



Bayesian Machine Learning

03/05/21 - François HU

Outline

1

Bayesian statistics

- define a probabilistic model
- apply bayesian inference
- conjugate priors

2

Latent variable models

- define latent variable and apply them to simplify probabilistic model
- cluster data with latent models like GMM
- train probabilistic models with EM-algorithm

3

Variational Inference

- apply variational inference for probabilistic models
- understand variational interpretation of LDA
- application of LDA to text mining

4

Markov Chain Monte Carlo

- train / do inference almost any probabilistic model with MCMC
- pros and cons of MCMC / VI

5

Extensions and oral presentations

PREREQUISITE

THEORY

1. Notions of **probability & statistics**
2. **Statistical Learning :**
supervised & unsupervised learning
3. **Information theory :**
Entropy, KL-divergence, ...
4. **Monte Carlo & Markov Chain**

APPLICATION

Python (or at least R)

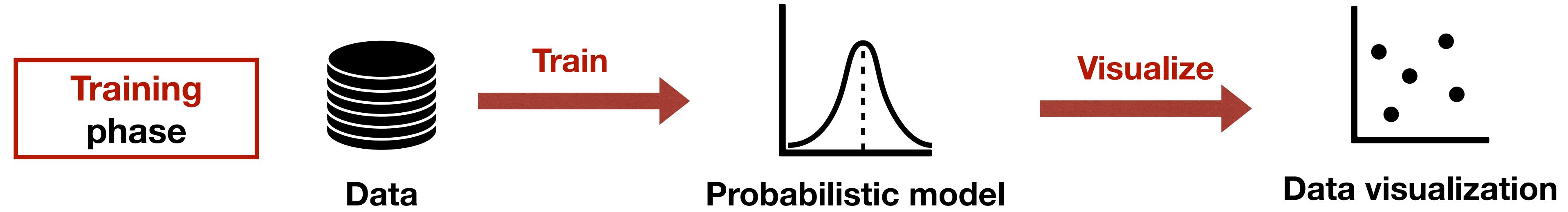
ALGORITHM

Some « classical » supervised & unsupervised models

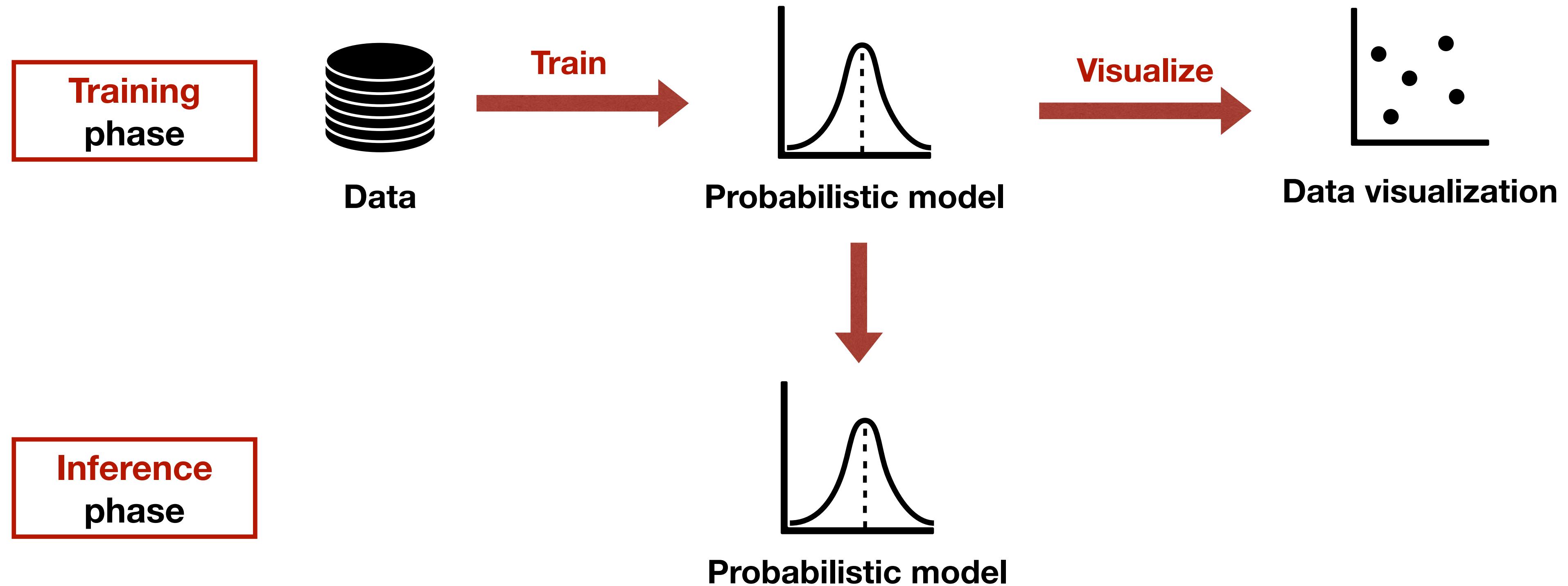
0

Gentle introduction to statistical learning

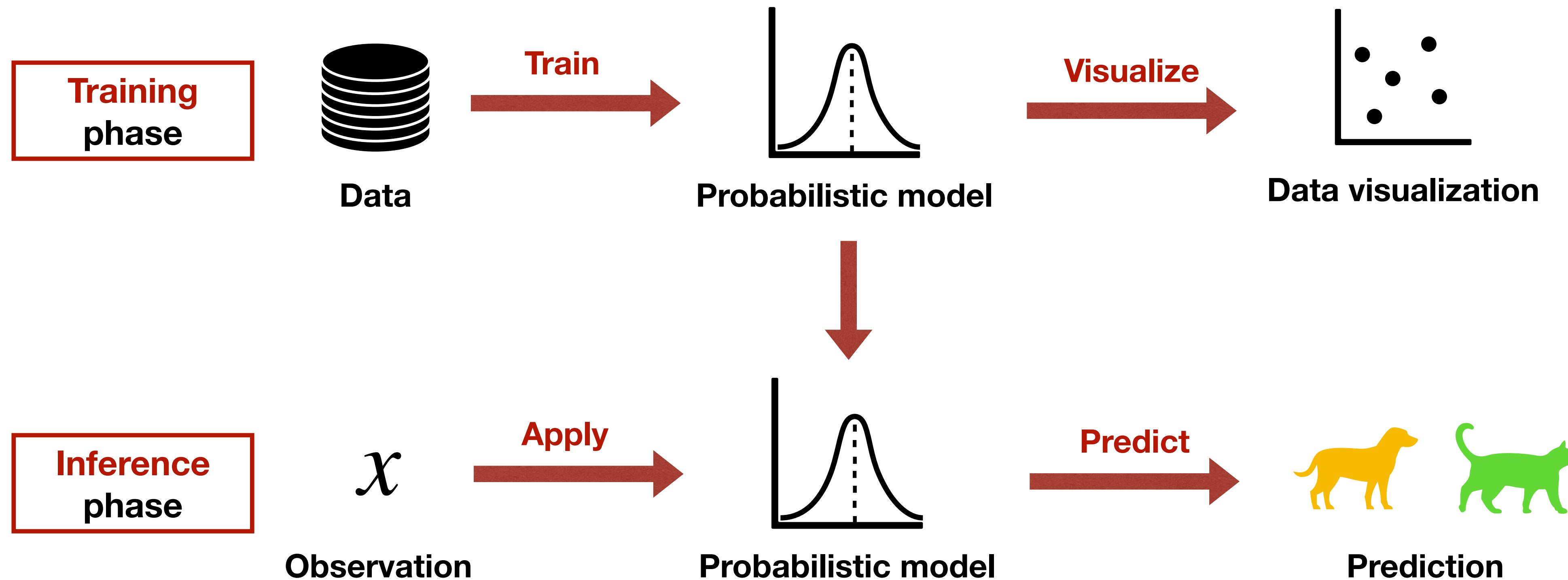
Simplified statistical learning process



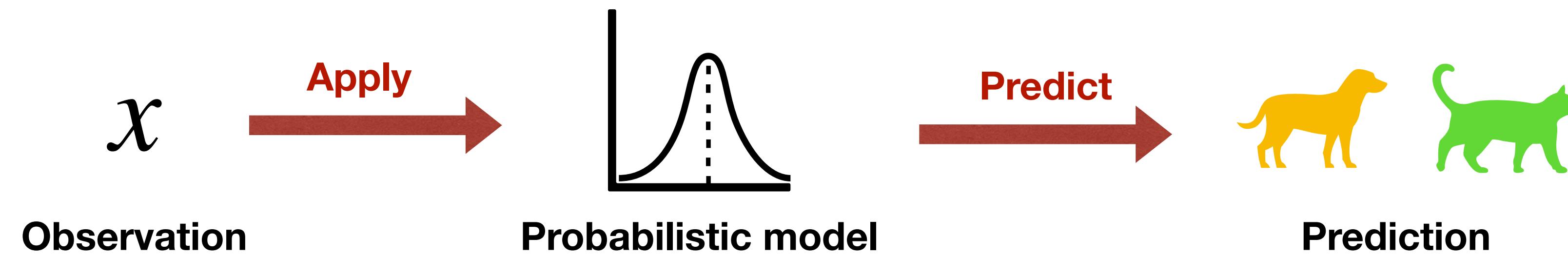
Simplified statistical learning process



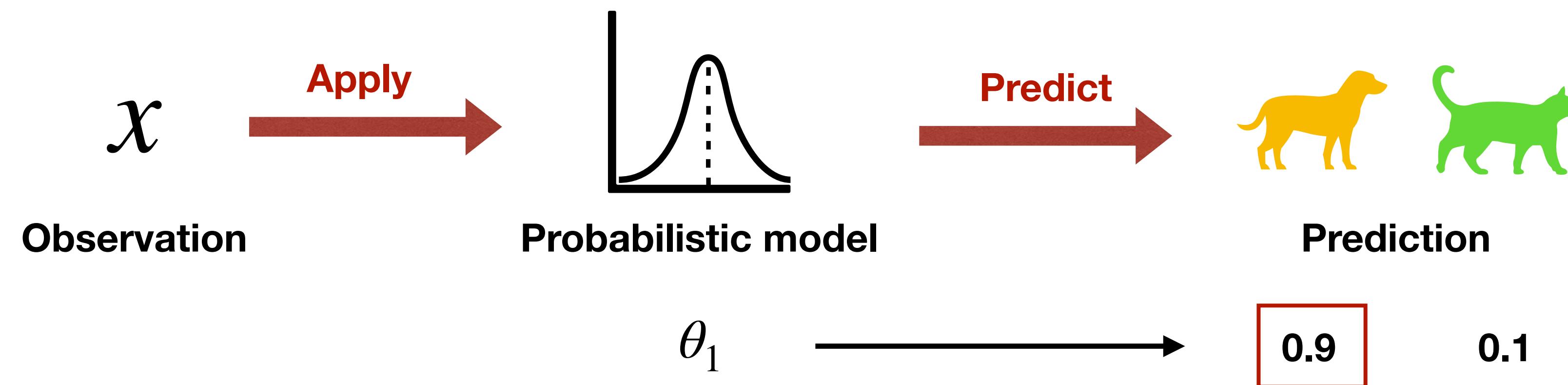
Simplified statistical learning process



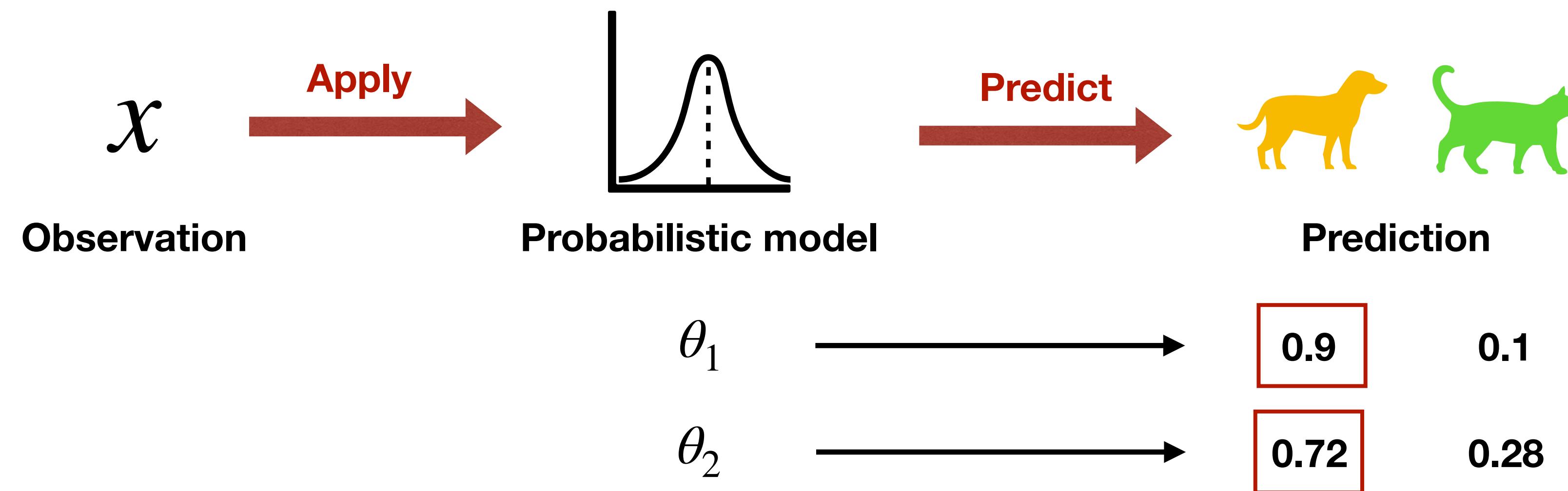
Inference phase



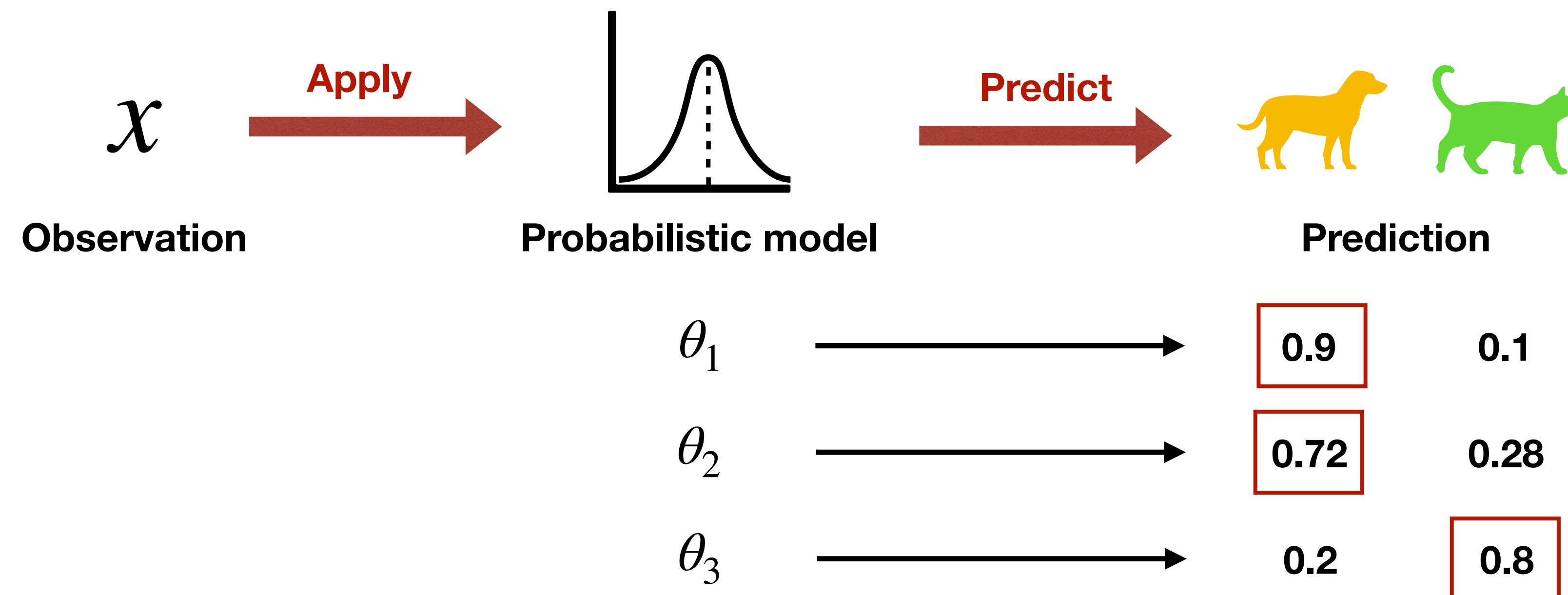
Inference phase



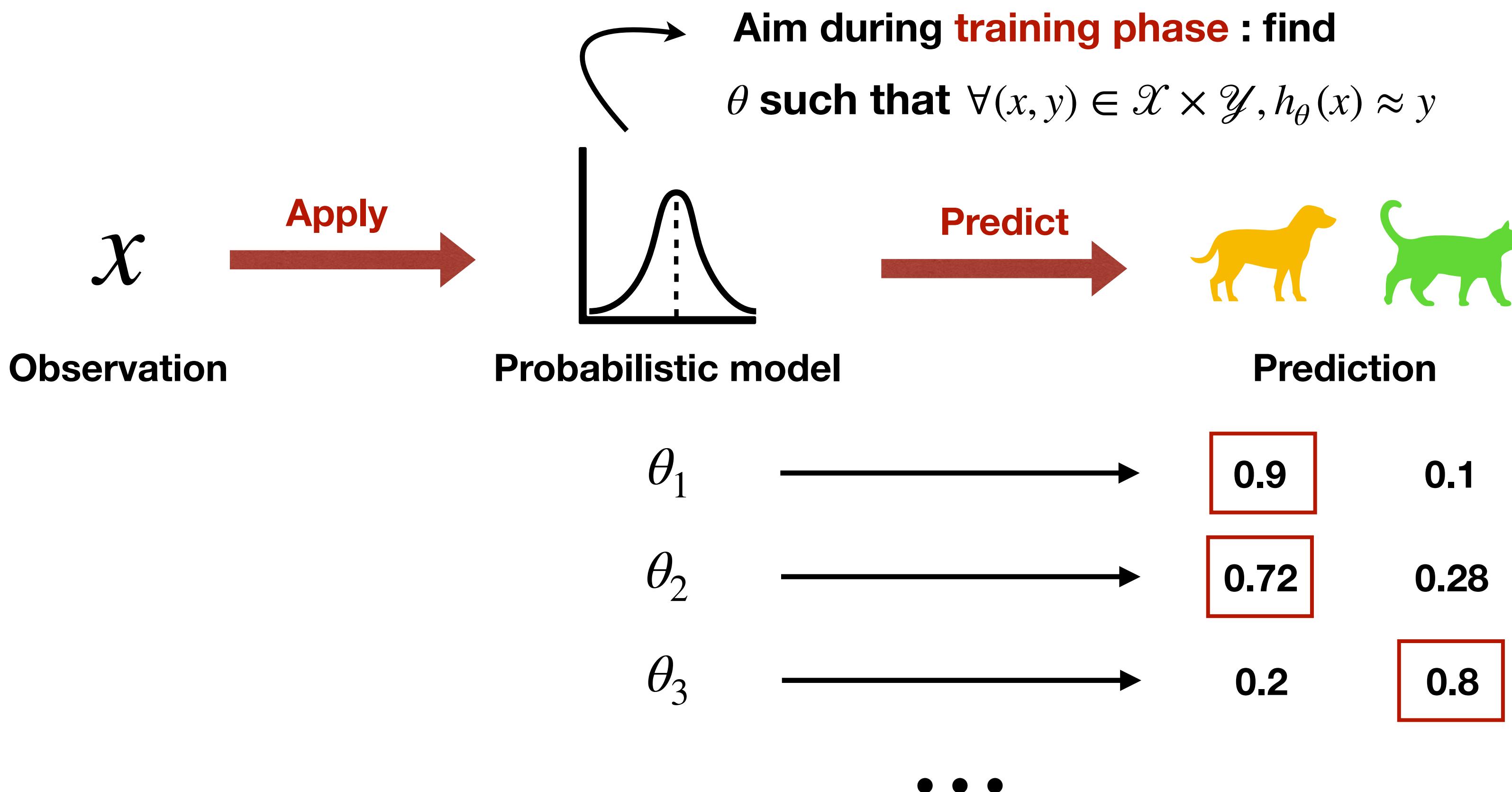
Inference phase



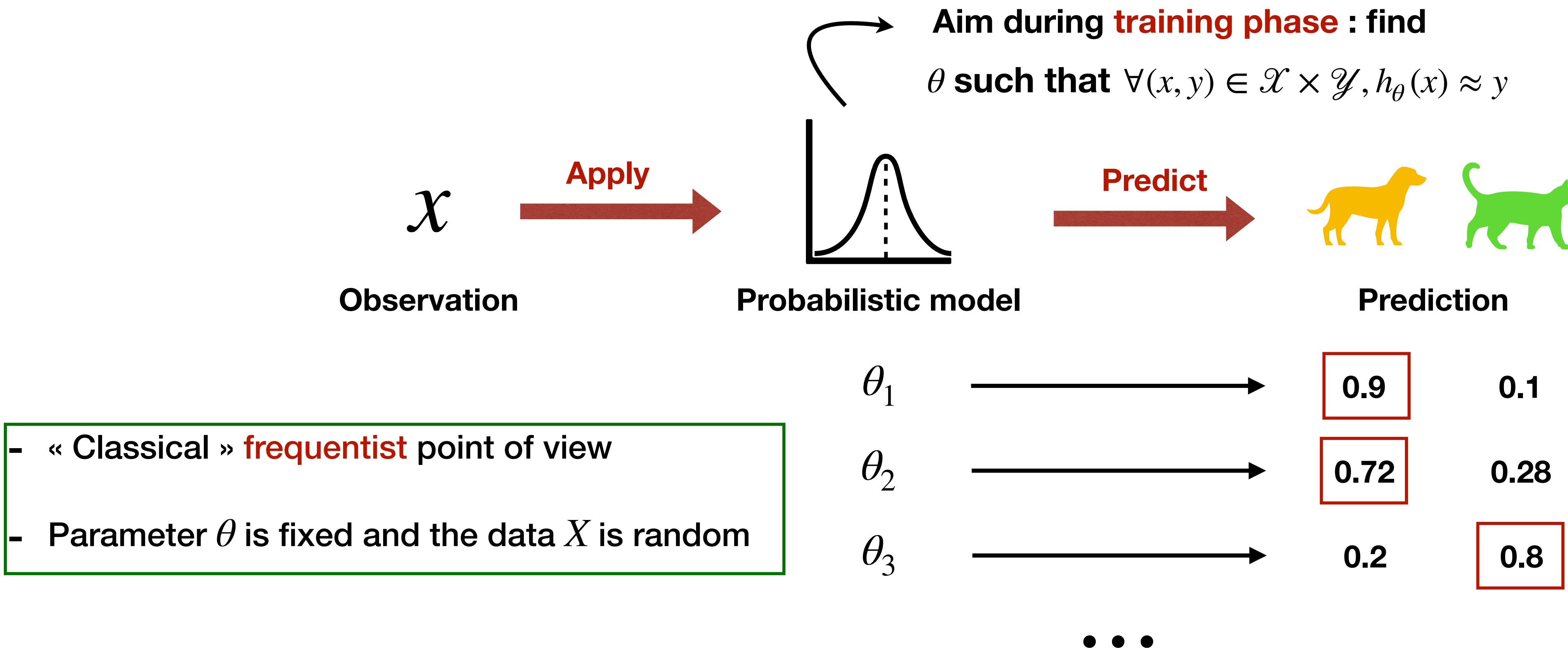
Inference phase



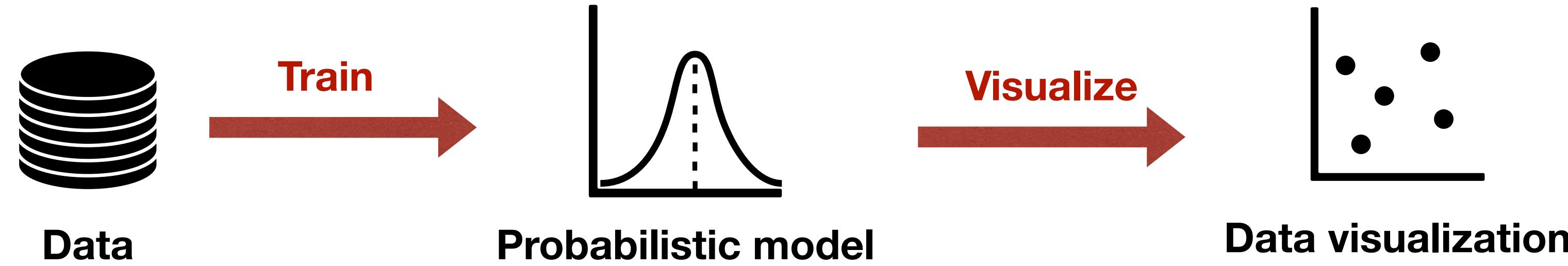
Inference phase



Inference phase



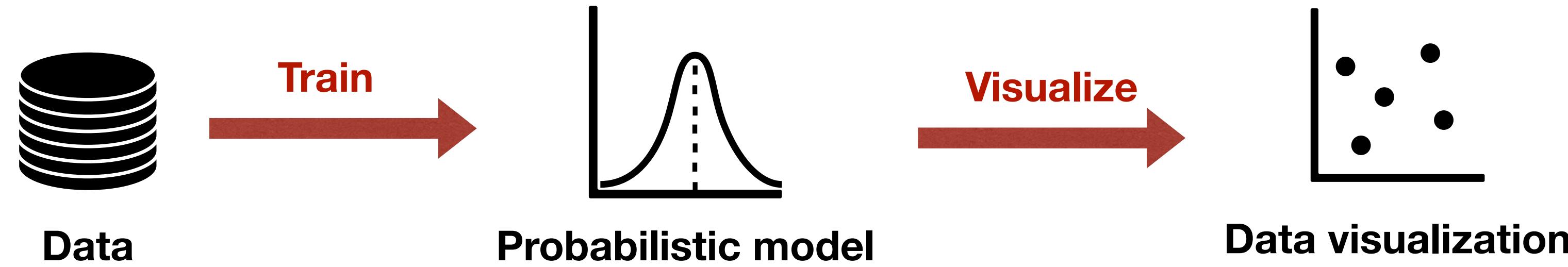
Training phase



Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Training phase



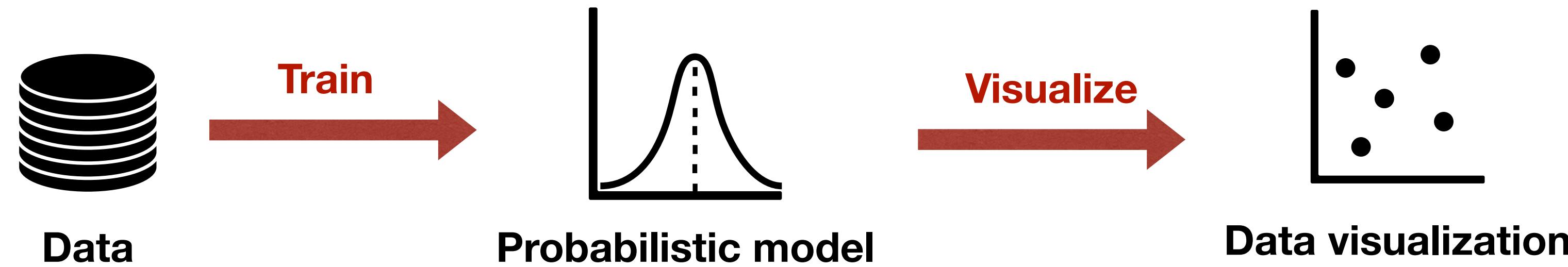
Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Training phase



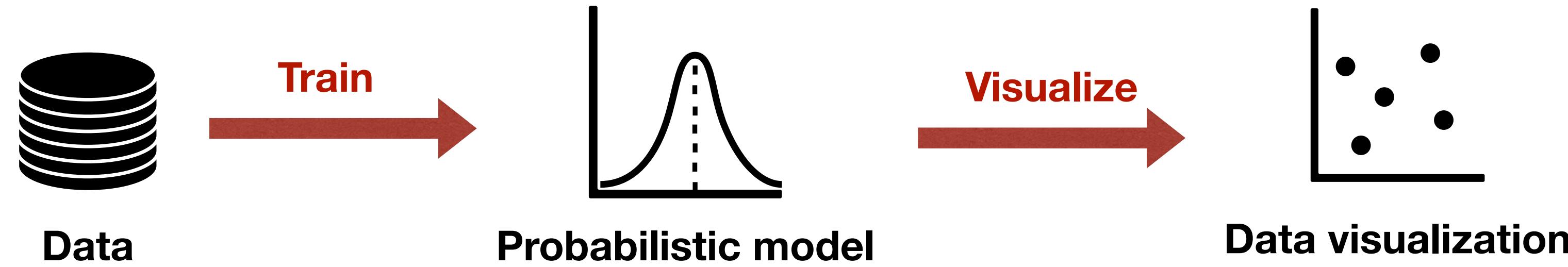
Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Training phase



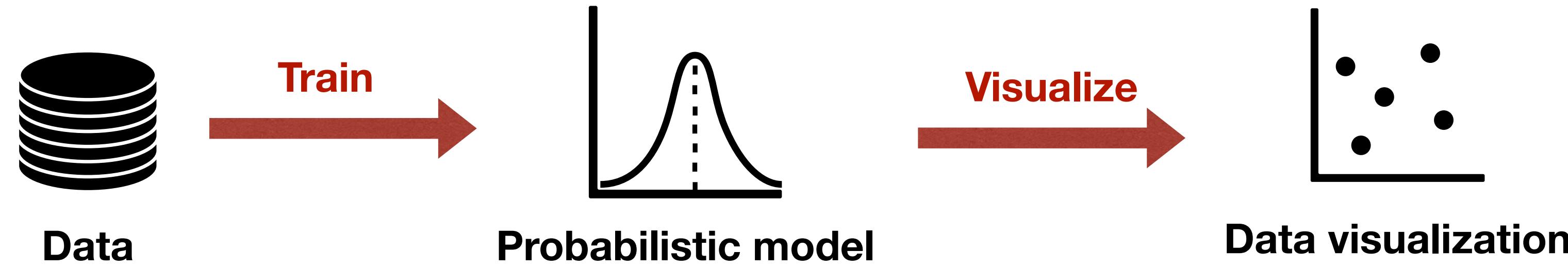
Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Training phase



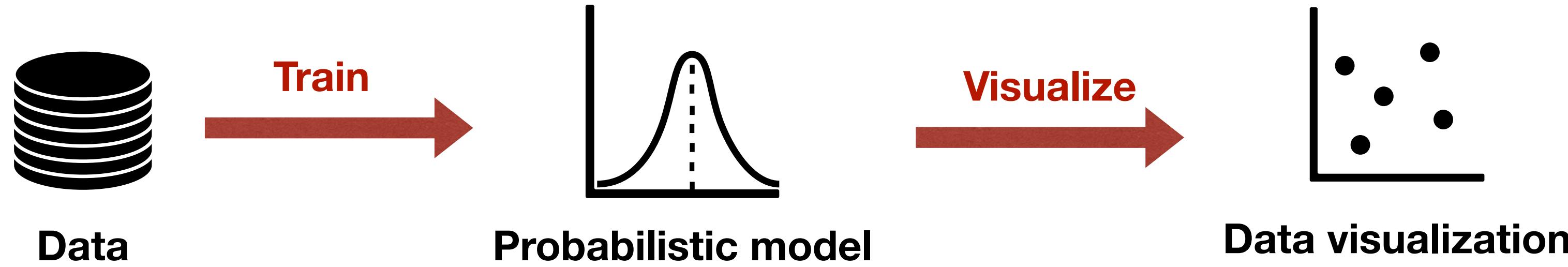
Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Training phase



Aim during **training phase** : find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

Usually in frequentist statistics we use the MLE : **Maximum Likelihood Estimation** $\hat{\theta} = \arg \max_{\theta} P(X | \theta)$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Bayesian statistics

1

Introduction to bayesian statistics

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Example : a deck of 52 playing cards

$$P(\text{King of Hearts}) = \frac{1}{52}$$

$$P(\text{Spade}) = \frac{1}{4}$$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable

1. Introduction to bayesian statistics

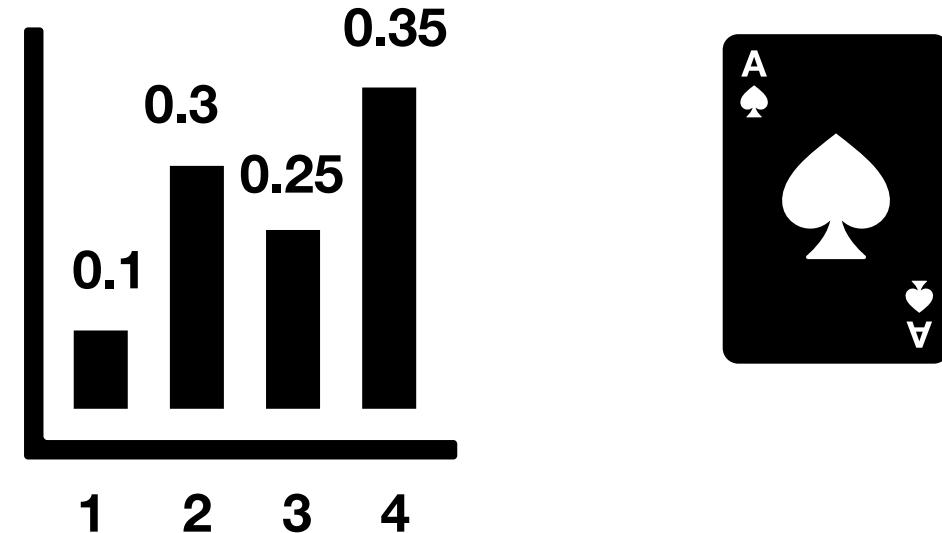
Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

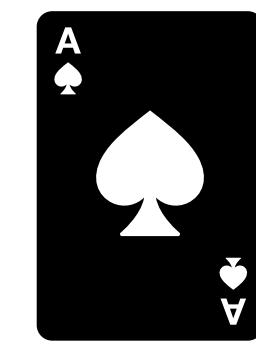
Random variable

- **Discrete** variable Example

Probability Mass Function
(PMF)



$P(X) = \begin{cases} 0.1 & \text{if } X = 1 \\ 0.3 & \text{if } X = 2 \\ 0.25 & \text{if } X = 3 \\ 0.35 & \text{if } X = 4 \end{cases}$



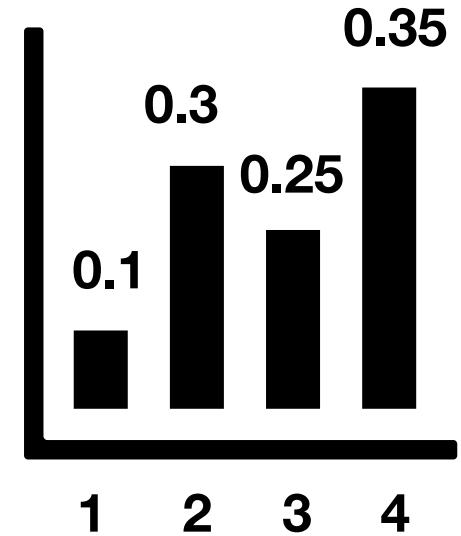
1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

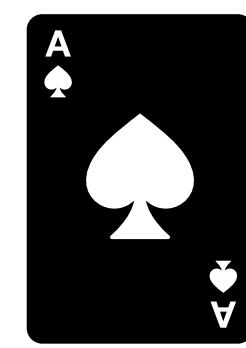
Random variable

- **Discrete** variable Example

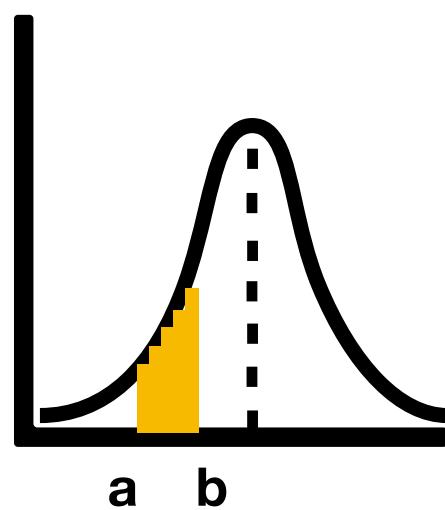


Probability Mass Function
(PMF)

$$P(X) = \begin{cases} 0.1 & \text{if } X = 1 \\ 0.3 & \text{if } X = 2 \\ 0.25 & \text{if } X = 3 \\ 0.35 & \text{if } X = 4 \end{cases}$$

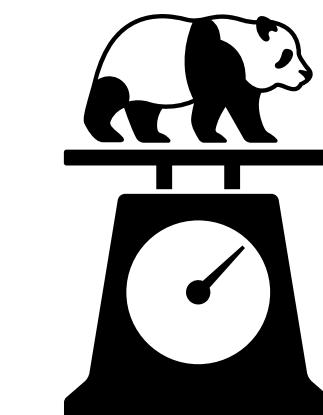


- **Continuous** variable Example



Probability Density function
(PDF)

$$P(X \in [a, b]) = \int_a^b p(s)ds$$



