k, \ \pi_{ijk} \}$$
 ,  $P(Y = k) =: \pi_k \ P(X_i = j \mid Y = k) =: \pi_{ijk}$ 

**Observations**: a set of messages with classification

$$D = \{d^{(1)} = \{Y^{(1)} = 1, X_1^{(1)} = 1, \dots, X_n^{(1)} = 0\},\$$

$$\dots,$$

$$d^{(N)} = \{Y_2^{(N)} = y^{(N)}, X_1^{(N)} = x_1^{(N)}, \dots, X_n^{(N)} = x_n^{(N)}\}\}$$

#### Likelihood Function

$$L(\theta|D) \ = \ P(D|\theta) \ = \ P(\{d^{(m)}\}|\{\pi_k,\pi_{ijk}\}) \ = \ \prod_m P(d^{(m)}|\{\pi_k,\pi_{ijk}\})$$
 (data items are IID) 
$$= \ \prod_m P(\{Y^{(m)} = y^{(m)},X_i{}^{(m)} = x_i{}^{(m)}\}|\{\pi_k,\pi_{ijk}\})$$
 (factorization) 
$$= \ \prod_m P(Y^{(m)} = y^{(m)}|\{\pi_k,\pi_{ijk}\}) \ P(\{X_i{}^{(m)} = x_i{}^{(m)}\}|Y^{(m)} = y^{(m)},\{\pi_k,\pi_{ijk}\})$$
 (cond. independence) 
$$= \ \prod_m P(Y^{(m)} = y^{(m)}|\{\pi_k\}) \ P(\{X_i{}^{(m)} = x_i{}^{(m)}\}|Y^{(m)} = y^{(m)},\{\pi_{ijk}\})$$
 
$$= \ \prod_m P(Y^{(m)} = y^{(m)}|\{\pi_k\}) \ \prod_i P(X_i{}^{(m)} = x_i{}^{(m)}|Y^{(m)} = y^{(m)},\{\pi_{ijk}\})$$

$$P(Y, X_1, \dots, X_n) = P(Y) \prod_{i=1}^n P(X_i \mid Y)$$

$$X_1 \qquad X_2 \qquad \dots \qquad X_n$$

Log-Likelihood Function

$$\ell(\{\pi_k, \pi_{ijk}\}|D) = \sum_{m} \log P(Y^{(m)} = y^{(m)}|\{\pi_k\}) + \sum_{m} \sum_{i} \log P(X_i^{(m)} = x_i^{(m)}|Y^{(m)} = y^{(m)}, \{\pi_{ijk}\})$$

*Alternative form for P:* (rewritten using indicator functions)

$$P(Y = k | \{\pi_k\}) = \prod_{k} \pi_k^{[Y=k]}$$

$$P(X_i = j | Y = k, \{\pi_{ijk}\}) = \prod_{j} \prod_{k} \pi_{i,j,k}^{[X_i = j][Y=k]}$$

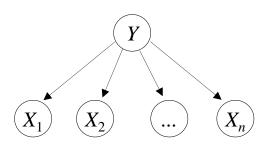
$$\ell(\{\pi_k, \pi_{ijk}\}|D) = \sum_{m} \sum_{k} [Y^{(m)} = k] \log \pi_k + \sum_{m} \sum_{i} \sum_{j} \sum_{k} [X_i^{(m)} = j][Y^{(m)} = k] \log \pi_{ijk}$$

Being both positive and depending on different variables, the two terms above can be optimized separately

Artificial Intelligence 2024–2025 Supervised Learning [13]

$$P(Y, X_1, \dots, X_n) = P(Y) \prod_{i=1}^n P(X_i \mid Y)$$

Lagrange multiplier



#### Maximum Likelihood Estimation

$$\ell(\{\pi_k, \pi_{ijk}\}|D) = \sum_{m} \sum_{k} [Y^{(m)} = k] \log \pi_k + \sum_{m} \sum_{i} \sum_{j} \sum_{k} [X_i^{(m)} = j][Y^{(m)} = k] \log \pi_{ijk}$$

Optimizing first term:

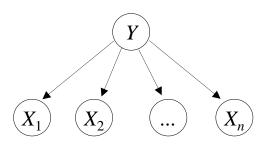
$$\ell^*(\{\pi_k\}|D) = \sum \sum [Y^{(m)} = k] \log \pi_k + \lambda (1 - \sum \pi_k)$$

$$\frac{\partial \ell^*}{\partial \pi_k} = \frac{\sum\limits_{m} [Y^{(m)} = k]}{\pi_k} - \lambda$$
 number of messages in  $D$  classified as  $k$  
$$\frac{\partial \ell^*}{\partial \pi_k} = 0 \quad \Rightarrow \quad \pi_k = \frac{N_{Y=k}}{\lambda}$$

$$\sum_{k} \pi_{k} = 1 \quad \Rightarrow \quad \sum_{k} \frac{N_{Y=k}}{\lambda} = 1 \quad \Rightarrow \quad \lambda = \sum_{k} N_{Y=k} = N$$

$$\pi_k^* = \frac{N_{Y=k}}{N}$$
 (Maximum Likelihood Estimator of  $\pi_k$ )

$$P(Y, X_1, \dots, X_n) = P(Y) \prod_{i=1}^n P(X_i \mid Y)$$



Maximum Likelihood Estimation

$$\ell(\{\pi_k, \pi_{ijk}\}|D) = \sum_{m} \sum_{k} [Y^{(m)} = k] \log \pi_k + \sum_{m} \sum_{i} \sum_{j} \sum_{k} [X_i^{(m)} = j][Y^{(m)} = k] \log \pi_{ijk}$$

Optimizing second term:

$$\ell^*(\{\pi_{ijk}\}|D) = \sum_{m} \sum_{i} \sum_{j} \sum_{k} [X_i^{(m)} = j][Y^{(m)} = k] \log \pi_{ijk} + \sum_{i} \sum_{k} \lambda_{ik} (1 - \sum_{j} \pi_{ijk})$$

$$\frac{\partial \ell^*}{\partial \pi_{ijk}} = \frac{\sum\limits_{m} [X_i^{(m)} = j][Y^{(m)} = k]}{\pi_{ijk}} - \lambda_{ik}$$

$$\frac{\partial \ell^*}{\partial \pi_{ijk}} = 0 \quad \Rightarrow \quad \pi_{ijk} = \frac{N_{X_i=j, Y=k}}{\lambda_{ik}}$$

$$\sum_{j} \pi_{ijk} = 1 \quad \Rightarrow \quad \sum_{j} \frac{N_{X_i=j, Y=k}}{\lambda_{ik}} = 1 \quad \Rightarrow \quad \lambda = \sum_{j} N_{X_i=j, Y=k} = N_{Y=k}$$

$$\pi^*_{ijk} = rac{N_{X_i=j, \ Y=k}}{N_{Y=k}}$$
 (Maximum Likelihood Estimator of  $\pi_{ijk}$ )

### MLE and discrete probability distributions

#### MLE as relative frequencies

The maximum likelihood estimator (MLE) for any discrete probability distribution over a dataset is calculated by determining the proportion of times each outcome occurs in the dataset

This method uses the *relative frequencies* of the outcomes to provide the best estimate of their probabilities

Artificial Intelligence 2024–2025 Supervised Learning [16]

### MLE and discrete probability distributions

#### MLE as relative frequencies (in subspaces)

The maximum likelihood estimator (MLE) for any discrete probability distribution over a dataset is calculated by determining the proportion of times each outcome occurs in the dataset

This method uses the *relative frequencies* of the outcomes to provide the best estimate of their probabilities

When considering *conditional probabilities*, the MLE can be applied within <u>specific subspaces</u> of the dataset by focusing on the subset of the data where that condition is met

Huh? What?

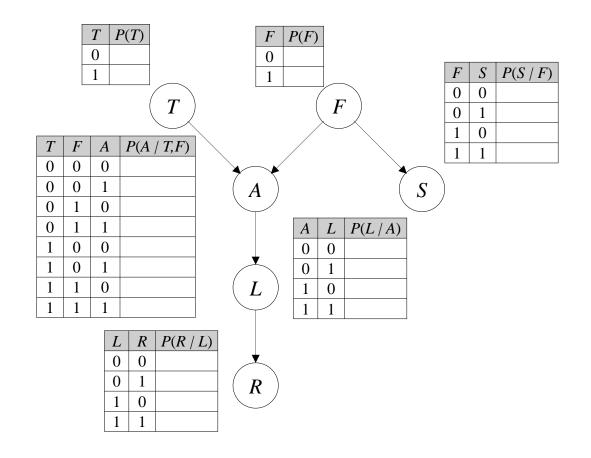
# MLE for Graphical Models: A Practical Rule

# Learning CPTs for a graphical model via MLE

**Model**: random variables plus the graph of dependencies

**Observations**: dataset of values, from <u>completely observed</u> outcomes

Parameters (to be determined): all conditional probabilities (i.e. all CPTs)