1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

joint probability

marginals

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable

- **Discrete** variable

- **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

Example :

joint probability

marginals

dependency : **one** deck of 52 playing cards

$$P(X_1 = \begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array}, X_2 = \begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array}) = 0$$

$$P(X_1 = \begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array}) \cdot P(X_2 = \begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array}) = \frac{1}{52^2}$$

independency : **two** dices

$$P(X_1 = \begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline\end{array}, X_2 = \begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline\end{array}) = \frac{1}{6^2}$$

$$P(X_1 = \begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline\end{array}) \cdot P(X_2 = \begin{array}{|c|c|}\hline \bullet & \bullet \\ \hline \bullet & \bullet \\ \hline\end{array}) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{6^2}$$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

Conditional probability : probability of X **given that Y happened**

conditional

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

joint probability
marginal

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

Conditional probability : probability of X **given that Y happened**

Example : one deck of 52 playing cards

conditional

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$
 joint probability
marginal

$$P(\begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array} | \begin{array}{|c|}\hline \text{A} \\ \hline \text{C} \\ \hline\end{array}) = \frac{P(\begin{array}{|c|}\hline \text{A} \\ \hline \text{H} \\ \hline\end{array}, \begin{array}{|c|}\hline \text{A} \\ \hline \text{C} \\ \hline\end{array})}{P(\begin{array}{|c|}\hline \text{A} \\ \hline \text{C} \\ \hline\end{array})}$$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules :

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if

$$P(X, Y) = P(X)P(Y)$$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules :

Chain rule $P(X_1, X_2) = P(X_1 | X_2) \times P(X_2)$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules :

Chain rule $P(X_1, X_2, X_3) = P(X_1 | X_2, X_3) \times P(X_2 | X_3) \times P(X_3)$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules :

Chain rule $P(X_1, \dots, X_n) = \prod_{k=1, \dots, n} P(X_k | X_1, \dots, X_{k-1})$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules :

Chain rule $P(X_1, \dots, X_n) = \prod_{k=1, \dots, n} P(X_k | X_1, \dots, X_{k-1})$

Sum rule $P(X) = \sum_{Y \in \mathcal{Y}} P(X, Y) \quad P(X) = \int_{Y \in \mathcal{Y}} P(X, Y) \cdot dY$

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules : **Chain rule** **Sum rule**



Bayes theorem :

1. Introduction to bayesian statistics

Probability & statistics : small review

Probability (non axiomatic definition) of an event : relative **frequency** of an event in an infinite trials

Random variable - **Discrete** variable - **Continuous** variable

Independence : two random variables X and Y are **independent** if $P(X, Y) = P(X)P(Y)$

Conditional probability : probability of X **given that Y happened**

$$P(X | Y) = \frac{P(X, Y)}{P(Y)}$$

Rules : **Chain rule** **Sum rule**



Bayes theorem :

θ Parameters

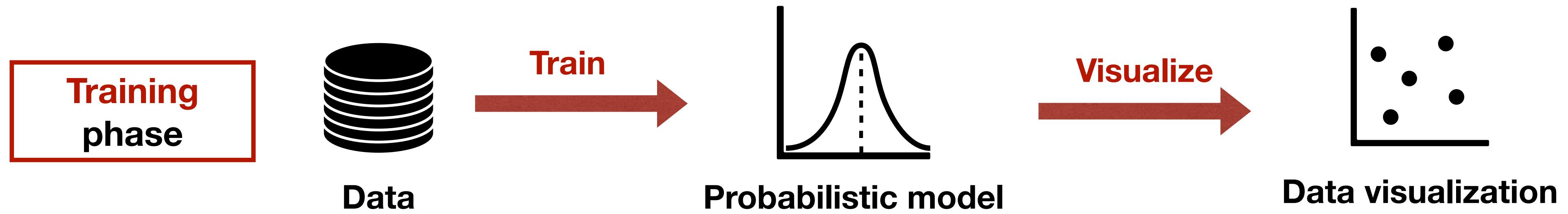
X Observations (data)

	Likelihood	Prior
Posterior	$P(X, \theta)$	$P(\theta) \times P(X)$
	Evidence	

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

1. Introduction to bayesian statistics

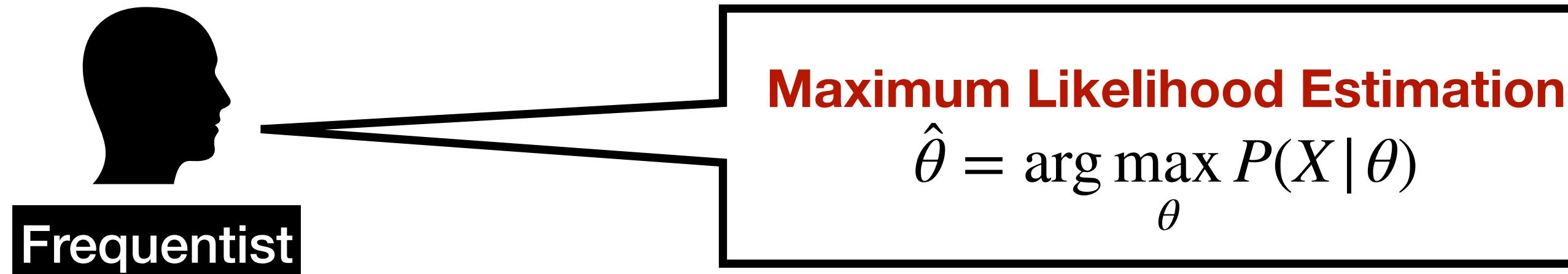
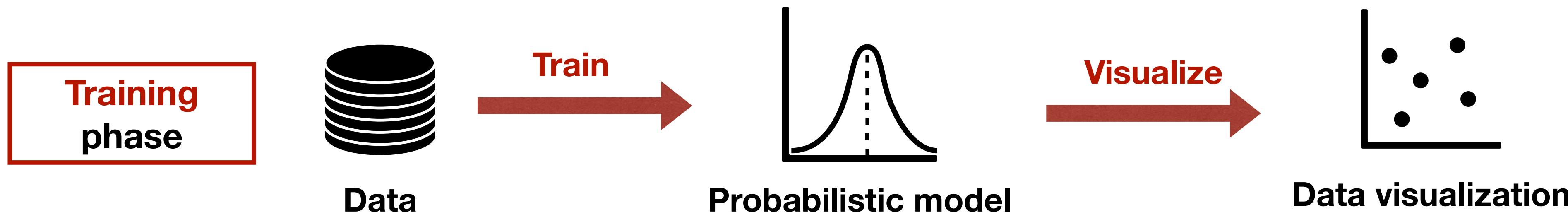
Frequentist VS Bayesian point of view



find θ such that $\forall (x, y) \in \mathcal{X} \times \mathcal{Y}, h_\theta(x) \approx y$

1. Introduction to bayesian statistics

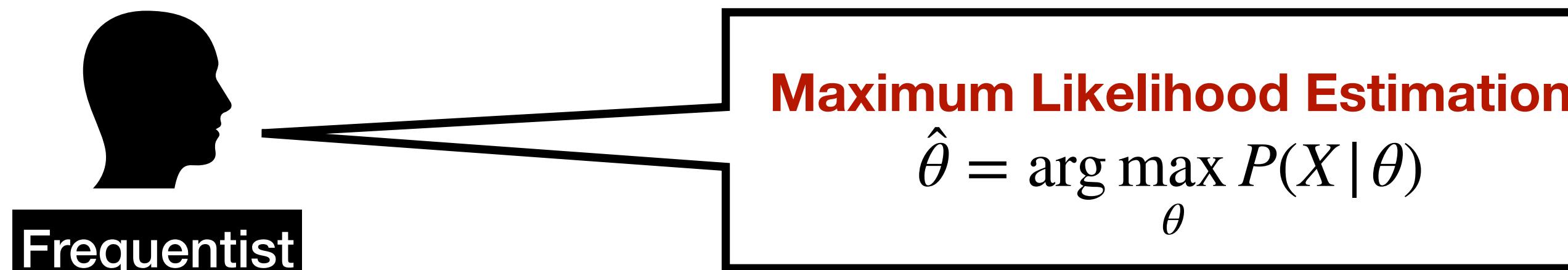
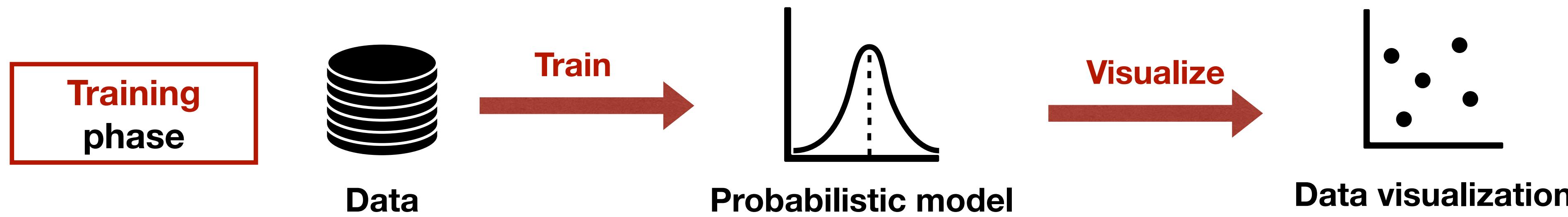
Frequentist VS Bayesian point of view



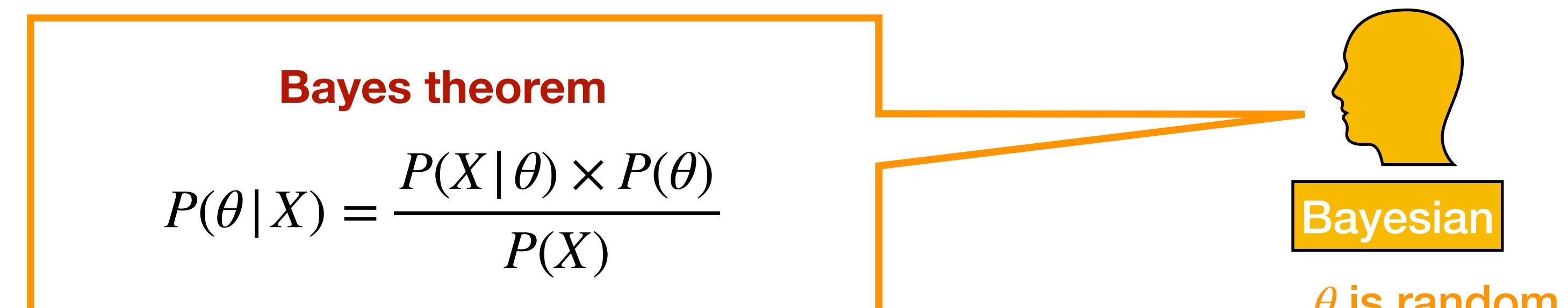
θ is fixed
 X is random

1. Introduction to bayesian statistics

Frequentist VS Bayesian point of view



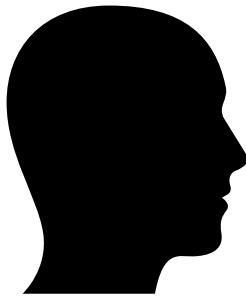
θ is fixed
 X is random



θ is random
 X is fixed

1. Introduction to bayesian statistics

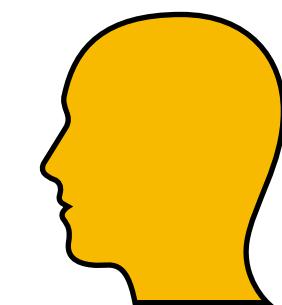
Frequentist VS Bayesian point of view



Frequentist

Maximum Likelihood Estimation

$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$



Bayesian

Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

1. Introduction to bayesian statistics

Frequentist VS Bayesian point of view



Frequentist

Maximum Likelihood Estimation

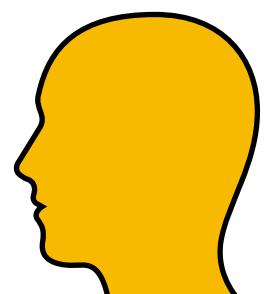
$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$

Problems :

- Only works well if we have big data : $|X| \gg |\theta|$
- Cannot start with a « belief » hence not practical nor flexible
- Cannot express uncertainty of estimated model parameters and predictions

Bayes theorem

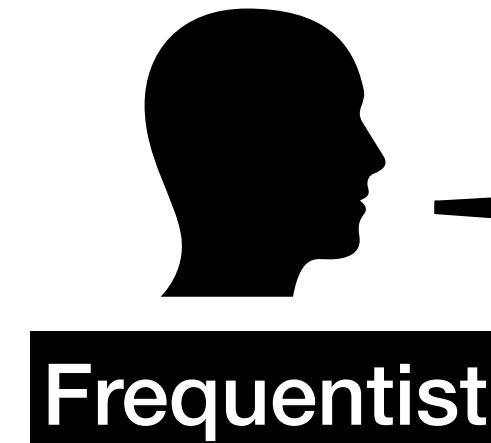
$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$



Bayesian

1. Introduction to bayesian statistics

Frequentist VS Bayesian point of view



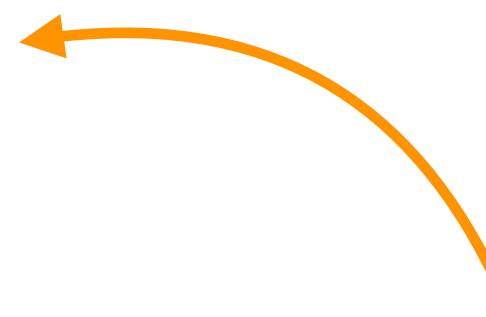
Maximum Likelihood Estimation

$$\hat{\theta} = \arg \max_{\theta} P(X | \theta)$$

Frequentist

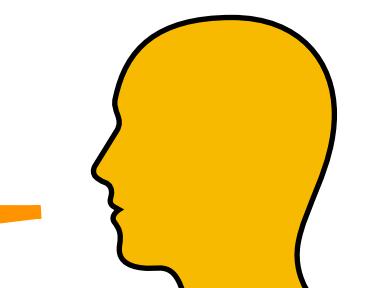
Problems :

- ~~Only works well if we have big data : $|X| \gg |\theta|$~~
- ~~Cannot start with a « belief » hence not practical nor flexible~~
- ~~Cannot express uncertainty of estimated model parameters and predictions~~



Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$



Bayesian

1. Introduction to bayesian statistics

Bayesian point of view : classification

Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Training phase

$$P(\theta | X_{train}, y_{train}) = \frac{P(y_{train} | X_{train}, \theta) \times P(\theta)}{P(y_{train} | X_{train})}$$

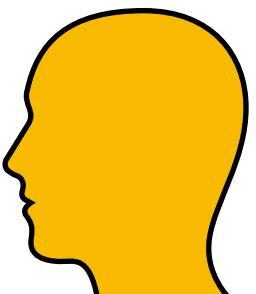


1. Introduction to bayesian statistics

Bayesian point of view : training

Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$



Bayesian

Training
phase

$$P(\theta | X_{train}, y_{train}) = \frac{P(y_{train} | X_{train}, \theta) \times P(\theta)}{P(y_{train} | X_{train})}$$

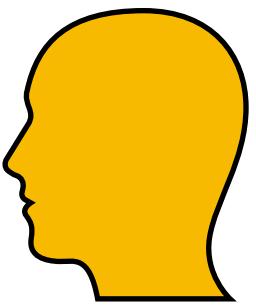
Can regularize your model when training on your data

1. Introduction to bayesian statistics

Bayesian point of view : inference

Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$



Bayesian

Training phase

$$P(\theta | X_{train}, y_{train}) = \frac{P(y_{train} | X_{train}, \theta) \times P(\theta)}{P(y_{train} | X_{train})}$$

Can regularize your model when training on your data

Inference phase

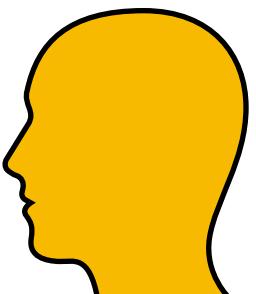
$$P(y_{new} | X_{new}, X_{train}, y_{train}) = \int P(y_{new} | X_{train}, \theta) \times P(\theta | X_{train}, y_{train}) d\theta$$

1. Introduction to bayesian statistics

Bayesian point of view : online learning

Bayes theorem

$$P(\theta | X) = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$



Bayesian

Training phase

$$P(\theta | X_{train}, y_{train}) = \frac{P(y_{train} | X_{train}, \theta) \times P(\theta)}{P(y_{train} | X_{train})}$$

Can regularize your model when training on your data

Inference phase

$$P(y_{new} | X_{new}, X_{train}, y_{train}) = \int P(y_{new} | X_{train}, \theta) \times P(\theta | X_{train}, y_{train}) d\theta$$

Online learning

$$P_{new}(\theta) = P(\theta | x_{new}) = \frac{P(x_{new} | \theta) \times P_{old}(\theta)}{P(x_{new})}$$

New prior

Posterior



2

Probabilistic models

2. Probabilistic model

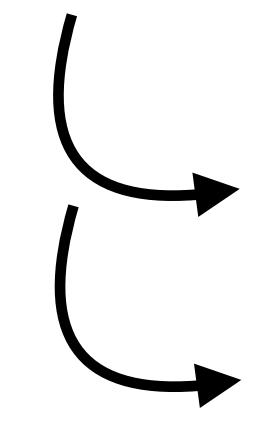
Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

- 
- **Nodes** : random variables
 - **Links** : probabilistic relationships

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks
(Directed graphical models)

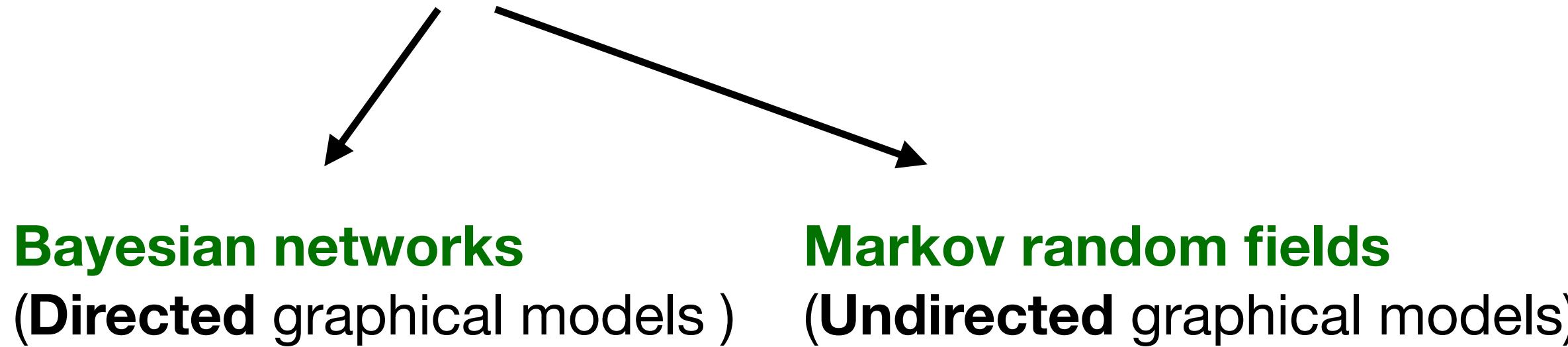


- **Nodes** : random variables
- **Links** : probabilistic relationships

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions



- **Nodes** : random variables
- **Links** : probabilistic relationships

2. Probabilistic model

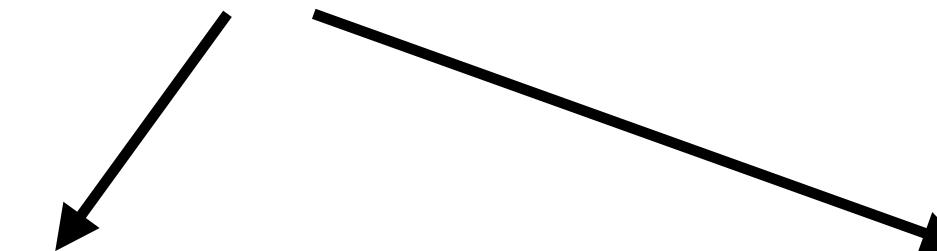
Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks
(Directed graphical models)

Markov random fields
(Undirected graphical models)

The focus of our course !



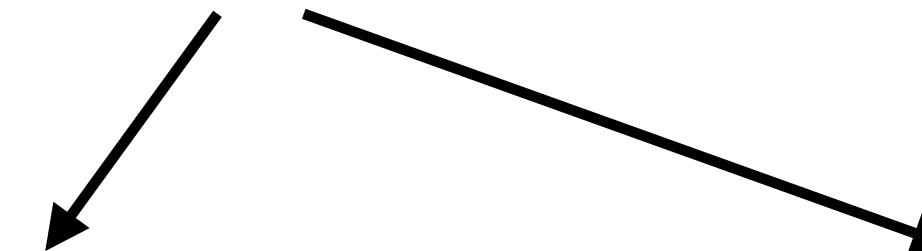
- **Nodes** : random variables
- **Links** : probabilistic relationships

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks
(Directed graphical models)



Markov random fields
(Undirected graphical models)

The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \dots$$

- **Nodes** : random variables
- **Links** : probabilistic relationships

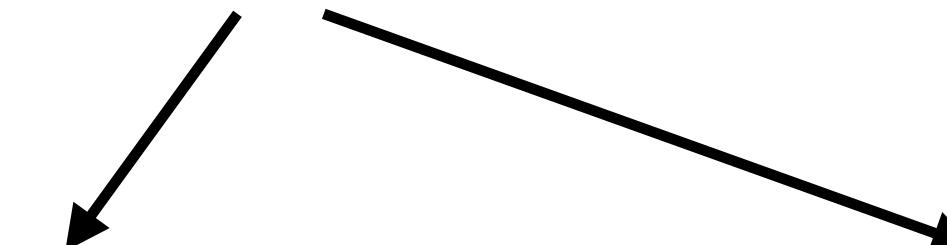
2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks

(Directed graphical models)



Markov random fields

(Undirected graphical models)

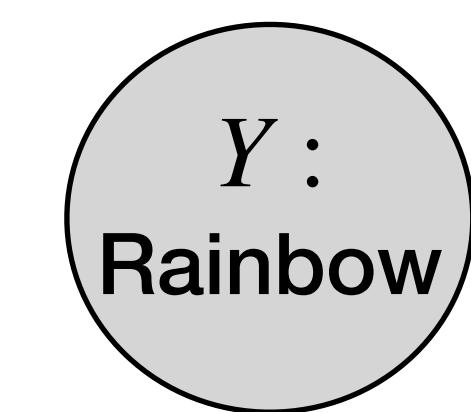
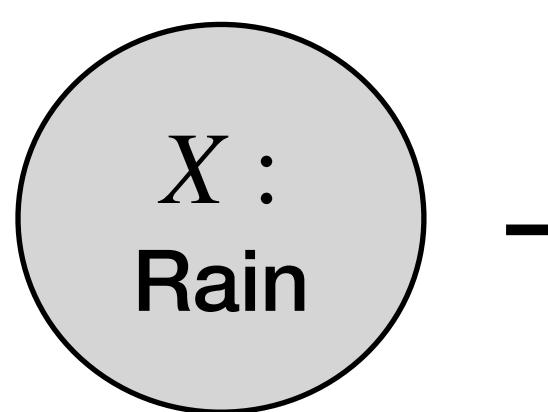
The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \dots$$

- **Nodes** : random variables
- **Links** : probabilistic relationships

Example :



$$P(X, Y) = \dots$$

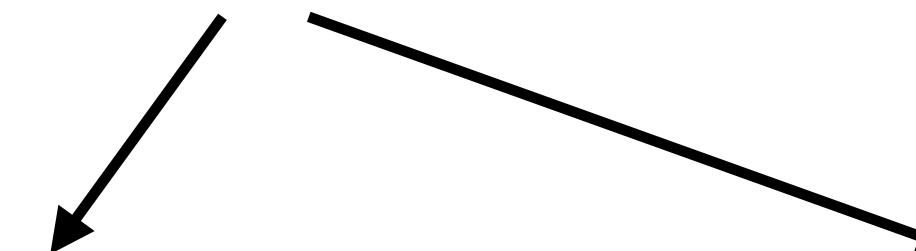
2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks

(Directed graphical models)



Markov random fields

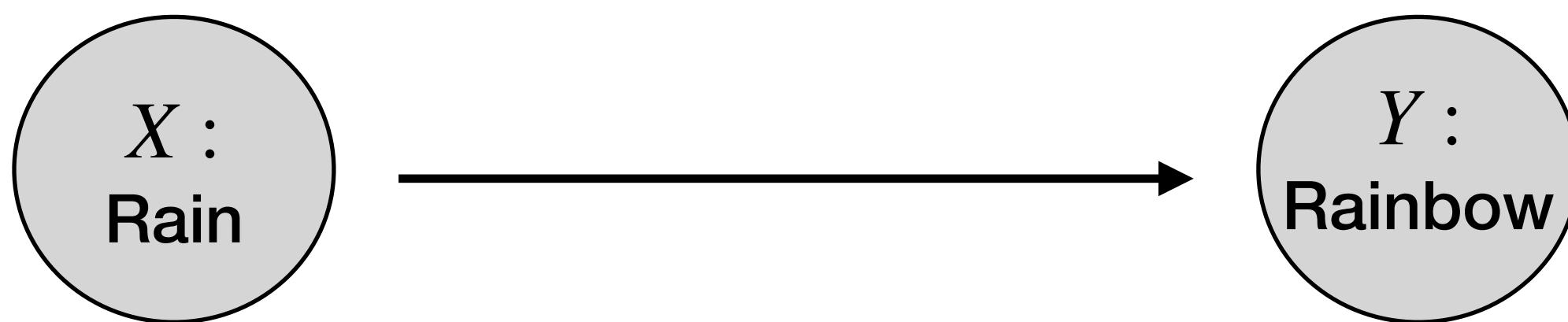
(Undirected graphical models)

The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \dots$$

Example :



- **Nodes** : random variables
- **Links** : probabilistic relationships

$$P(X, Y) = P(Y|X) \times P(X)$$

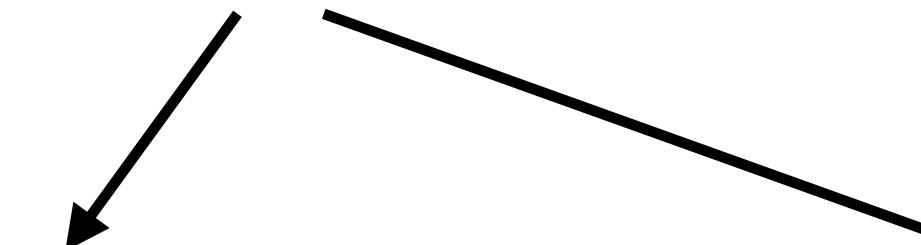
2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks

(Directed graphical models)