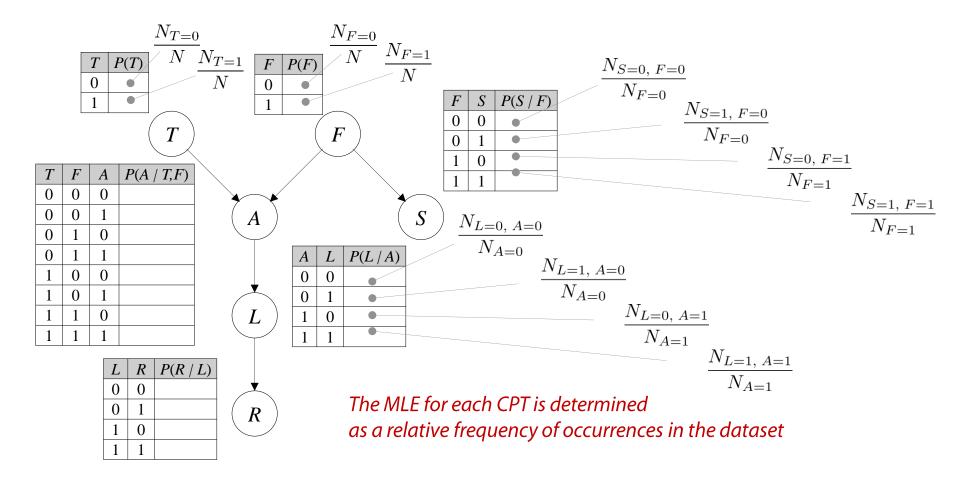
Artificial Intelligence 2024–2025 Supervised Learning [19]

# Learning CPTs for a graphical model via MLE

**Model**: random variables plus the graph of dependencies

**Observations**: dataset of values, from <u>completely observed</u> outcomes

Parameters (to be determined): all conditional probabilities (i.e. all CPTs)

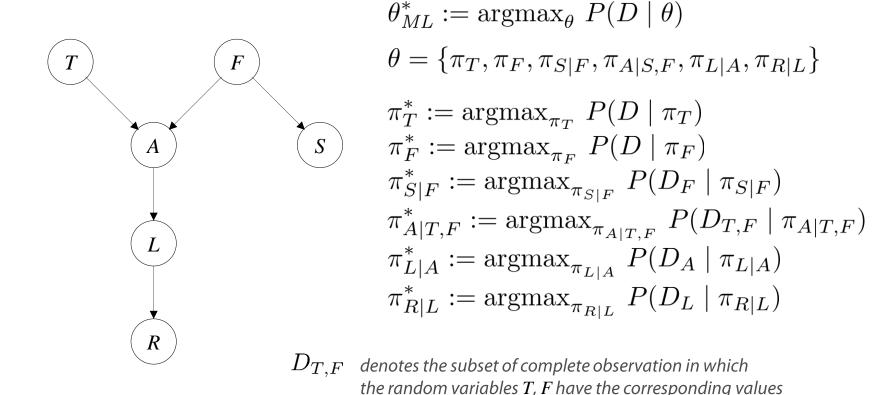


Artificial Intelligence 2024–2025 Supervised Learning [20]

### Learning CPTs for a graphical model via MLE

#### *More in general:*

The MLE of a (directed) graphical model is the MLE of each node (in each corresponding observation subset)



Artificial Intelligence 2024-2025 Supervised Learning [21]

### Bayesian Learning: Maximum a Posteriori (MAP) estimator

## Bayesian learning

#### Maximum a Posteriori Estimation (MAP)

Instead of a likelihood function, the a posteriori probability is maximized

$$P(\theta|D) = \frac{P(D|\theta) P(\theta)}{P(D)} = \frac{P(D|\theta) P(\theta)}{\sum_{\theta} P(D|\theta) P(\theta)}$$

Which is equivalent to optimize, w.r.t.  $\theta$ :

$$P(D|\theta) P(\theta)$$

$$\theta_{MAP}^* := \operatorname{argmax}_{\theta} P(D|\theta) P(\theta)$$

#### Advantages:

- Regularization: not all possible combinations of values might be present in D
- A formula for incremental learning:
   a priori terms could represent what was known before observations D

#### **Problem:**

• Which *prior* distribution?  $P(\theta)$ 

Artificial Intelligence 2024–2025 Supervised Learning [23]

### Beta distribution

Gamma function (n integer > 0)

$$\Gamma(n) := (n-1)!$$

Beta function ( $\alpha$  and  $\beta$  integers > 0)

$$B(\alpha,\beta) := \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)} = \frac{(\alpha-1)!(\beta-1)!}{(\alpha+\beta-1)!} \quad \text{The definition is more complex when $\alpha$ and $\beta$ are not integers (see Wikipedia)}$$