Markov random fields

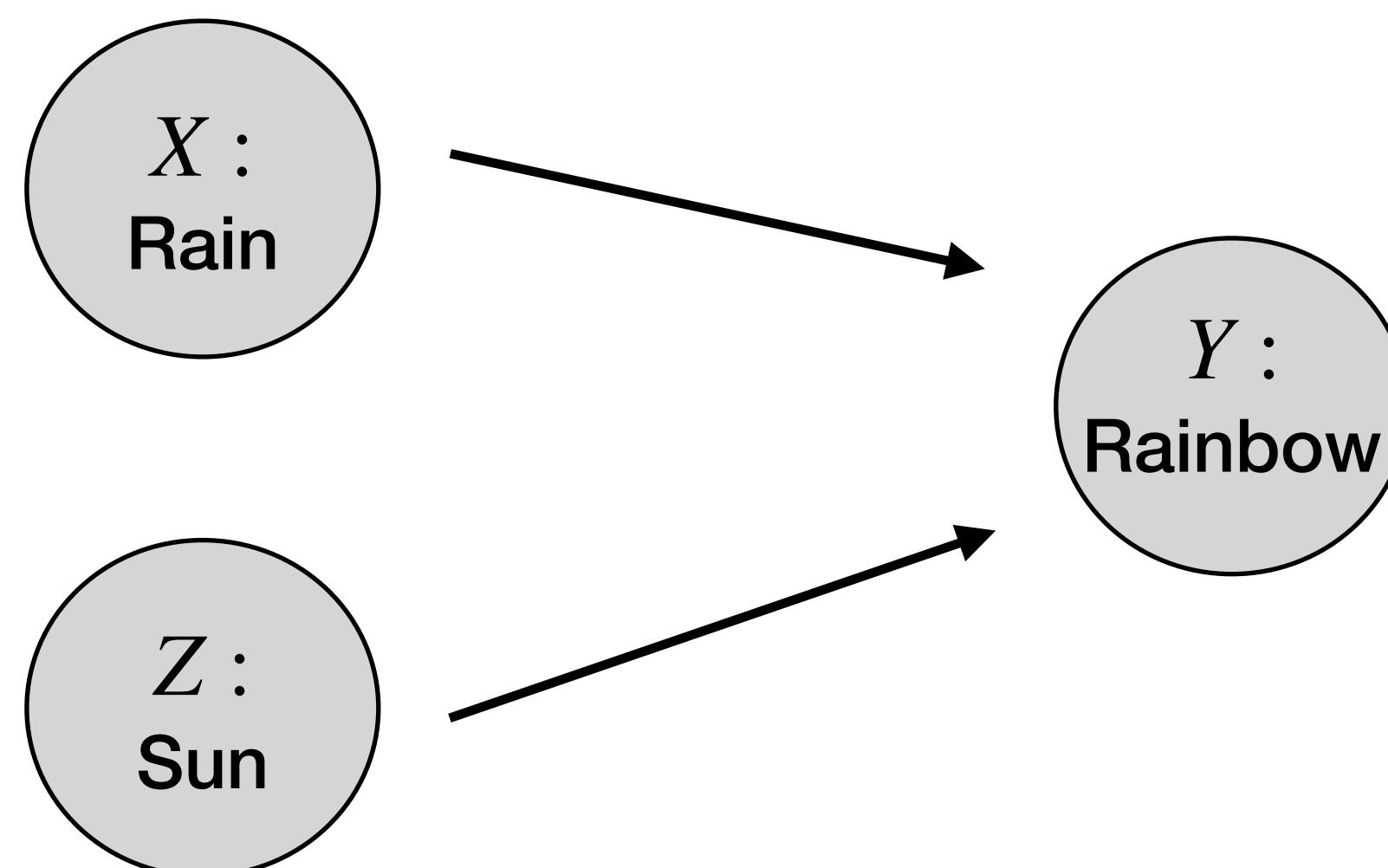
(Undirected graphical models)

The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \dots$$

Example :



$$P(X, Y, Z) = \dots$$

Features

X, Z conditionally
independent given Y

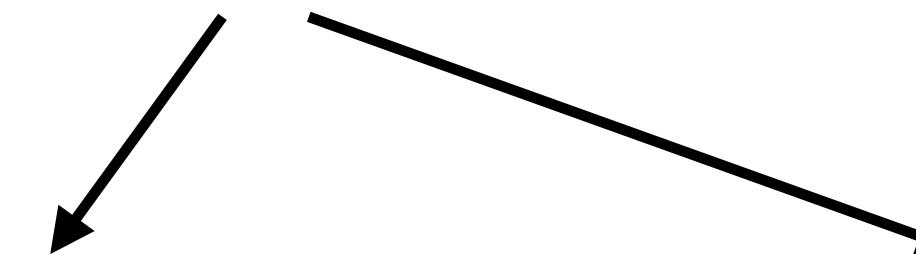
2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks

(Directed graphical models)



Markov random fields

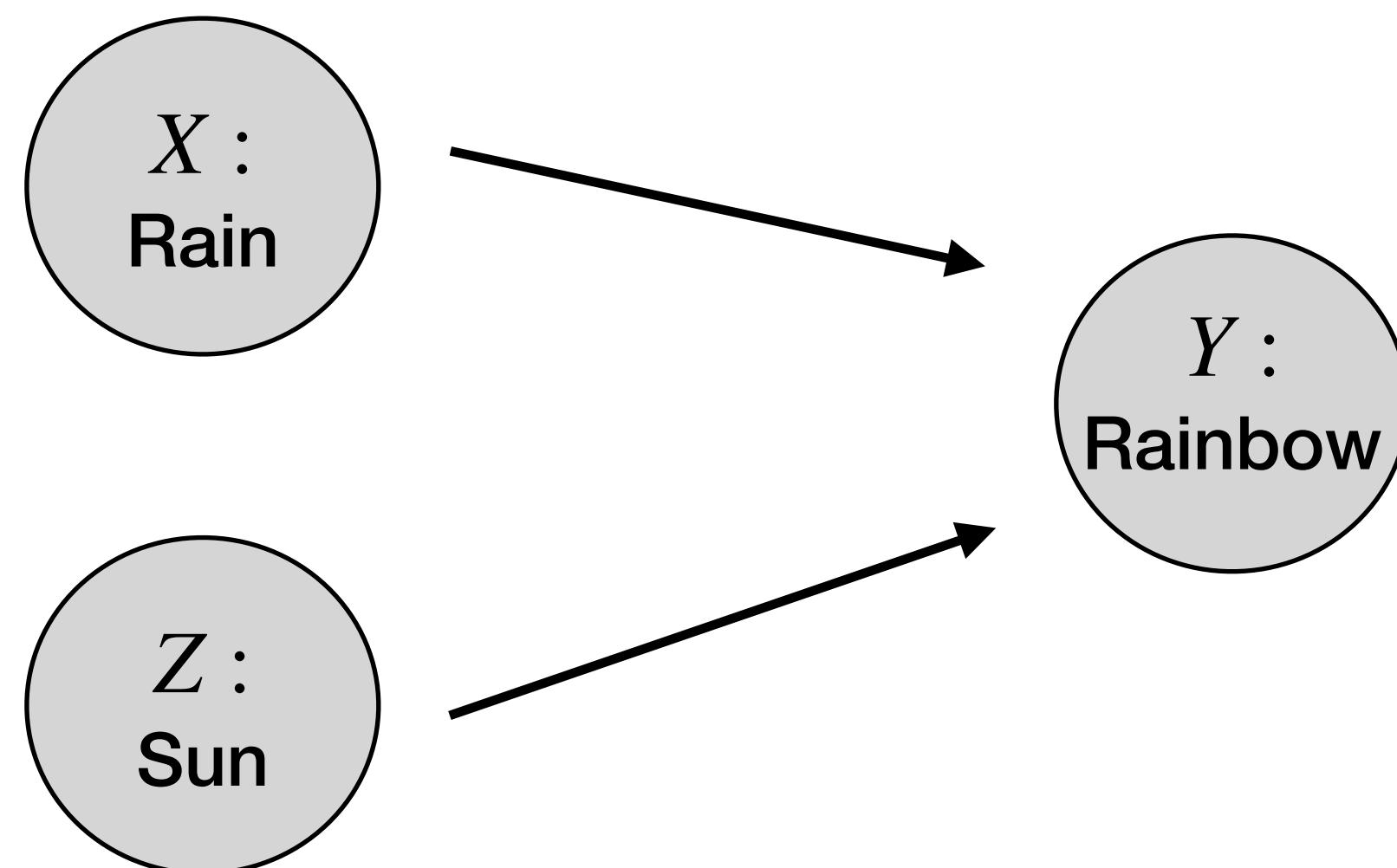
(Undirected graphical models)

The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \dots$$

Example :



$$\begin{aligned} P(X, Y, Z) &= P(Y|X, Z) \times P(X, Z) \\ &= P(Y|X, Z) \times P(X) \times P(Z) \end{aligned}$$

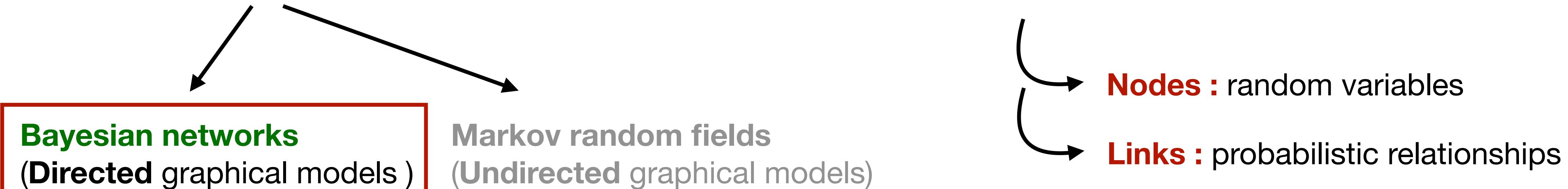
Features

X , Z conditionally
independent given Y

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

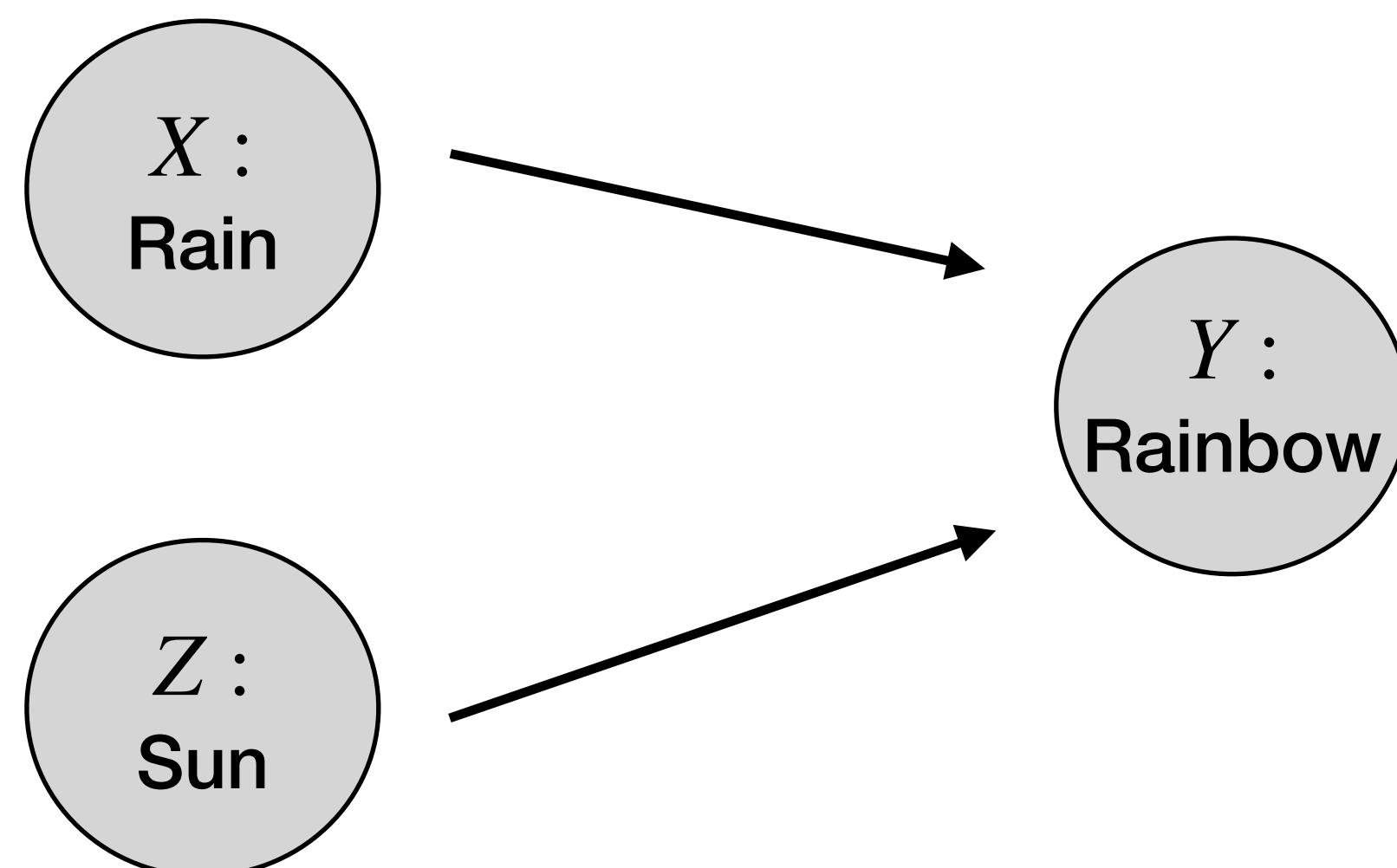


The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \prod_{i=1, \dots, N} P(X_i \mid \text{parents}(X_i))$$

Example :



Features
X , Z conditionally
independent given Y

$$\begin{aligned} P(X, Y, Z) &= P(Y|X, Z) \times P(X, Z) \\ &= P(Y|X, Z) \times P(X) \times P(Z) \end{aligned}$$

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks
(Directed graphical models)

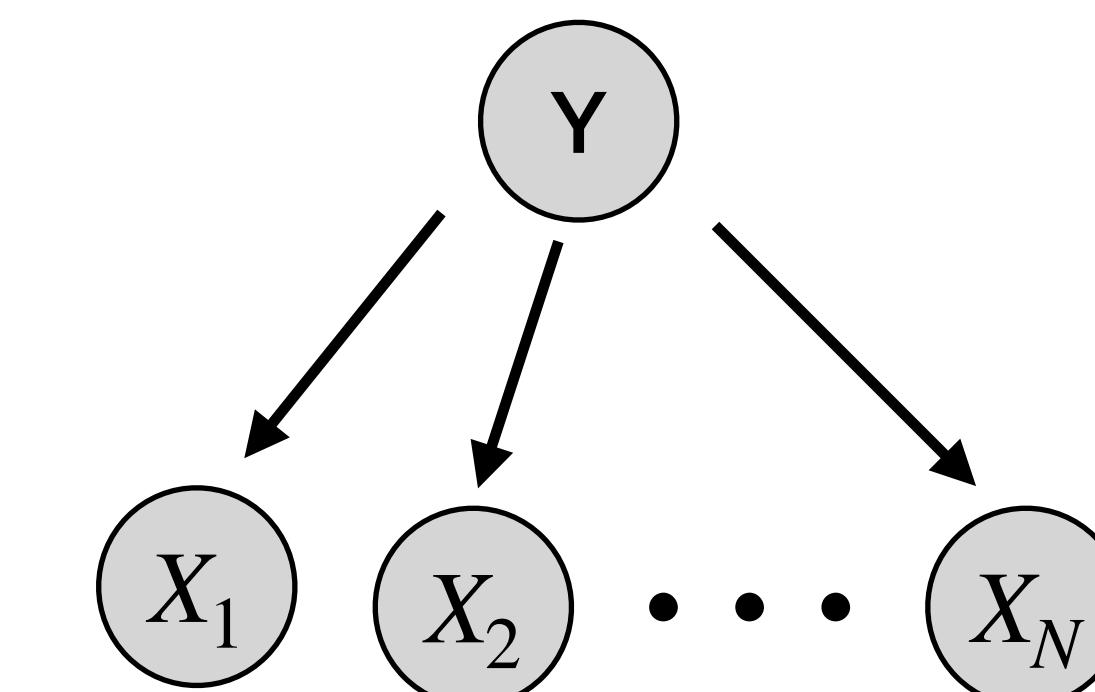
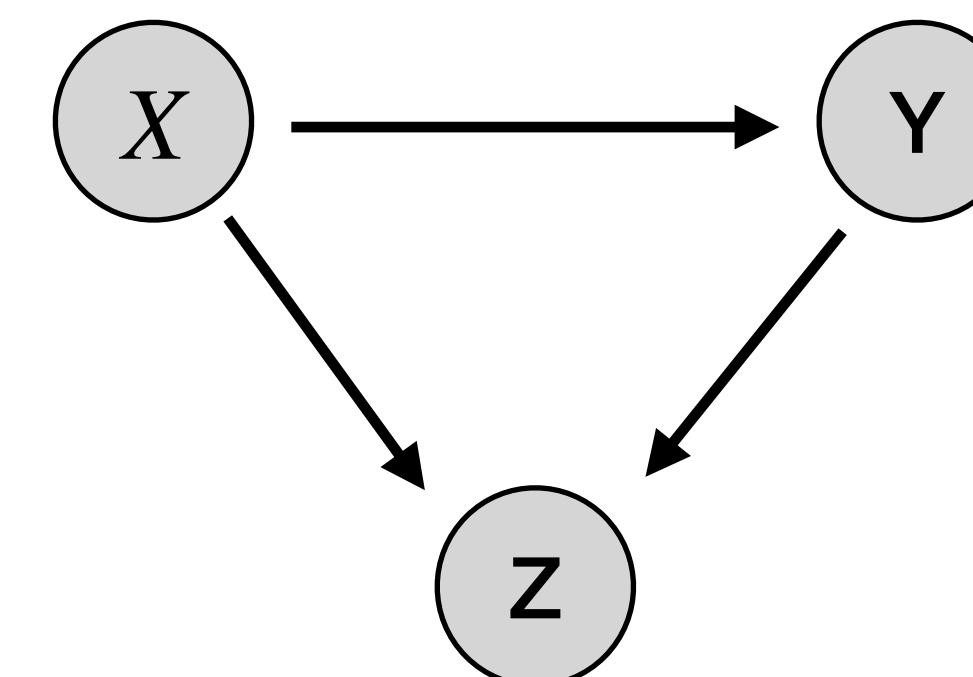
Markov random fields
(Undirected graphical models)

The focus of our course !

Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \prod_{i=1, \dots, N} P(X_i \mid \text{parents}(X_i))$$

Example :



$$P(X, Y) = P(Y|X) \times P(X)$$

$$P(X, Y, Z) = \dots$$

$$P(Y, X_1, \dots, X_N) =$$

- **Nodes** : random variables
- **Links** : probabilistic relationships

2. Probabilistic model

Probabilistic Graphical Model (PGM)

Probabilistic graphical models : analysis using **diagrammatic representations** of probability distributions

Bayesian networks
(Directed graphical models)

Markov random fields
(Undirected graphical models)

The focus of our course !

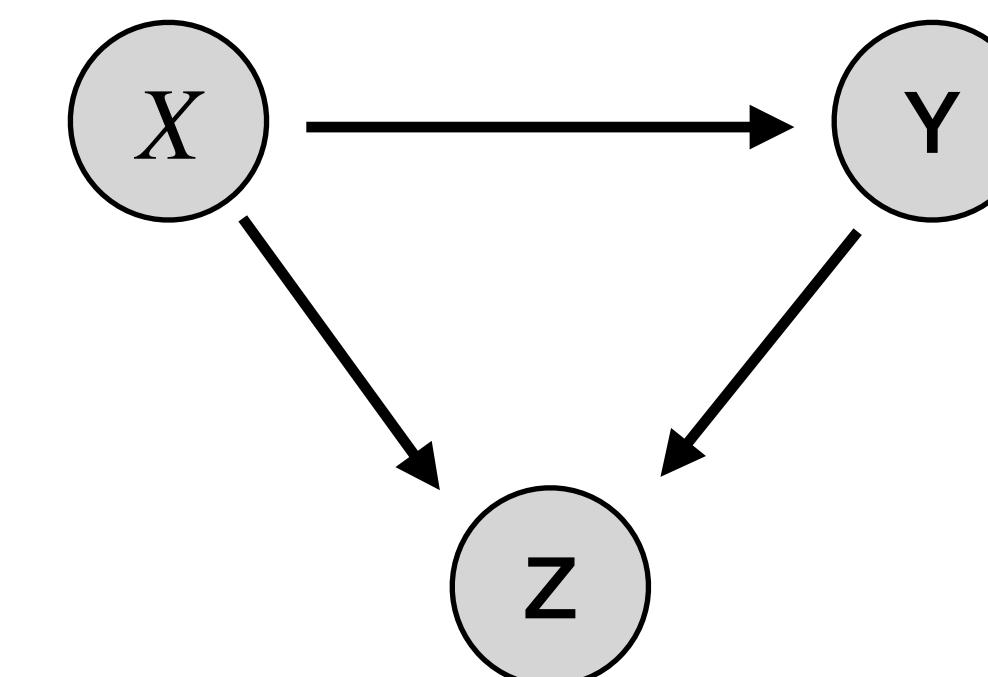
Model : joint probability over all variables

$$P(X_1, \dots, X_N) = \prod_{i=1, \dots, N} P(X_i | \text{parents}(X_i))$$

Example :

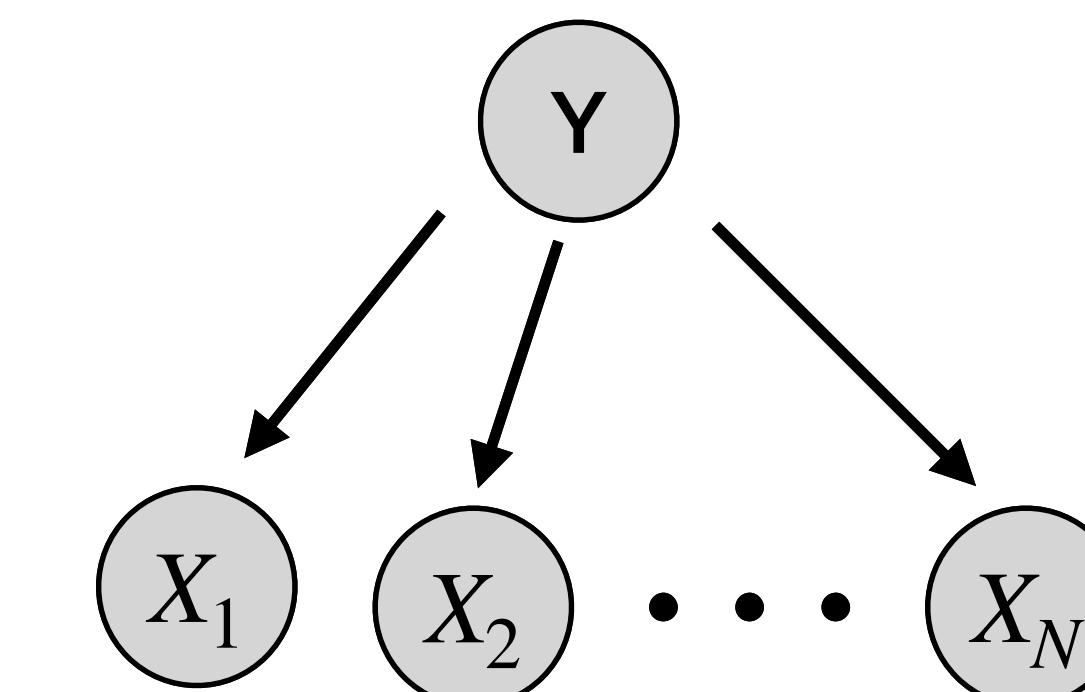


$$P(X, Y) = P(Y|X) \times P(X)$$



$$P(X, Y, Z) = P(Z|X, Y) \times P(X, Y)$$

$$= P(Z|X, Y) \times P(Y|X) \times P(X)$$



$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1, \dots, N} P(X_i | Y)$$

Diagrammatic representations:

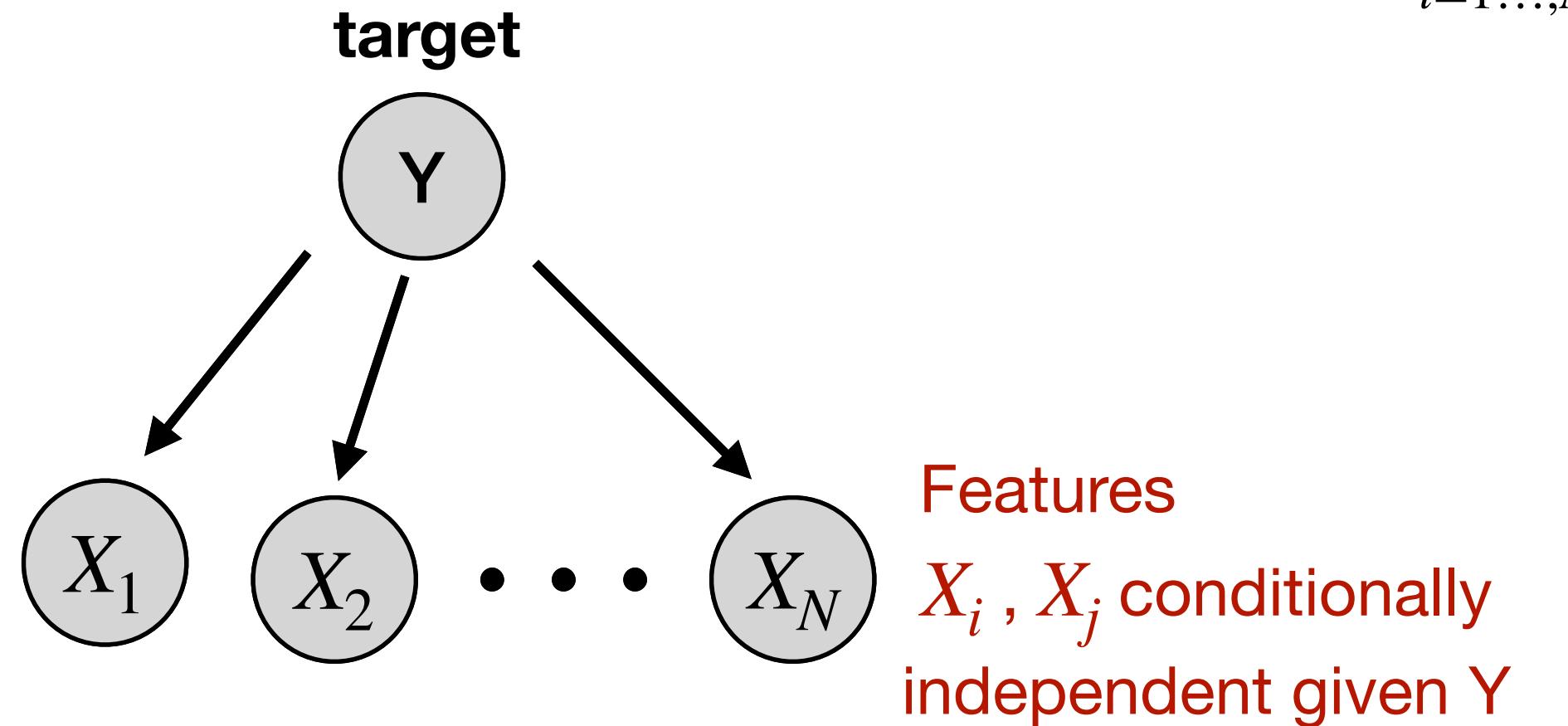
- **Nodes** : random variables
- **Links** : probabilistic relationships

2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$

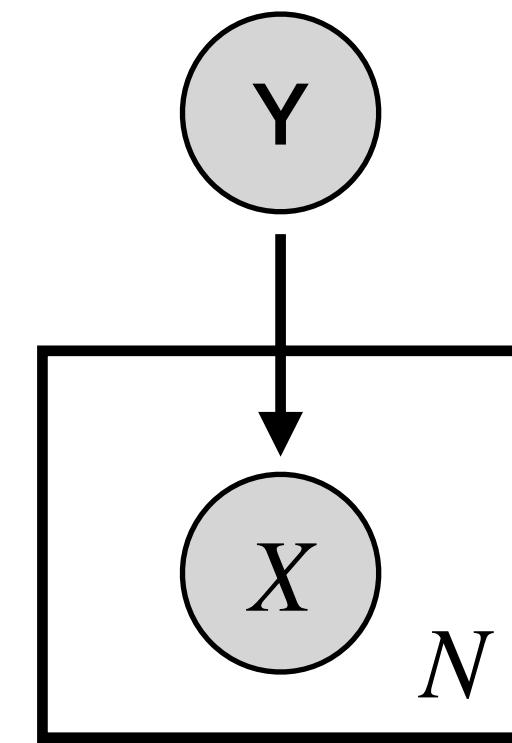
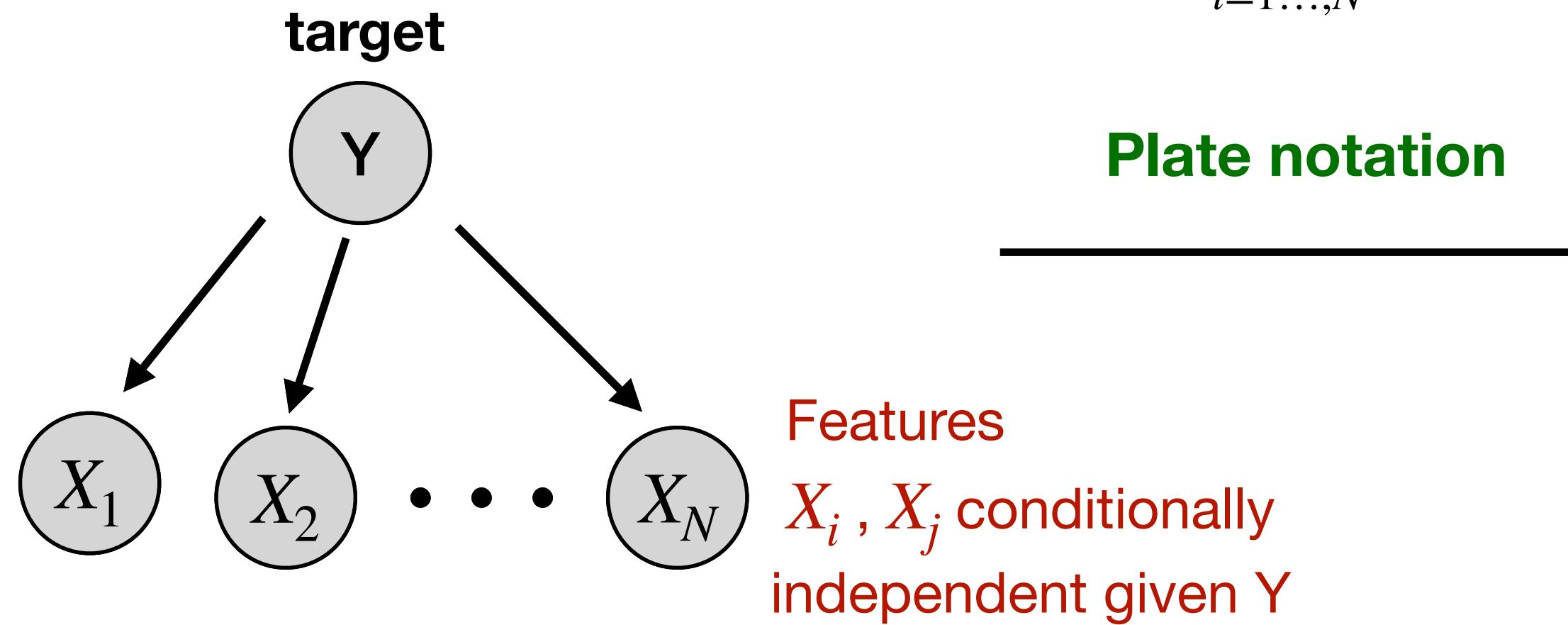


2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$

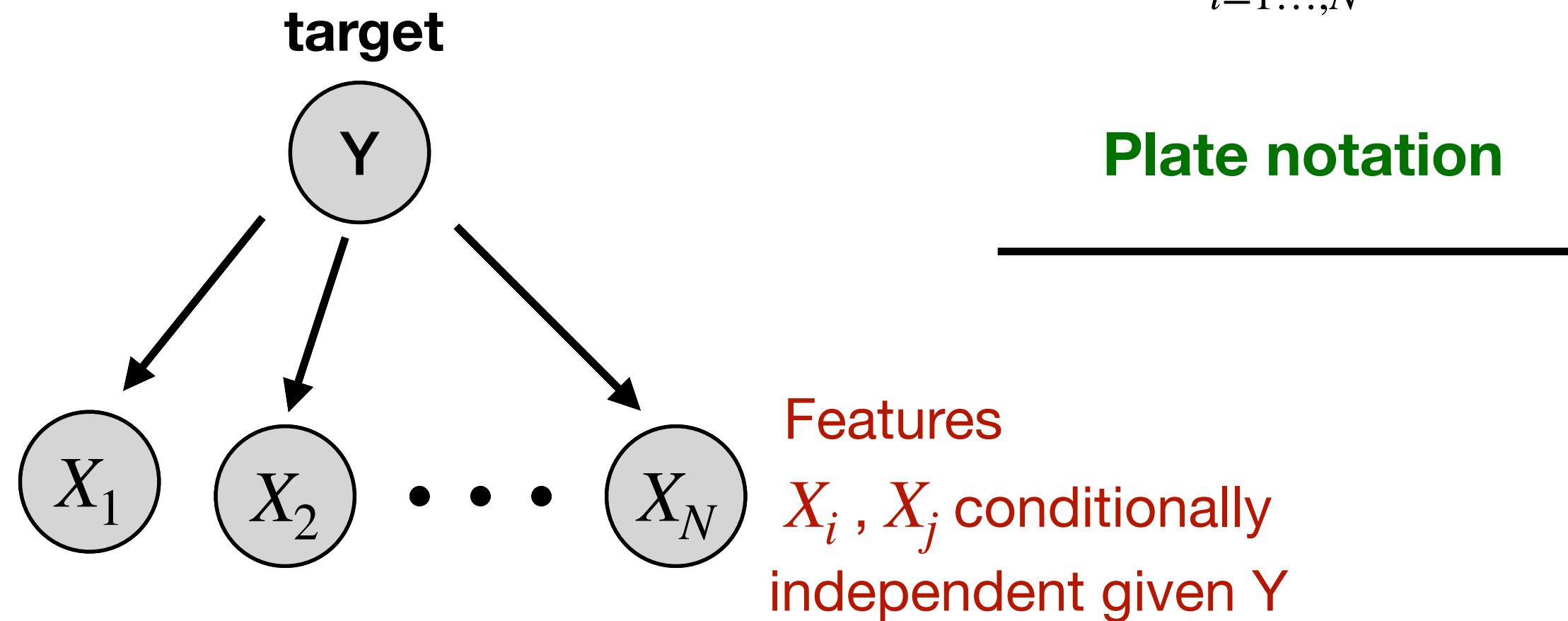


2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$



Reminder : Frequentist linear regression $x_i \in \mathbb{R}^d$ $y_1 \in \mathbb{R}$

Scalar notation :

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$y_i = x_i^T \theta + \epsilon_i$$

Matrix notation :

$$\mathbf{X} = (x_1, \dots, x_n) \text{ and } \mathbf{y} = (y_1, \dots, y_n)$$

$$y = \mathbf{X}^T \theta + \epsilon$$

2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$

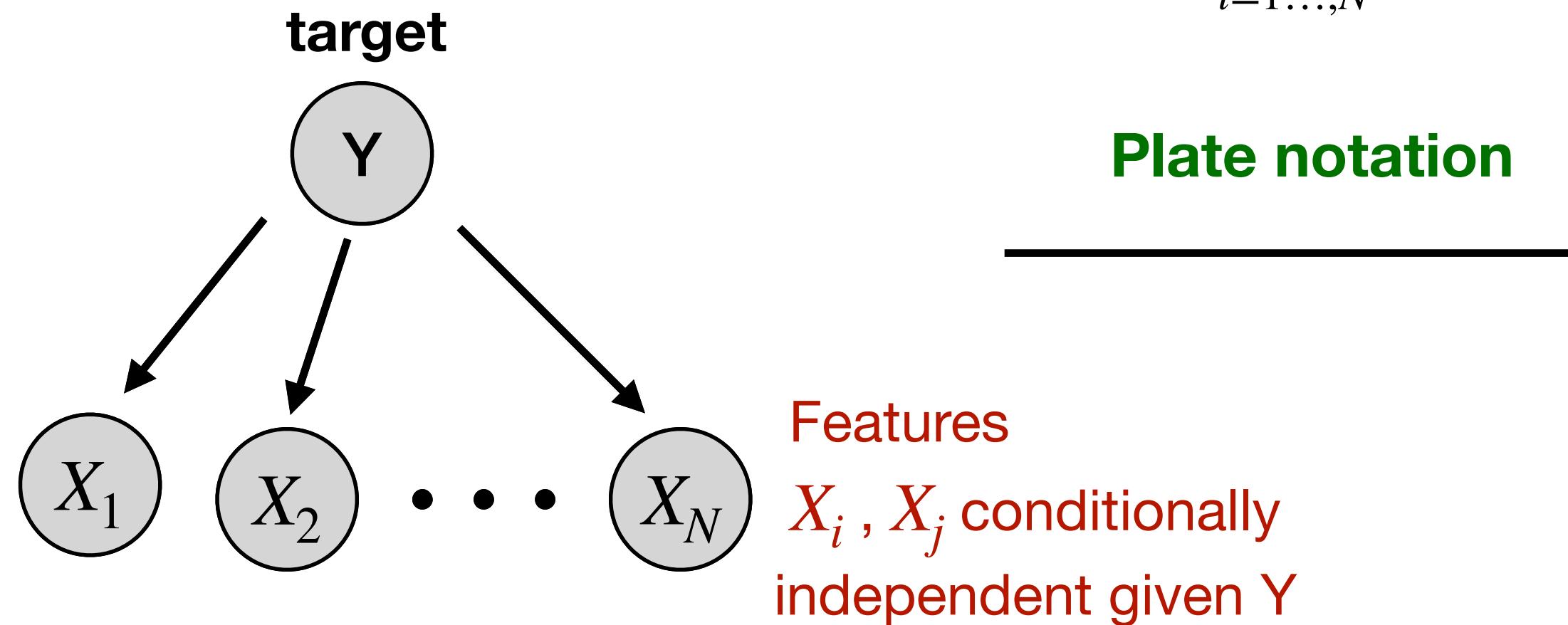
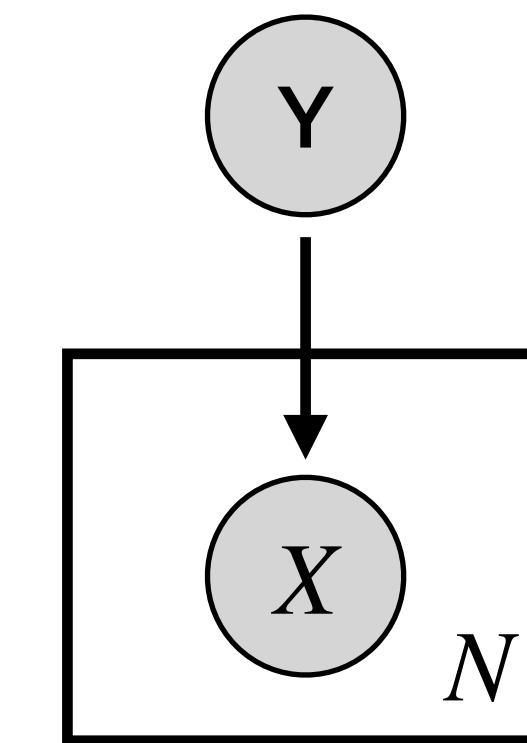


Plate notation



Reminder : Frequentist linear regression

$$x_i \in \mathbb{R}^d \quad y_i \in \mathbb{R}$$

Scalar notation :

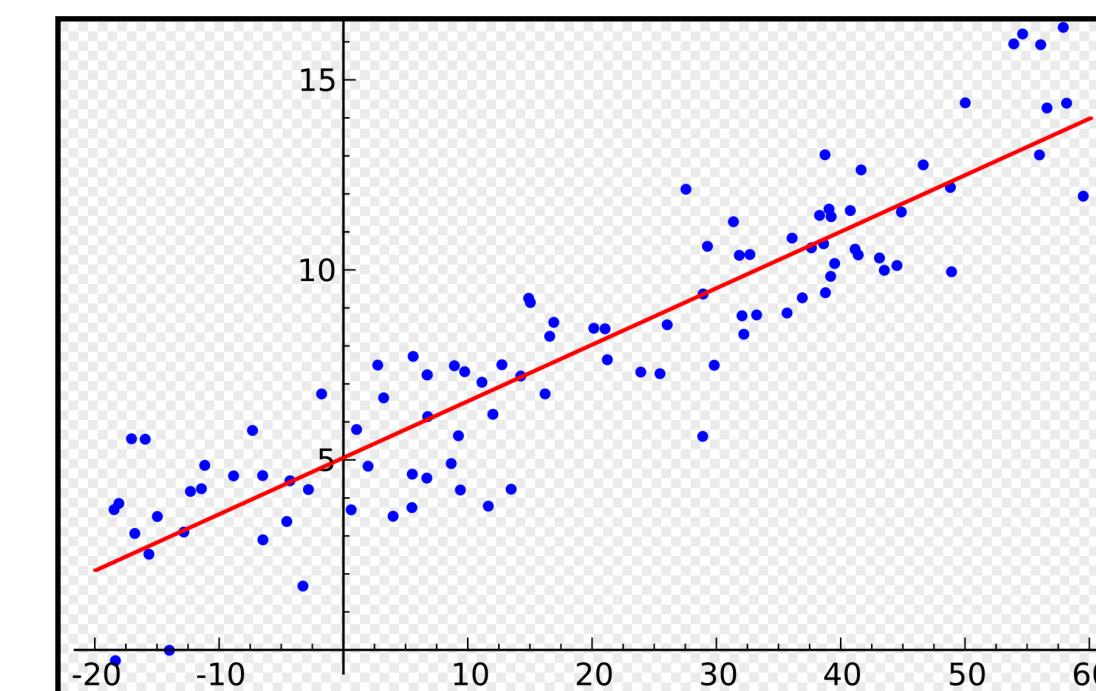
$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$

$$y_i = x_i^T \theta + \epsilon_i$$

Matrix notation :

$$\mathbf{X} = (x_1, \dots, x_n) \text{ and } \mathbf{y} = (y_1, \dots, y_n)$$

$$y = \mathbf{X}^T \theta + \epsilon$$



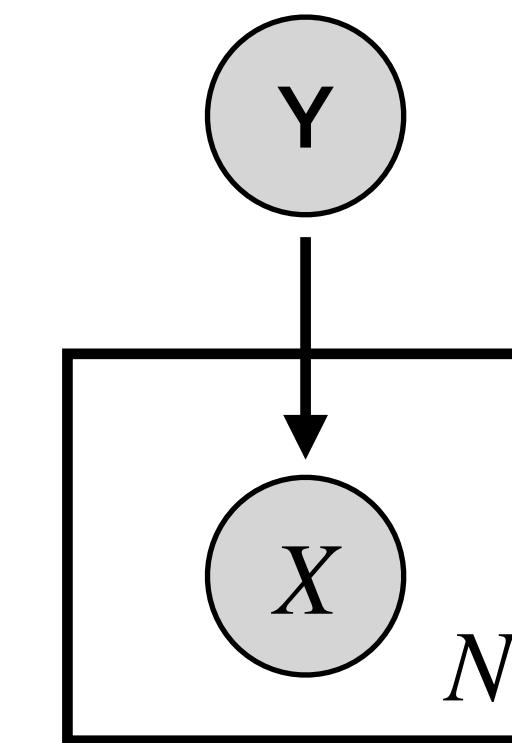
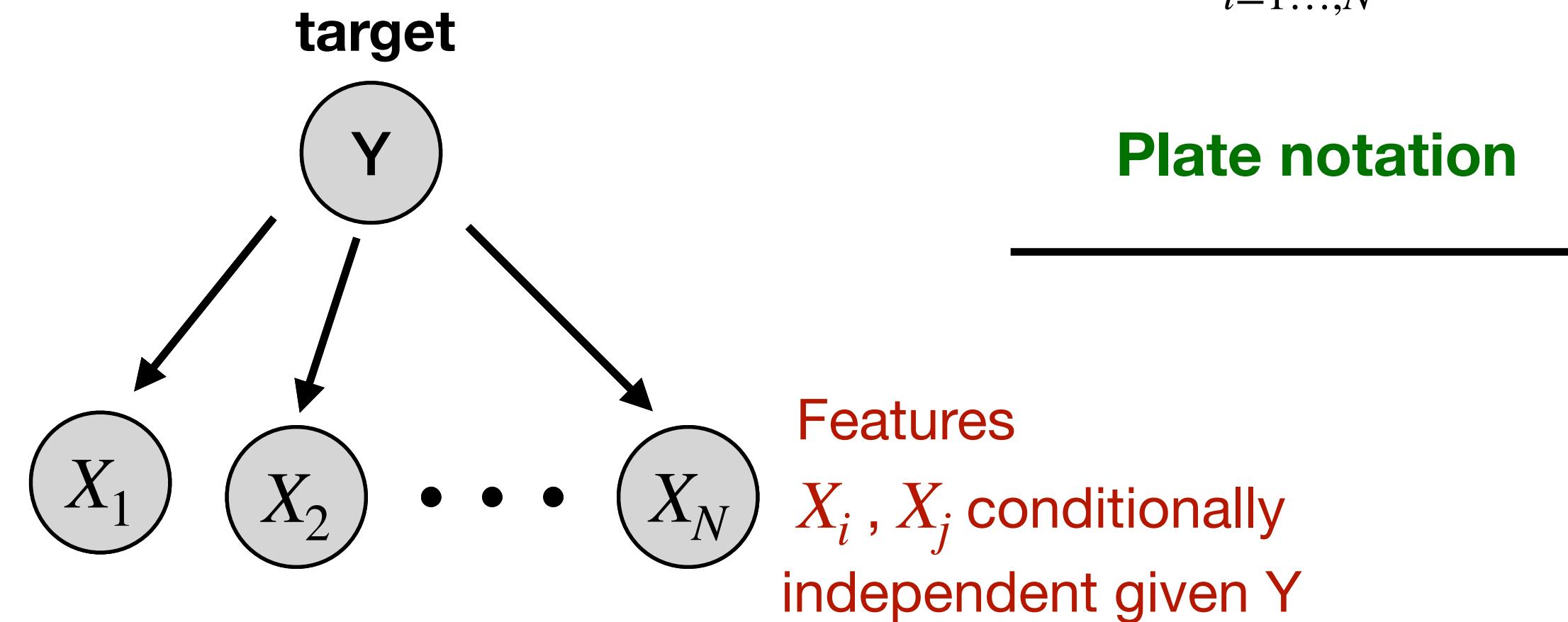
$$\min_{\theta} \|\theta^T \mathbf{X} - \mathbf{y}\|^2$$