■ Beta probability density function (pdf) ( $\alpha$  and  $\beta$  integers > 0)

$$\operatorname{Beta}(\theta;\alpha,\beta) := \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{\operatorname{B}(\alpha,\beta)} \quad \text{The maximum occurs at:} \quad \theta = \frac{\alpha-1}{\alpha+\beta-2}$$

$$\operatorname{Beta}(\theta;1,1) \quad \operatorname{Beta}(\theta;2,2) \quad \operatorname{Beta}(\theta;10,10) \quad \operatorname{Beta}(\theta;2,3)$$

Artificial Intelligence 2024-2025

### Conjugate prior distributions

Coin tossing (i.e. Binomial)

 $\alpha_D$  and  $\beta_D$  are the result counts (i.e. heads and tails)

$$P(D|\theta) = {\binom{\alpha_D + \beta_D}{\alpha_D}} \prod_i P(X_i|\theta) = {\binom{\alpha_D + \beta_D}{\alpha_D}} \theta^{\alpha_D} (1 - \theta)^{\beta_D}$$

A posteriori probability with Beta prior

 $\alpha_P$  and  $\beta_P$  are are the **hyperparameters** of the prior

$$P(D|\theta)P(\theta) = \begin{pmatrix} \alpha_D + \beta_D \\ \alpha_D \end{pmatrix} \theta^{\alpha_D} (1-\theta)^{\beta_D} \cdot \text{Beta}(\theta; \alpha_P, \beta_P) = \begin{pmatrix} \alpha_D + \beta_D \\ \alpha_D \end{pmatrix} \theta^{\alpha_D} (1-\theta)^{\beta_D} \cdot \frac{\theta^{\alpha_P - 1} (1-\theta)^{\beta_P - 1}}{B(\alpha_P, \beta_P)}$$
$$= \begin{pmatrix} \alpha_D + \beta_D \\ \alpha_D \end{pmatrix} \frac{\theta^{\alpha_D + \alpha_P - 1} (1-\theta)^{\beta_D + \beta_P - 1}}{B(\alpha_P, \beta_P)} = \begin{pmatrix} \alpha_D + \beta_D \\ \alpha_D \end{pmatrix} \frac{B(\alpha_D + \alpha_P, \beta_D + \beta_P)}{B(\alpha_P, \beta_P)} \text{Beta}(\theta; \alpha_D + \alpha_P, \beta_D + \beta_P)$$

this factor is a positive constant (for  $\theta$ )

Artificial Intelligence 2024–2025 Supervised Learning [25]

### Conjugate prior distributions

Coin tossing (i.e. Binomial)

 $\alpha_D$  and  $\beta_D$  are the result counts (i.e. heads and tails)

$$P(D|\theta) = {\binom{\alpha_D + \beta_D}{\alpha_D}} \prod_i P(X_i|\theta) = {\binom{\alpha_D + \beta_D}{\alpha_D}} \theta^{\alpha_D} (1 - \theta)^{\beta_D}$$

A posteriori probability with Beta prior

$$P(D|\theta)P(\theta) = \binom{\alpha_D + \beta_D}{\alpha_D} \frac{\mathrm{B}(\alpha_D + \alpha_P, \beta_D + \beta_P)}{\mathrm{B}(\alpha_P, \beta_P)} \mathrm{Beta}(\theta; \alpha_D + \alpha_P, \beta_D + \beta_P)$$

$$/ \text{"is proportional to"}$$

$$P(D|\theta)P(\theta) \propto \mathrm{Beta}(\theta; \alpha_D + \alpha_P, \beta_D + \beta_P)$$

#### *Optimization:*

$$\theta_{MAP}^* = \operatorname{argmax}_{\theta} \operatorname{Beta}(\theta; \alpha_D + \alpha_P, \beta_D + \beta_P) = \frac{\alpha_D + \alpha_P - 1}{\alpha_D + \alpha_P + \beta_D + \beta_P - 2}$$

which is the same as MLE but with the addition of  $\alpha_P + \beta_P$  pseudo-observations

Being a **conjugate prior**  $P(\theta)$  of a distribution  $P(D|\theta)$  is in the same family of  $P(\theta)$ 

Artificial Intelligence 2024–2025 Supervised Learning [26]

### Conjugate prior distributions

Coin tossing (i.e. a specific observation i)

$$P(D_i|\theta) = \theta^{[X_i=1]} (1 - \theta)^{[X_i=0]}$$

*Likelihood (of a dataset)* 

$$P(D|\theta) = \binom{N}{N_{X=1}} \prod_{i} P(D_i|\theta) = \binom{N}{N_{X=1}} \theta^{N_{X=1}} (1-\theta)^{N_{X=0}}$$

A posteriori probability with Beta prior

"is proportional to" 
$$P(D|\theta)P(\theta) \propto \mathrm{Beta}(\theta,\ N_{X=1}+\alpha_P,\ N_{X=0}+\beta_P)$$

Therefore

$$\theta_{MAP}^* = \operatorname{argmax}_{\theta} \operatorname{Beta}(\theta, N_{X=1} + \alpha_P, N_{X=0} + \beta_P) = \frac{N_{X=1} + \alpha_P - 1}{N + \alpha_P + \beta_P - 2}$$

which is the same as MLE but with the addition of  $\alpha_P + \beta_P$  pseudo-observations

Being a **conjugate prior**  $P(\theta)$  of a distribution  $P(D|\theta)$  in the above sense means that the posterior  $P(D|\theta)P(\theta)$  is in the same family of  $P(\theta)$ 

### Anti-spam filter

$$P(Y, X_1, \dots, X_n) = P(Y) \prod_{i=1}^n P(X_i \mid X_{i-1})$$

$$X_1 \qquad X_2 \qquad \dots \qquad X_n$$

Maximum a Posteriori (MAP) Estimation

The adapted computations for:

$$\theta_{MAP}^* := \operatorname{argmax}_{\theta} P(D|\theta) P(\theta)$$

yield:

$$\pi_k^* = \frac{\alpha_k + N_{Y=k} - 1}{\alpha_k + \beta_k + N - 2} \qquad (\textit{MAP Estimator of } \pi_k)$$

$$\pi_{ijk}^* = \frac{\alpha_{ijk} + N_{X_i=j, \ Y=k} - 1}{\alpha_{ijk} + \beta_{ijk} + N_{V-k} - 2} \qquad (\textit{MAP Estimator of } \pi_{ijk})$$

where the

$$\alpha_k, \beta_k, \alpha_{ijk}, \beta_{ijk}$$

are the *hyperparameters* of the prior distribution representing the *pseudo-observations* made *before* the arrival of new, actual observations *D* 

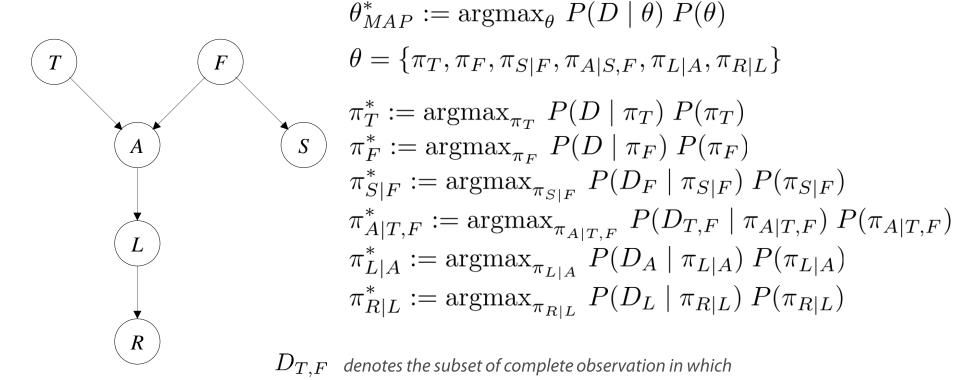
### Bayesian Learning: MAP for Graphical Models

# Learning CPTs for a graphical model

As Maximum a Posteriori Estimation

#### *More in general:*

The MAP of a (directed) graphical model is the MAP of each node (in each corresponding observation subset)



the random variables T,F have the corresponding value

# Machine Learning

Artificial Intelligence 2024–2025 Supervised Learning [31]

### ML = Representation + Evaluation + Optimization

#### Representation

The objective is learning a specific distribution

$$P({X_r};\theta)$$

where  $\{X_r\}$  are all the random variables of interest and  $\theta$  is a *set* of parameters

Which kind of distribution (i.e. the *model* or also the *learner*) do we select?

Example: assume we select the anti-spam filter (i.e. Naïve Bayesian Classifier) as the model the parameters in such case are the numerical probabilities in the CPTs

#### Evaluation

Given a dataset D, how well does a specific set of parameter values  $\hat{\theta}$  make the distribution P fit the dataset?

An estimator, i.e. a scoring function of some sort, must be selected

#### Optimization

How can we find the optimal set of parameter values  $\theta^*$  with respect to the *estimator* of choice?

In general, this is an optimization problem

Artificial Intelligence 2024–2025 Supervised Learning [32]