2. Probabilistic model

Plates and examples of probabilistic model

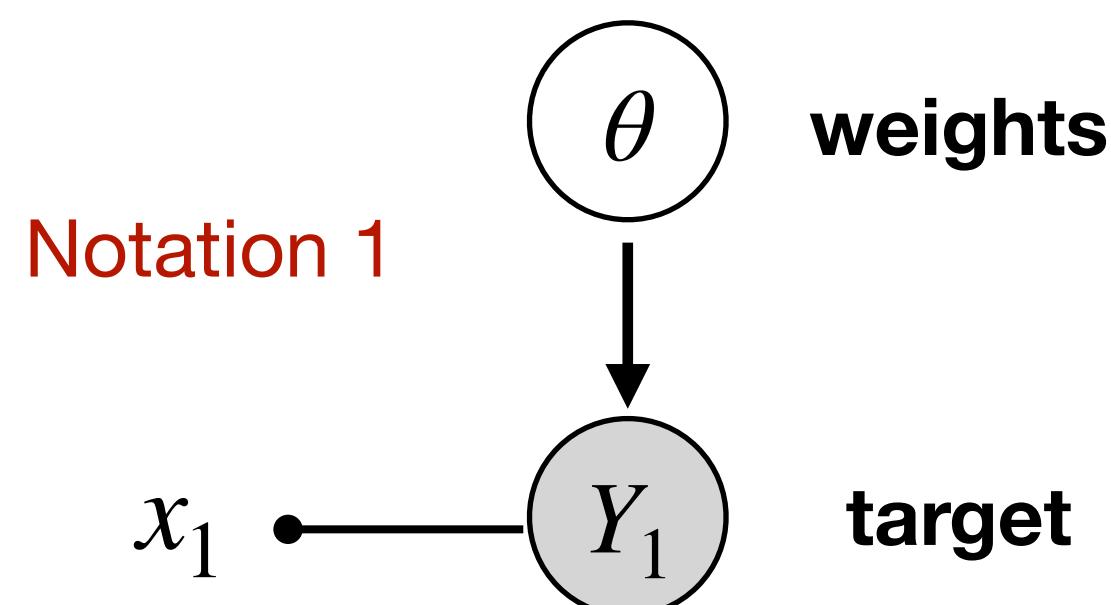
Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1\dots,N} P(X_i | Y)$$



Bayesian Linear regression

$$(x_1, y_1) \text{ with } x_1 \in \mathbb{R}^d, y_1 \in \mathbb{R}$$



$$\begin{aligned} P(\theta, y_1 | x_1) &= P(y_1 | \theta, x_1) \times P(\theta) \\ P(y_1 | \theta, x_1) &= \dots \\ P(\theta) &= \dots \end{aligned}$$

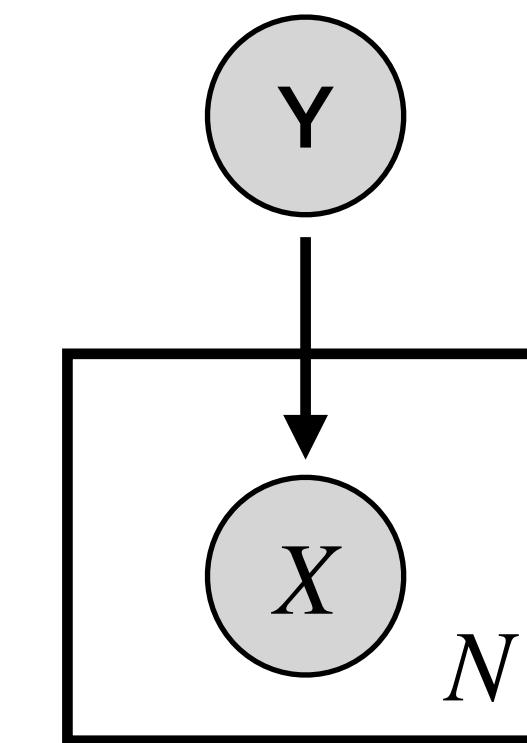
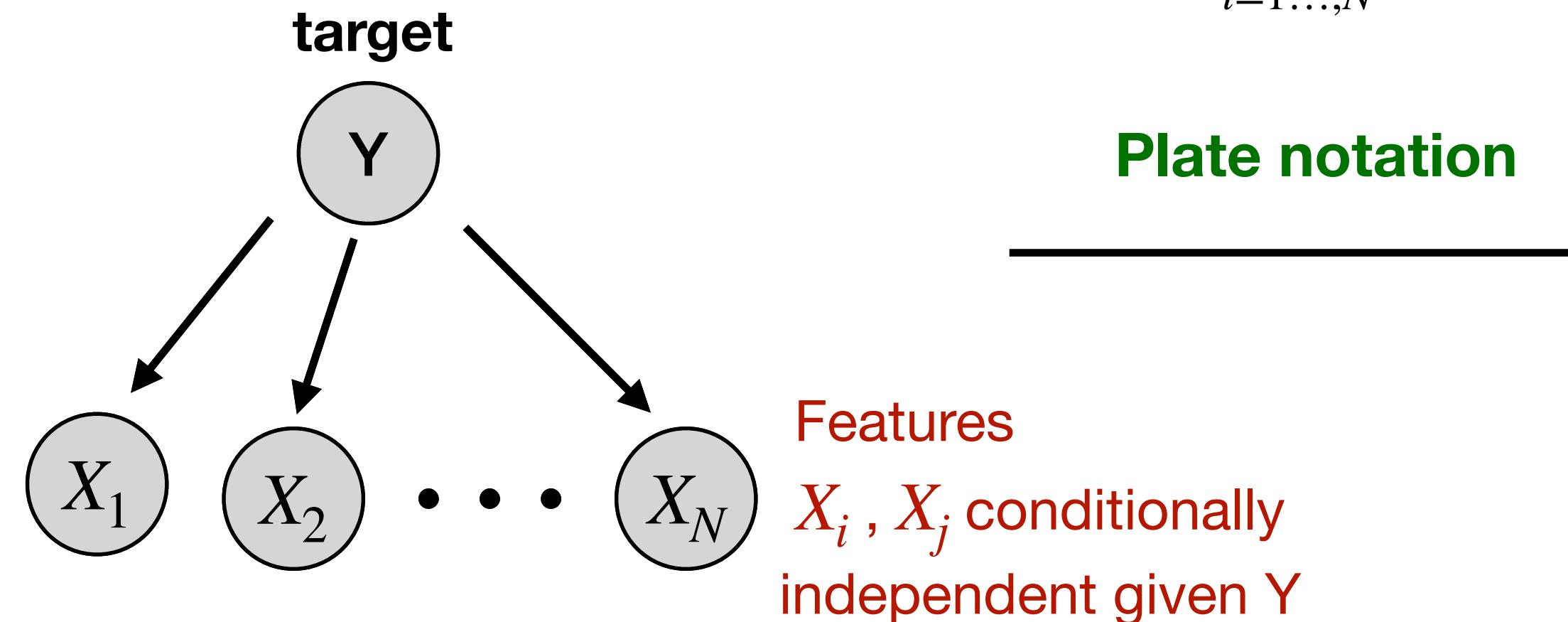


2. Probabilistic model

Plates and examples of probabilistic model

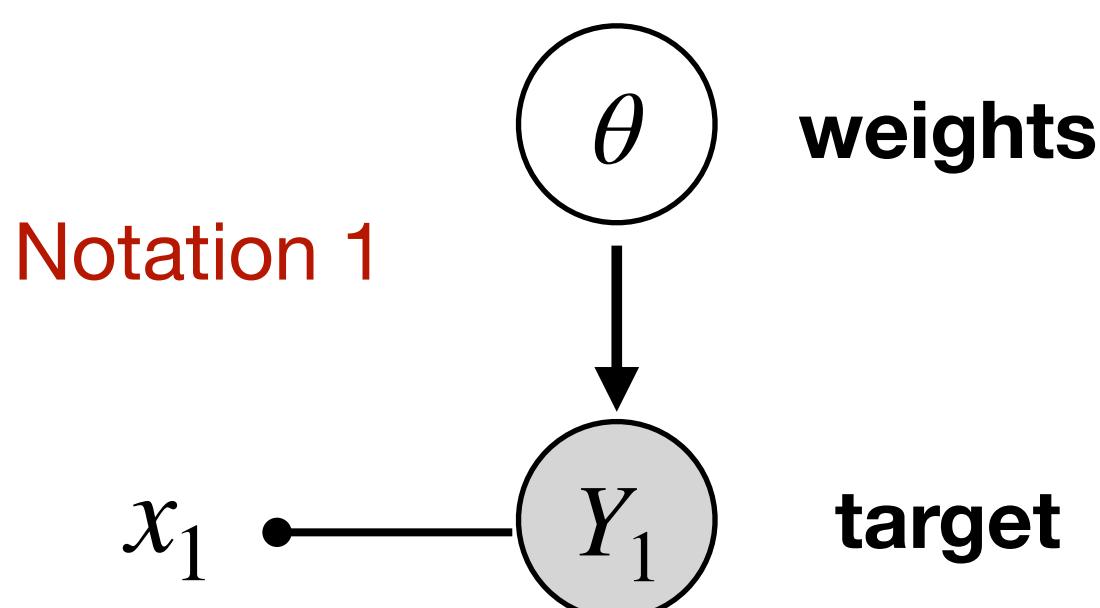
Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$



Bayesian Linear regression

$$(x_1, y_1) \text{ with } x_1 \in \mathbb{R}^d, y_1 \in \mathbb{R}$$



$$\begin{aligned} P(\theta, y_1 | x_1) &= P(y_1 | \theta, x_1) \times P(\theta) \\ P(y_1 | \theta, x_1) &= \mathcal{N}(y_1 | \theta^T x_1, \sigma^2) \\ P(\theta) &= \mathcal{N}(\theta | 0, \gamma^2) \end{aligned}$$

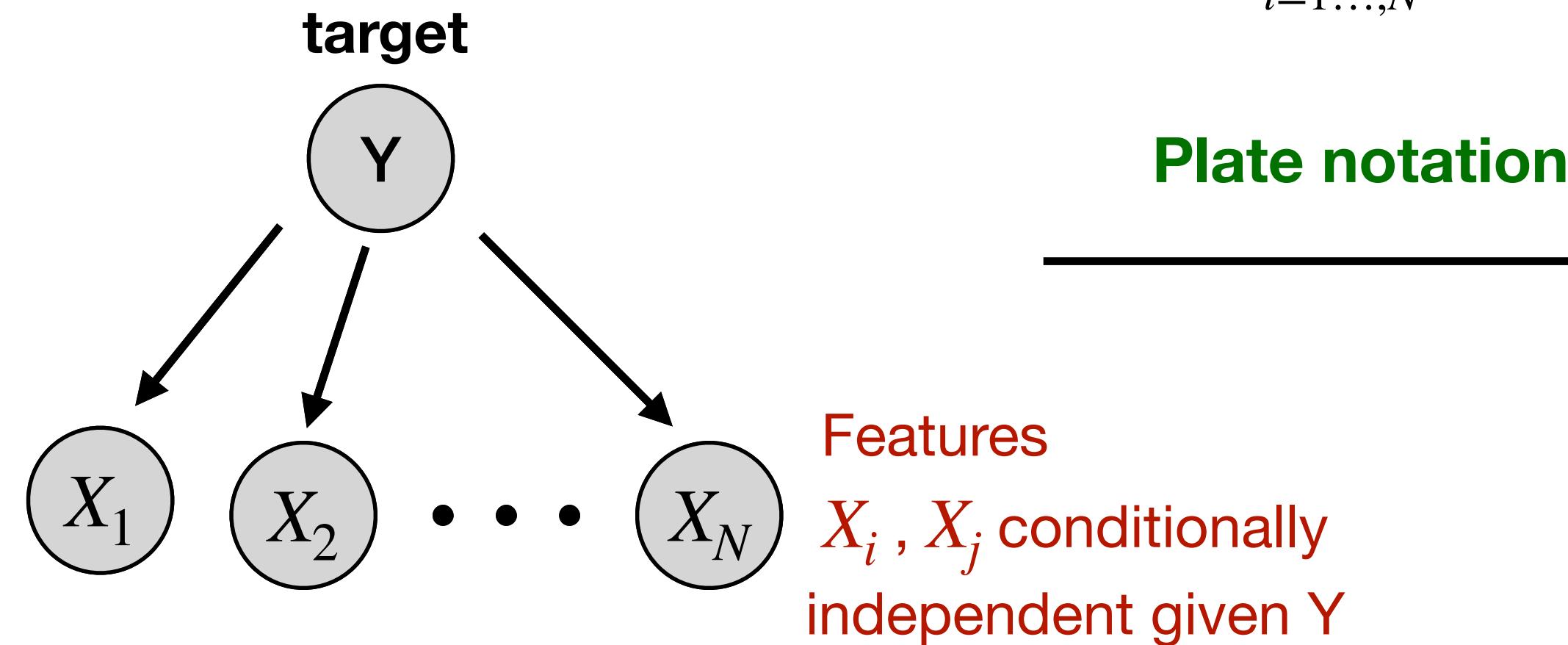


2. Probabilistic model

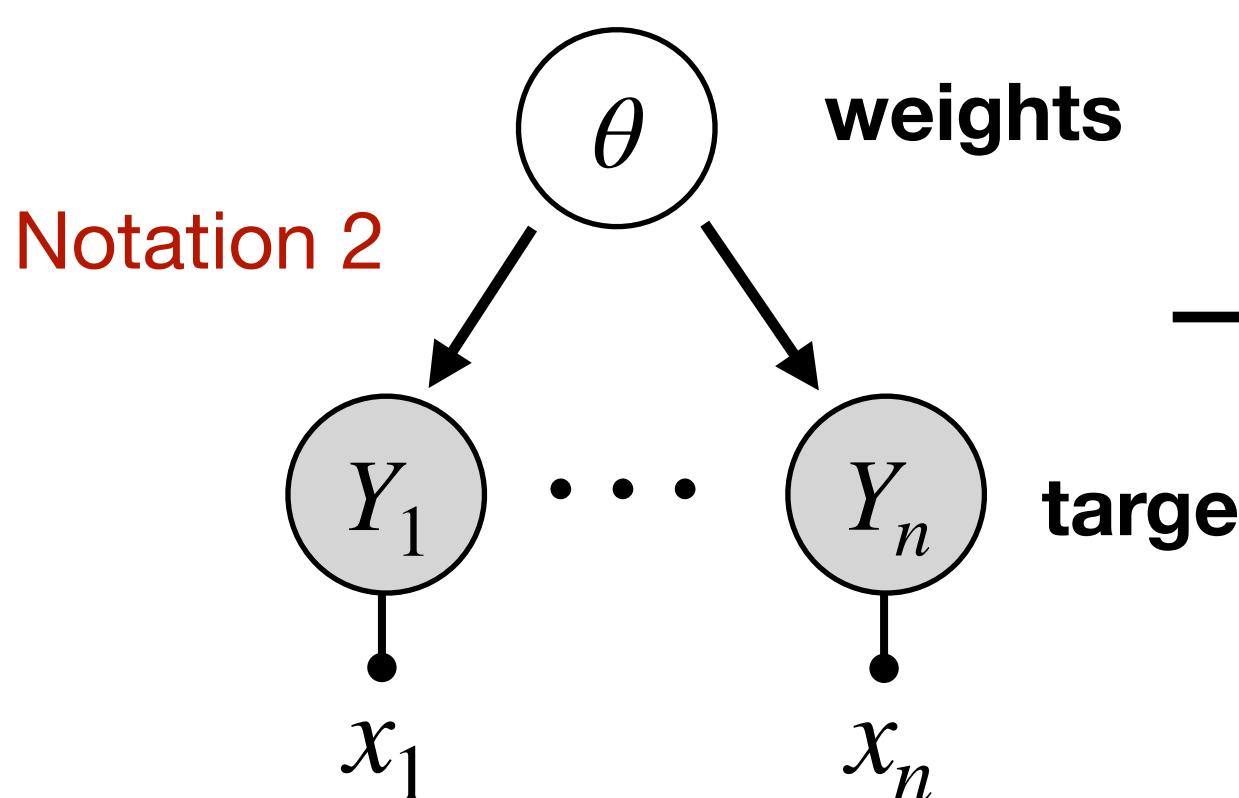
Plates and examples of probabilistic model

Naive Bayes Classifier

$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$



Bayesian Linear regression



$$D = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \text{with} \quad x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

$$P(\theta, y_i | x_i) = P(y_i | \theta, x_i) \times P(\theta)$$

$$P(y_i | \theta, x_i) = \mathcal{N}(y_i | \theta^T x_i, \sigma^2)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2)$$

Legend :

— Fixed variable

○ Hidden (latent) r.v.

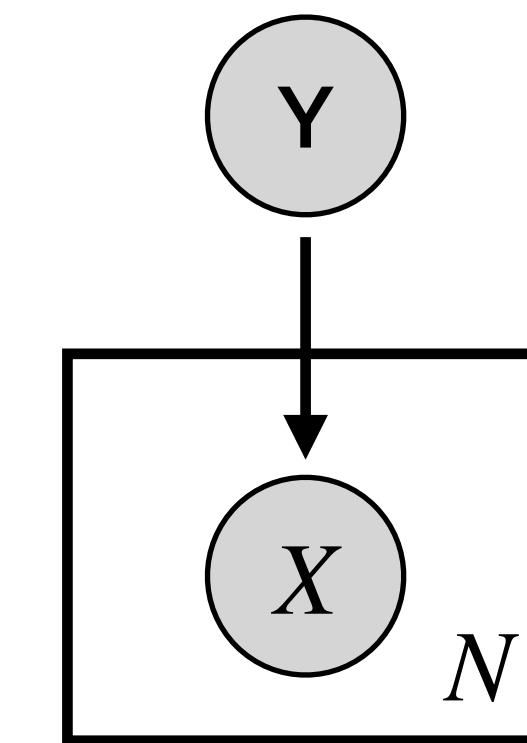
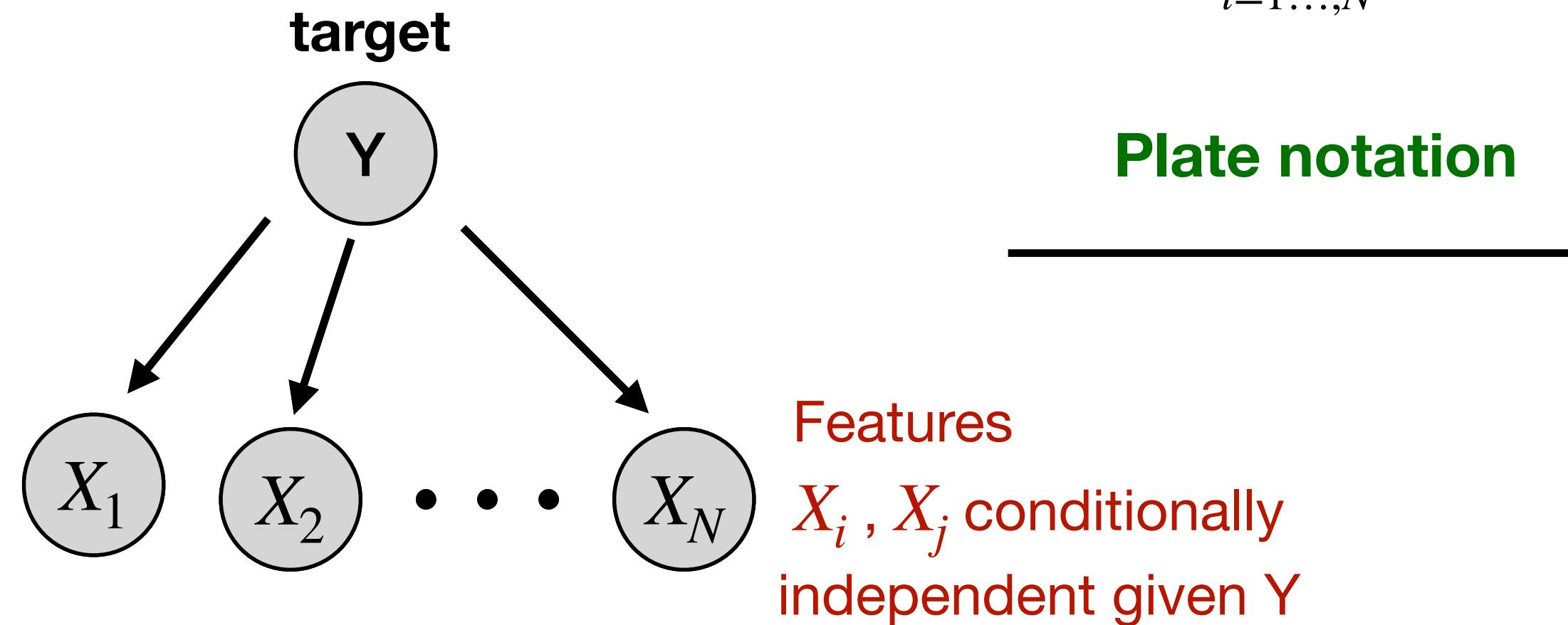
● Observed r.v.

2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

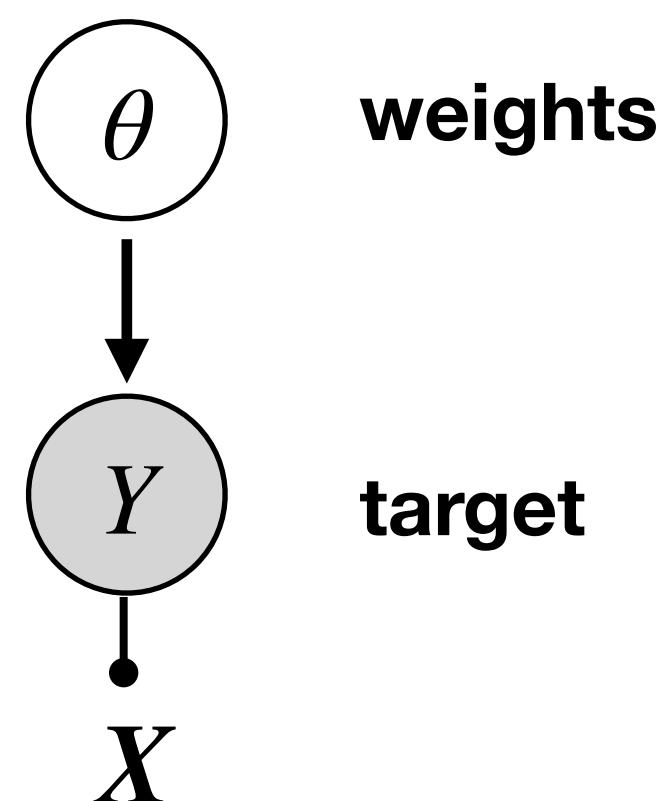
$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$



Bayesian Linear regression

$$\mathbf{X} = (x_1, \dots, x_n) \text{ and } \mathbf{y} = (y_1, \dots, y_n) \text{ with } x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

Notation 3



$$P(\theta, y | x) = P(y | \theta, x) \times P(\theta)$$

$$P(y | \theta, x) = \dots$$

$$P(\theta) = \dots$$

Legend :

— Fixed variable

○ Hidden (latent) r.v.

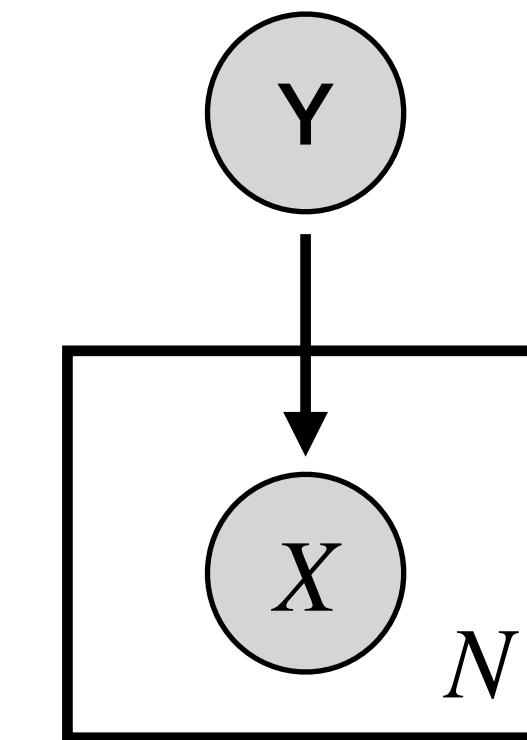
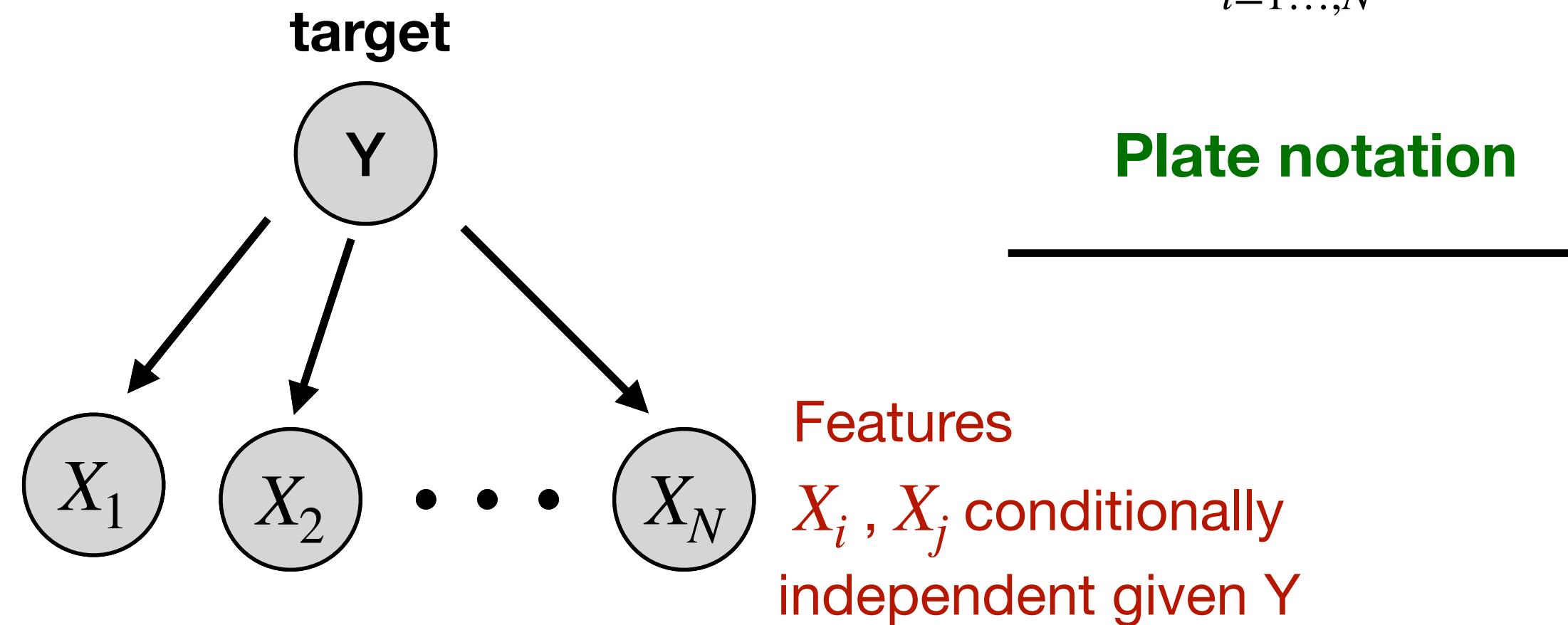
● Observed r.v.

2. Probabilistic model

Plates and examples of probabilistic model

Naive Bayes Classifier

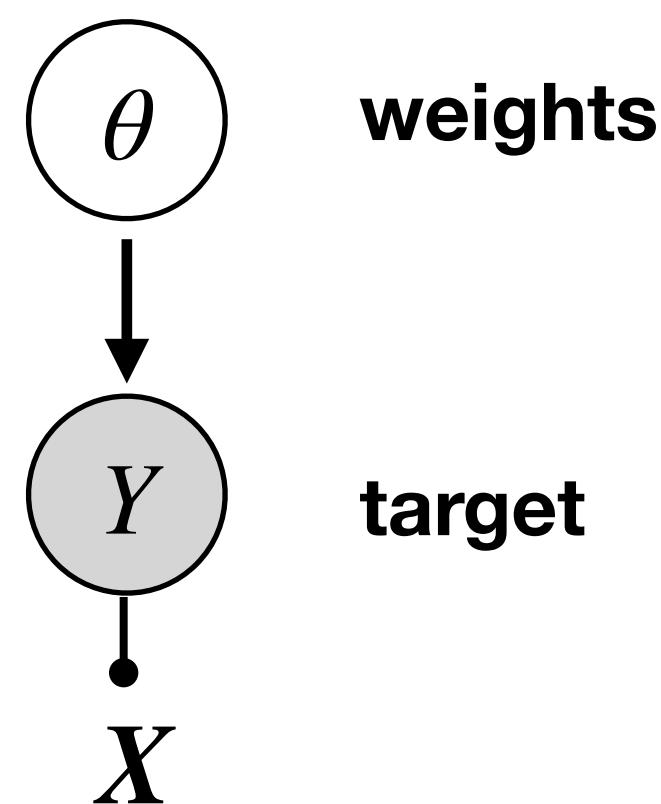
$$P(Y, X_1, \dots, X_N) = P(Y) \prod_{i=1 \dots, N} P(X_i | Y)$$



Bayesian Linear regression

$$\mathbf{X} = (x_1, \dots, x_n) \text{ and } \mathbf{y} = (y_1, \dots, y_n) \text{ with } x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

Notation 3



$$P(\theta, y | x) = P(y | \theta, x) \times P(\theta)$$

$$P(y | \theta, x) = \mathcal{N}(y | \theta^T x, \sigma^2 I_n)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2 I_n)$$

Legend :

| — Fixed variable

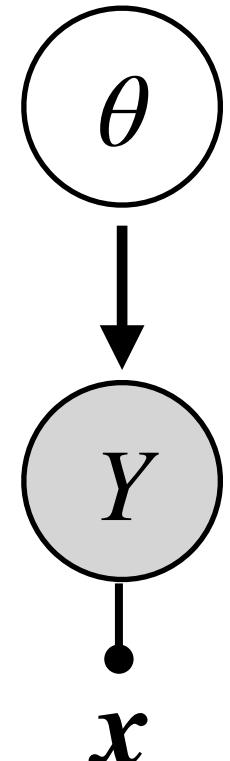
○ Hidden (latent) r.v.

● Observed r.v.

2. Probabilistic model

Linear regression

Bayesian Linear regression



$$P(\theta, y | X) = P(y | \theta, X) \times P(\theta)$$

$$P(y | \theta, X) = \mathcal{N}(y | \theta^T X, \sigma^2 I_n)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2 I_n)$$

Objective :

Frequentist linear regression

Objective : $\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \|\theta^T X - y\|^2$

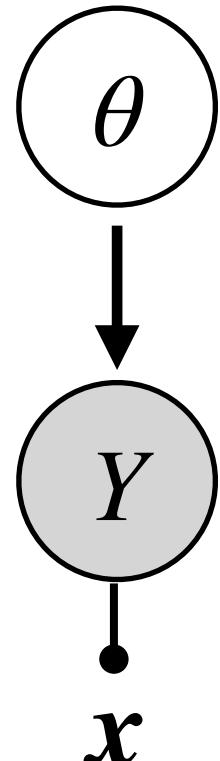
$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$$

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

2. Probabilistic model

Linear regression

Bayesian Linear regression



$$P(\theta, y | X) = P(y | \theta, X) \times P(\theta)$$

$$P(y | \theta, X) = \mathcal{N}(y | \theta^T X, \sigma^2 I_n)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2 I_n)$$

Objective : $\arg \max_{\theta} P(\theta | X, y) = \arg \max_{\theta} P(\theta, y | X)$

Frequentist linear regression

Objective : $\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \|\theta^T X - y\|^2$

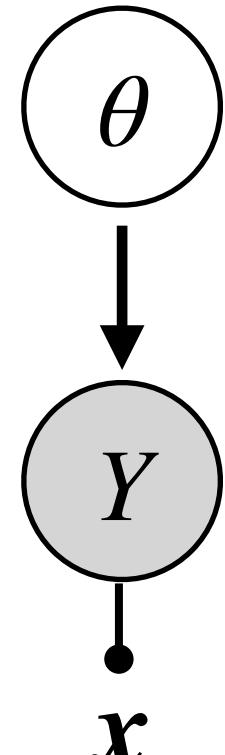
$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$$

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

2. Probabilistic model

Linear regression

Bayesian Linear regression



$$P(\theta, y | X) = P(y | \theta, X) \times P(\theta)$$

$$P(y | \theta, X) = \mathcal{N}(y | \theta^T X, \sigma^2 I_n)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2 I_n)$$

Objective : $\arg \max_{\theta} P(\theta | X, y) = \arg \max_{\theta} P(\theta, y | X)$

Frequentist linear regression

Objective : $\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \|\theta^T X - y\|^2$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$$

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

Theorem :

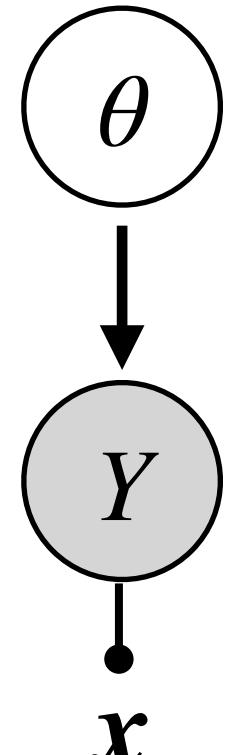
There exists $\lambda \in \mathbb{R}$ such that : $\arg \max_{\theta} P(\theta | X, y) = \arg \min_{\theta} \left\{ \|\theta^T X - y\|^2 + \lambda \|\theta\|^2 \right\}$

So by adding a normal prior on the weight we turned this problem into a L_2 regularised problem

2. Probabilistic model

Linear regression

Bayesian Linear regression



$$P(\theta, y | X) = P(y | \theta, X) \times P(\theta)$$

$$P(y | \theta, X) = \mathcal{N}(y | \theta^T X, \sigma^2 I_n)$$

$$P(\theta) = \mathcal{N}(\theta | 0, \gamma^2 I_n)$$

Objective : $\arg \max_{\theta} P(\theta | X, y) = \arg \max_{\theta} P(\theta, y | X)$

Frequentist linear regression

Objective : $\min_{\theta} \mathcal{L}(\theta) = \min_{\theta} \|\theta^T X - y\|^2$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta)$$

$$\hat{\theta} = (X^T X)^{-1} X^T y$$

Theorem :

There exists $\lambda \in \mathbb{R}$ such that : $\arg \max_{\theta} P(\theta | X, y) = \arg \min_{\theta} \left\{ \|\theta^T X - y\|^2 + \lambda \|\theta\|^2 \right\}$

So by adding a normal prior on the weight we turned this problem into a L_2 regularised problem

Proof : see the whiteboard in class or left as an exercise

3

Analytical Inference

3. Analytical Inference

Reminder of posterior distribution

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

3. Analytical Inference

Reminder of posterior distribution

Posterior distribution

The diagram illustrates the formula for the posterior distribution:

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

The components are labeled as follows:

- Likelihood**: Fixed by model
- Prior**: Fixed by us
- Evidence**: Fixed by data

3. Analytical Inference

Reminder of posterior distribution

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Posterior

Fixed by model Likelihood Prior Fixed by us

Evidence
HARD TO COMPUTE

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

Fixed by data

3. Analytical Inference

Maximum a posteriori (MAP) : definition & remarks

Posterior distribution

The diagram illustrates the formula for the posterior distribution:

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

The components are labeled as follows:

- Likelihood**: Fixed by model
- Prior**: Fixed by us
- Evidence**: Fixed by data
- HARD TO COMPUTE**: A yellow box containing the term $P(X)$.

Below the formula, the evidence is shown as:

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

Remarks

- We have to **avoid computing** the evidence
- Naive approach : **maximum a posteriori** ,
$$\hat{\theta}_{MAP} = \arg \max_{\theta} \left\{ \frac{P(\theta | X) \cdot P(\theta)}{P(X)} \right\}$$
$$= \arg \max_{\theta} P(X | \theta) \cdot P(\theta)$$
- This maximization can be done with numerical **optimization** problem

3. Analytical Inference

Maximum a posteriori (MAP) : limitations

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Posterior

Likelihood **Prior**

Evidence
HARD TO COMPUTE

Fixed by model

Fixed by us

Fixed by data

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

Remarks

- We have to **avoid computing** the evidence
- Naive approach : **maximum a posteriori** ,
$$\hat{\theta}_{MAP} = \arg \max_{\theta} \left\{ \frac{P(\theta | X) \cdot P(\theta)}{P(X)} \right\}$$
$$= \arg \max_{\theta} P(X | \theta) \cdot P(\theta)$$
- This maximization can be done with numerical **optimization** problem

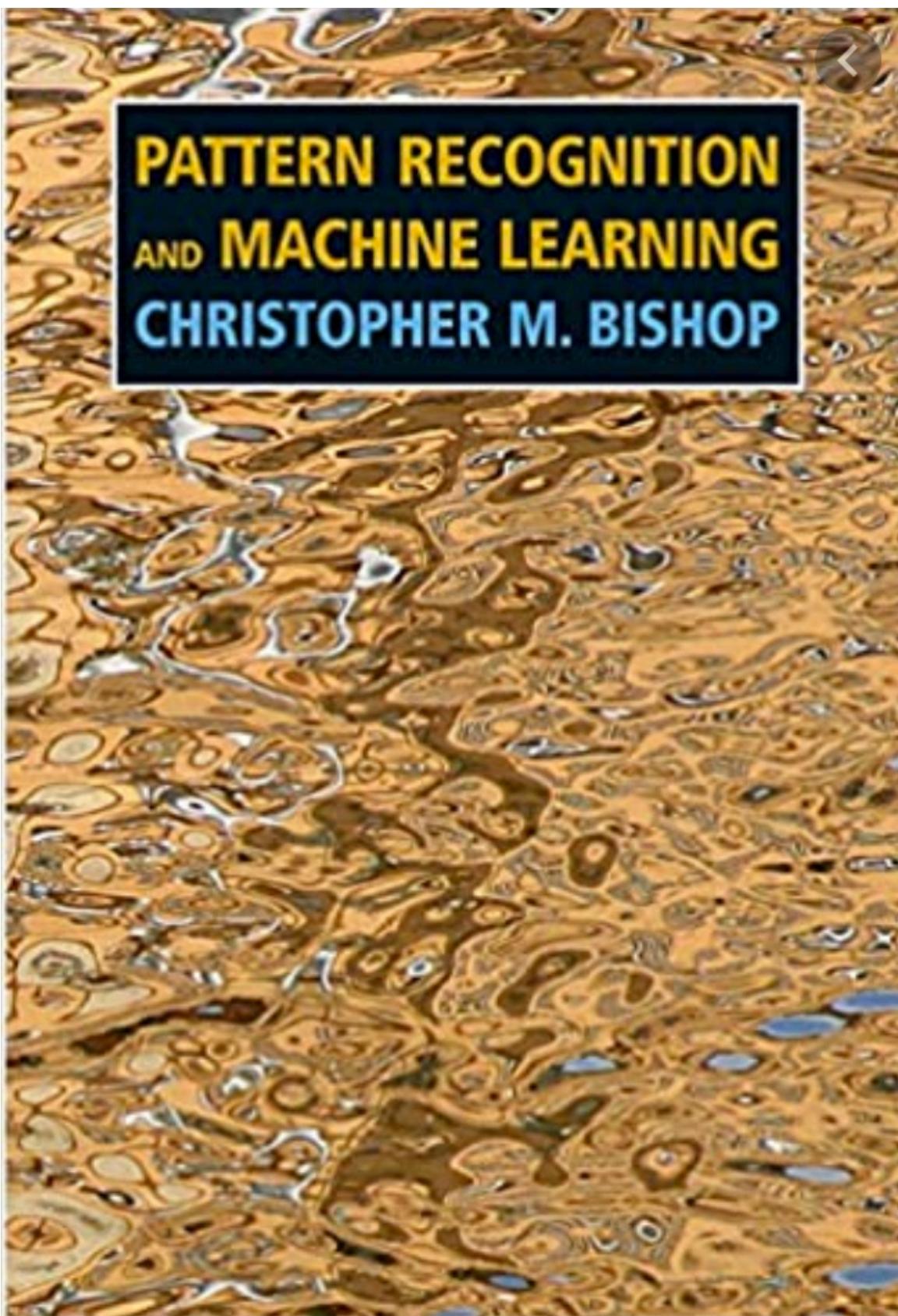
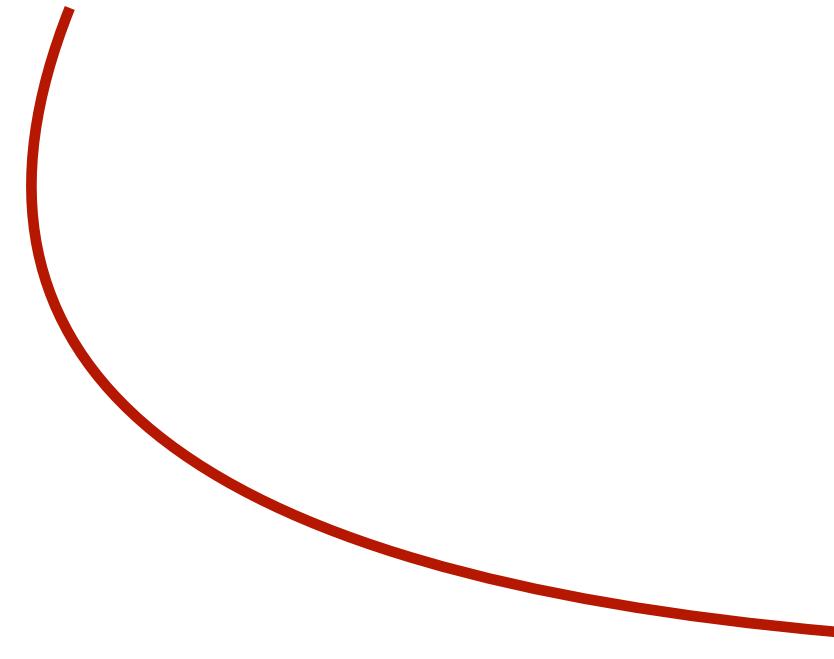
Limitations (among many others)

1. in general, **not representative of bayesian** methods : $\hat{\theta}_{MAP}$ is a point estimate like $\hat{\theta}_{MLE}$
 - can't compute **credible intervals** because it doesn't return a pdf/pmf (not a bayesian inference)
2. **can't use online learning** : the prior is not well updated

3. Analytical Inference

Maximum a posteriori (MAP) : book

For more theoretical details (and example on analytical inference) :



4

Conjugate distributions

4. Conjugate distributions

Conjugate distributions : avoid computing evidence

Posterior distribution

The diagram illustrates the formula for the posterior distribution:

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

The components are labeled as follows:

- Likelihood**: Fixed by model
- Prior**: Fixed by us
- Evidence**: Fixed by data

4. Conjugate distributions

Conjugate distributions : avoid computing evidence

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Posterior

Fixed by model

Likelihood

Prior

Fixed by us

Evidence

Fixed by data

Remarks

- We have to **avoid computing** the evidence
- We can choose a **convenient prior** which enable us to compute the posterior :
Conjugate prior

4. Conjugate distributions

Conjugate distributions : avoid computing evidence

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Posterior

The diagram illustrates the components of the posterior distribution formula. The term $P(X, \theta)$ is labeled "Fixed by model". The term $P(\theta)$ is labeled "Fixed by us". The term $P(X)$ is labeled "Fixed by data". The terms $P(X | \theta)$ and $P(\theta)$ are grouped together and labeled "Likelihood" and "Prior" respectively, while the term $P(X)$ is labeled "Evidence".

Remarks

- We have to **avoid computing** the evidence
- We can choose a **convenient prior** which enable us to compute the **posterior** :
Conjugate prior

Conjugate prior

$P(\theta)$ is **conjugate** to $P(X | \theta)$ if the $P(\theta)$ and $P(X | \theta)$ lie in the same family of distributions (gaussian for example)

4. Conjugate distributions

Conjugate distributions : avoid computing evidence

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X|\theta) \times P(\theta)}{P(X)}$$

Posterior

Likelihood Prior

Evidence

Fixed by model Fixed by us Fixed by data

Remarks

- We have to **avoid computing** the evidence
- We can choose a **convenient prior** which enable us to compute the posterior :
Conjugate prior

Conjugate prior

$P(\theta)$ is **conjugate** to $P(X|\theta)$ if the $P(\theta)$ and $P(X|\theta)$ lie in the same family of distributions (gaussian for example)

Example

$$P(\theta | X) = \frac{\mathcal{N}(\theta | \mu_{prior}, \sigma^2_{prior})}{P(X)} \times P(\theta)$$

$\mathcal{N}(\theta | \mu_{posterior}, \sigma^2_{posterior})$

In the context of a gaussian, the prior for the mean is a gaussian !

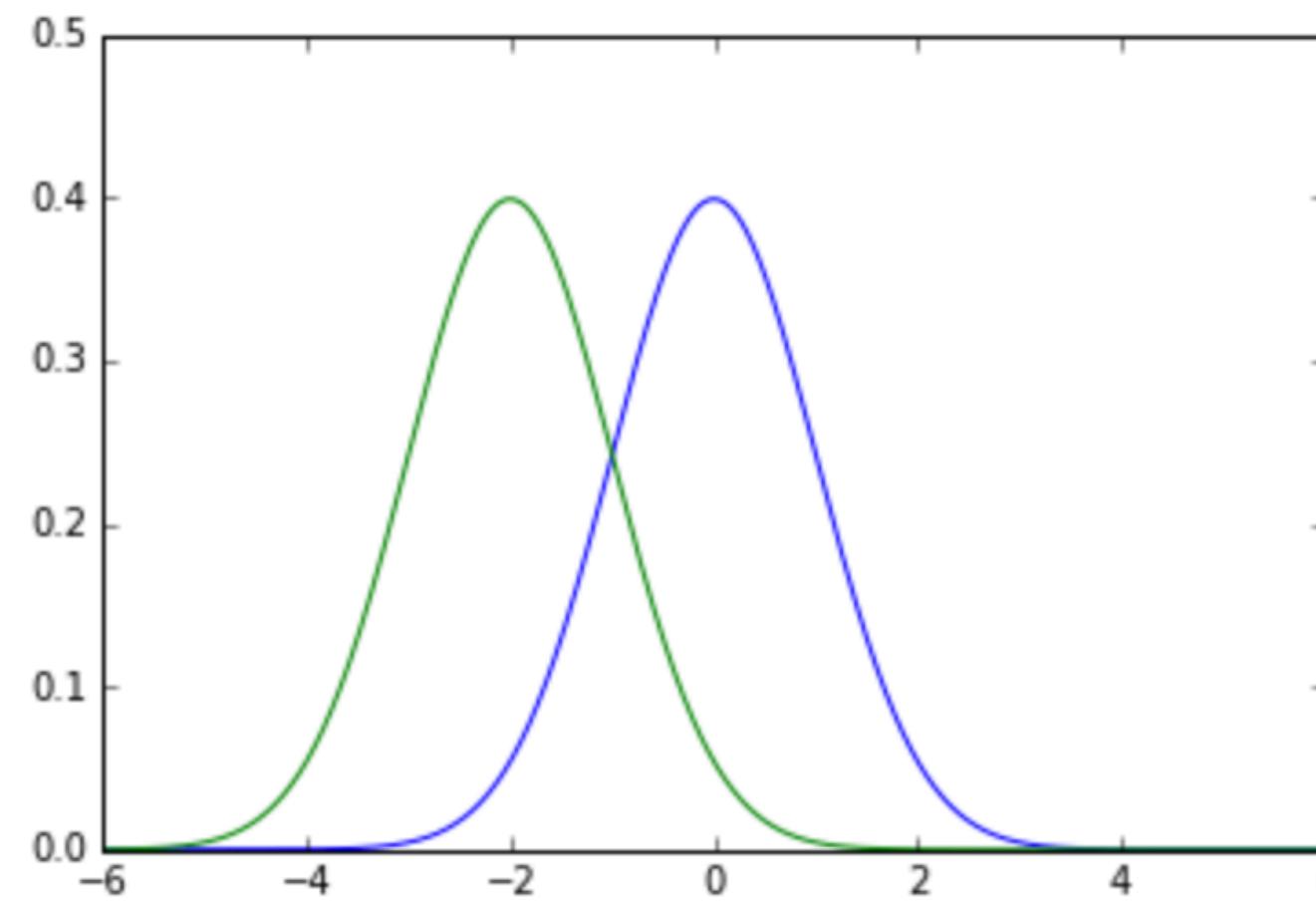
4. Conjugate distributions

Conjugate distributions : example

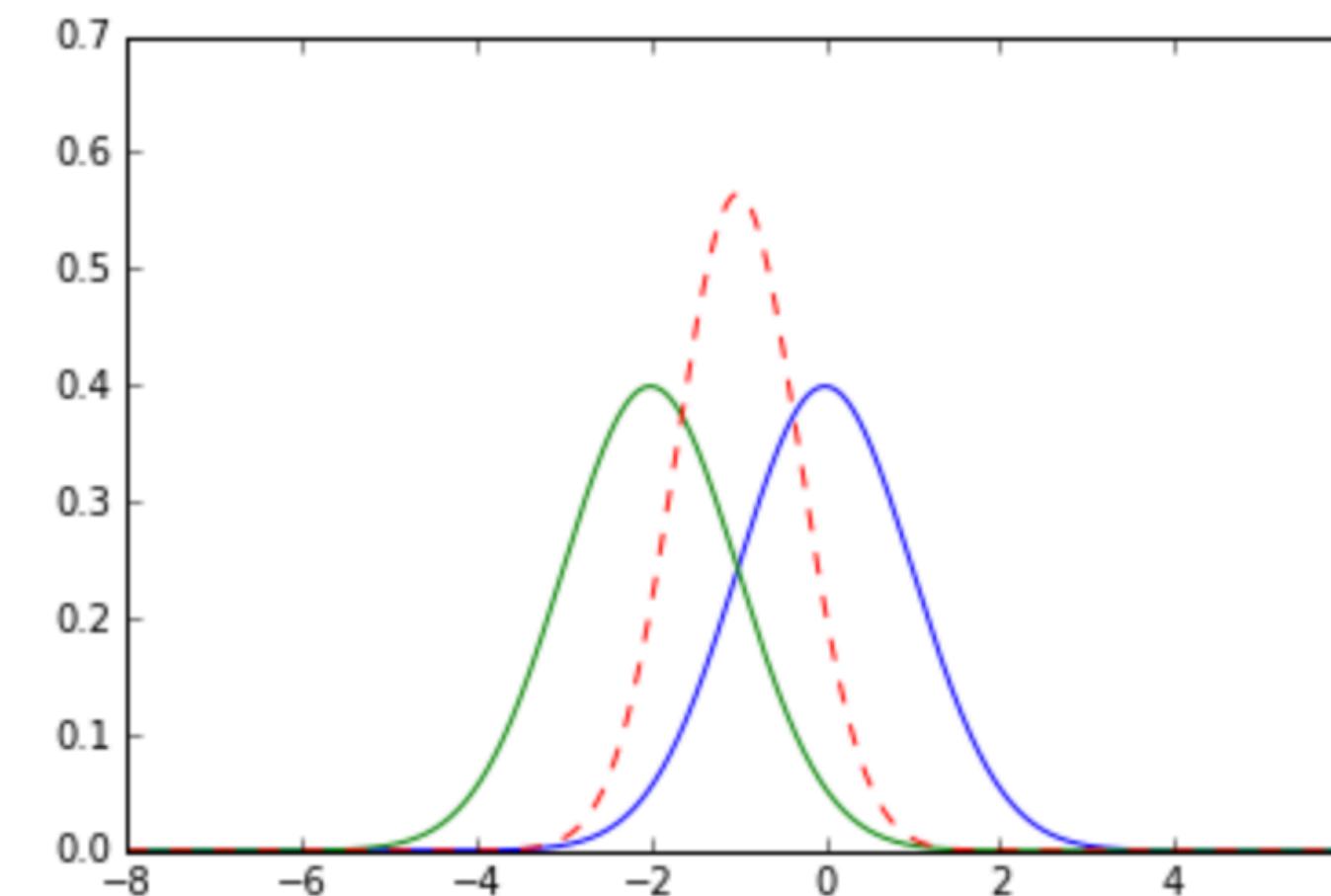
Example

$$P(\theta|X) = \frac{\mathcal{N}(X|\theta, \sigma^2) \times P(\theta)}{P(X)}$$

$\mathcal{N}(\theta|\mu_{prior}, \sigma^2_{prior})$



pointwise product



Exercice (left as an exercice, correction in the next lecture)

Show that $\mathcal{N}(\theta|x/2, 1/2) = \frac{\mathcal{N}(x|\theta, 1) \times \mathcal{N}(\theta|0, 1)}{P(x)}$

4. Conjugate distributions

Usual distributions : Gamma distribution

Gamma distribution

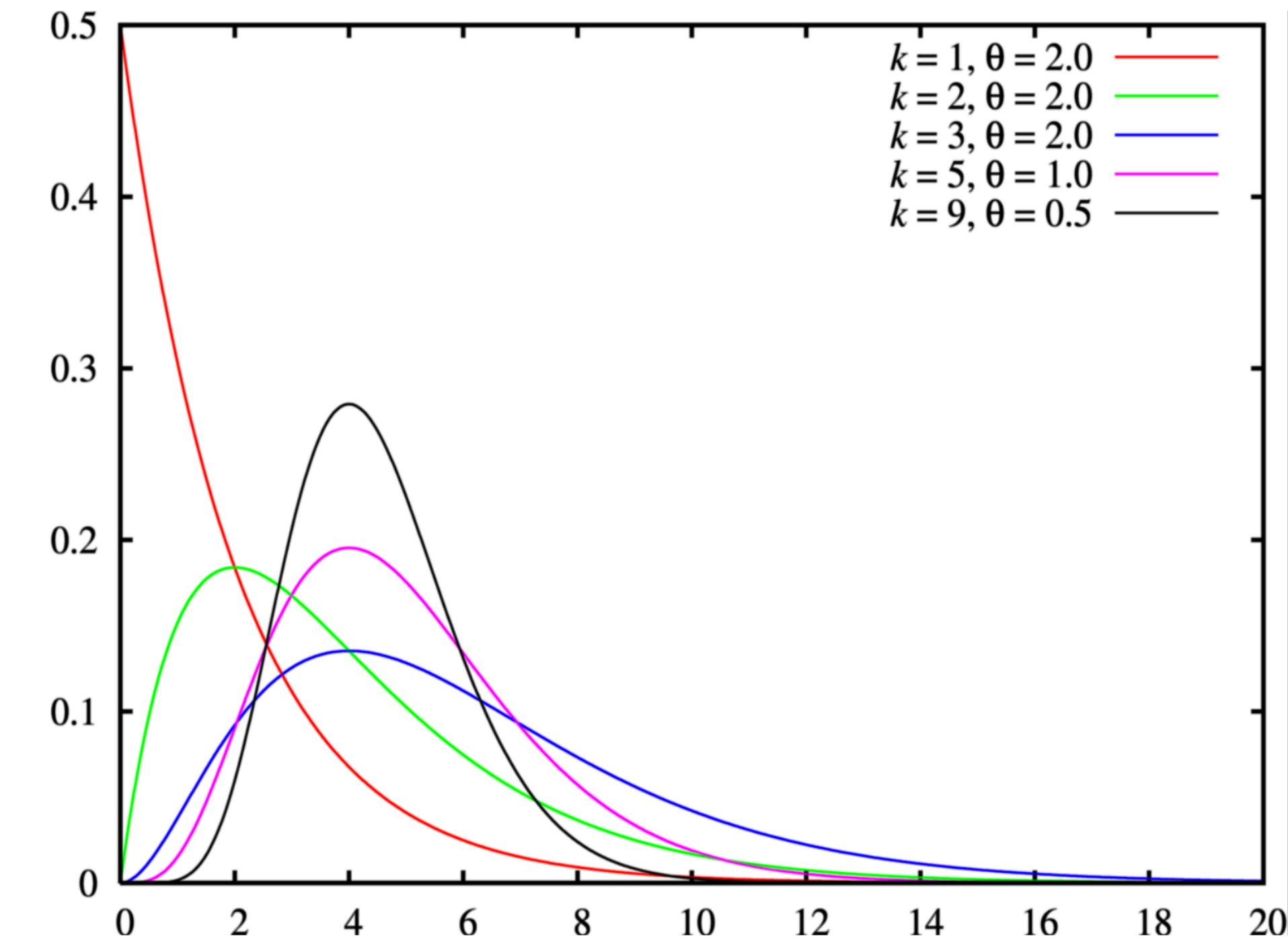
PDF : $\Gamma(x | k, \theta) = \frac{\theta^k}{\Gamma(k)} x^{k-1} e^{-\theta x}$ with $x, k, \theta > 0$

$$\Gamma(k) = (k - 1)!$$

mean : $\mathbb{E}[x] = \frac{k}{\theta}$

variance : $V(x) = \frac{k}{\theta^2}$

mode : $Mode [x] = \frac{k - 1}{\theta}$



Example

$$\Gamma(\gamma | k_{posterior}, \theta_{posterior}) \quad P(\gamma | x) = \frac{\mathcal{N}(x | \mu, \gamma^{-1}) \times P(\gamma)}{P(x)} \quad \Gamma(\gamma | k_{prior}, \theta_{prior})$$

In the context of a gaussian, the prior for the precision is a gamma !

4. Conjugate distributions

Usual distributions : Gamma distribution

Gamma distribution

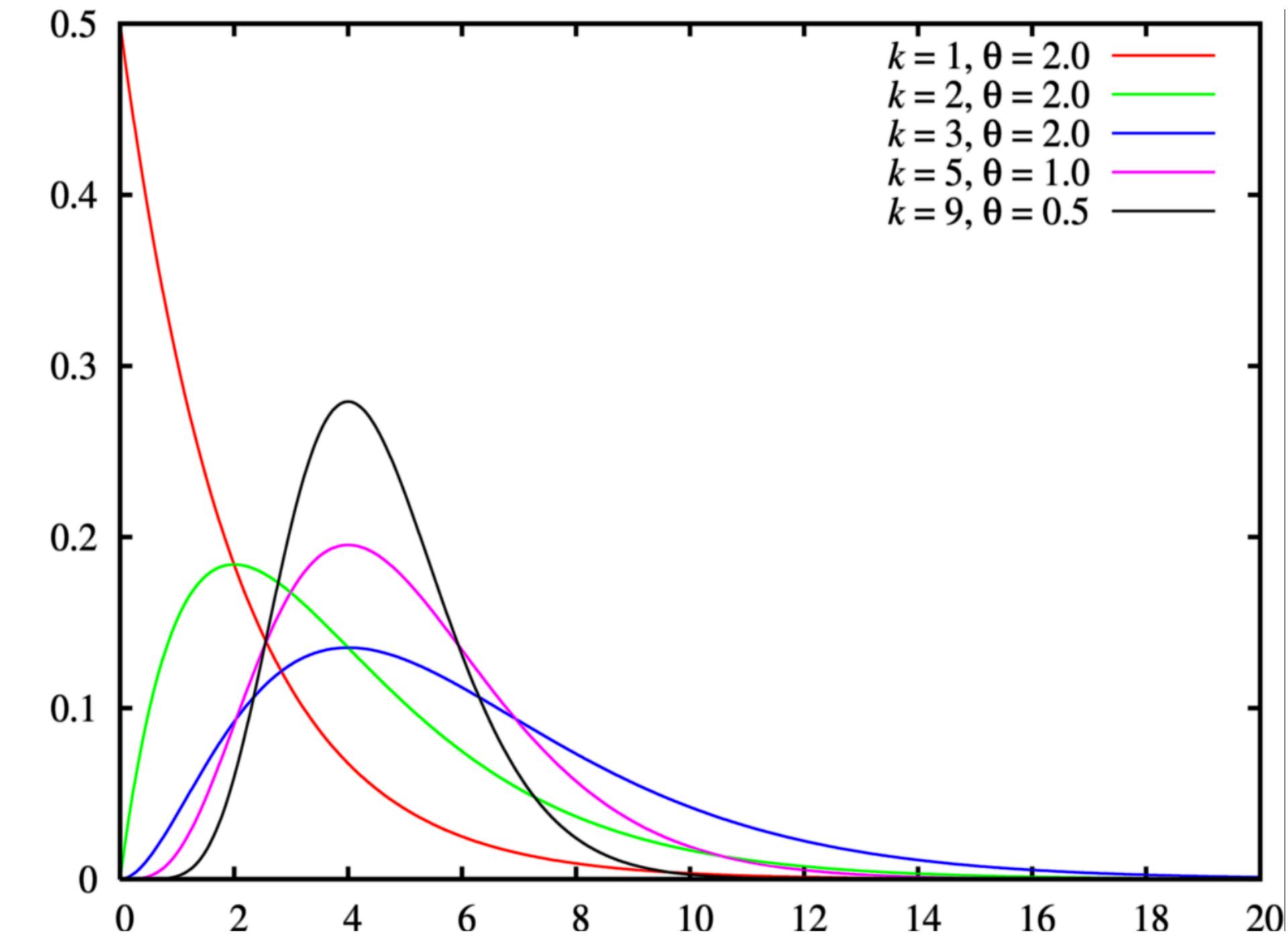
PDF : $\Gamma(x | k, \theta) = \frac{\theta^k}{\Gamma(k)} x^{k-1} e^{-\theta x}$ with $x, k, \theta > 0$

$$\Gamma(k) = (k - 1)!$$

mean : $\mathbb{E}[x] = \frac{k}{\theta}$

variance : $V(x) = \frac{k}{\theta^2}$

mode : $Mode [x] = \frac{k - 1}{\theta}$



Exercice (left as an exercice, correction in the next lecture)

$$P(\gamma | x) = \frac{\mathcal{N}(x | \mu, \gamma^{-1}) \times P(\gamma)}{P(x)}$$

$\Gamma(\gamma | k_{prior}, \theta_{prior})$

$\Gamma(\gamma | k_{posterior}, \theta_{posterior})$

$\Gamma(\gamma | k_{prior} + 1/2, \theta_{prior} + (x - \mu)^2/2)$

In the context of a gaussian, the prior for the precision is a gamma !

4. Conjugate distributions

Usual distributions : Beta distribution

Gamma distribution

PDF : $B(x | \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$ with $\alpha, \beta > 0$ and $x \in [0,1]$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha) \cdot \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

mean : $\mathbb{E}[x] = \frac{\alpha}{\alpha + \beta}$

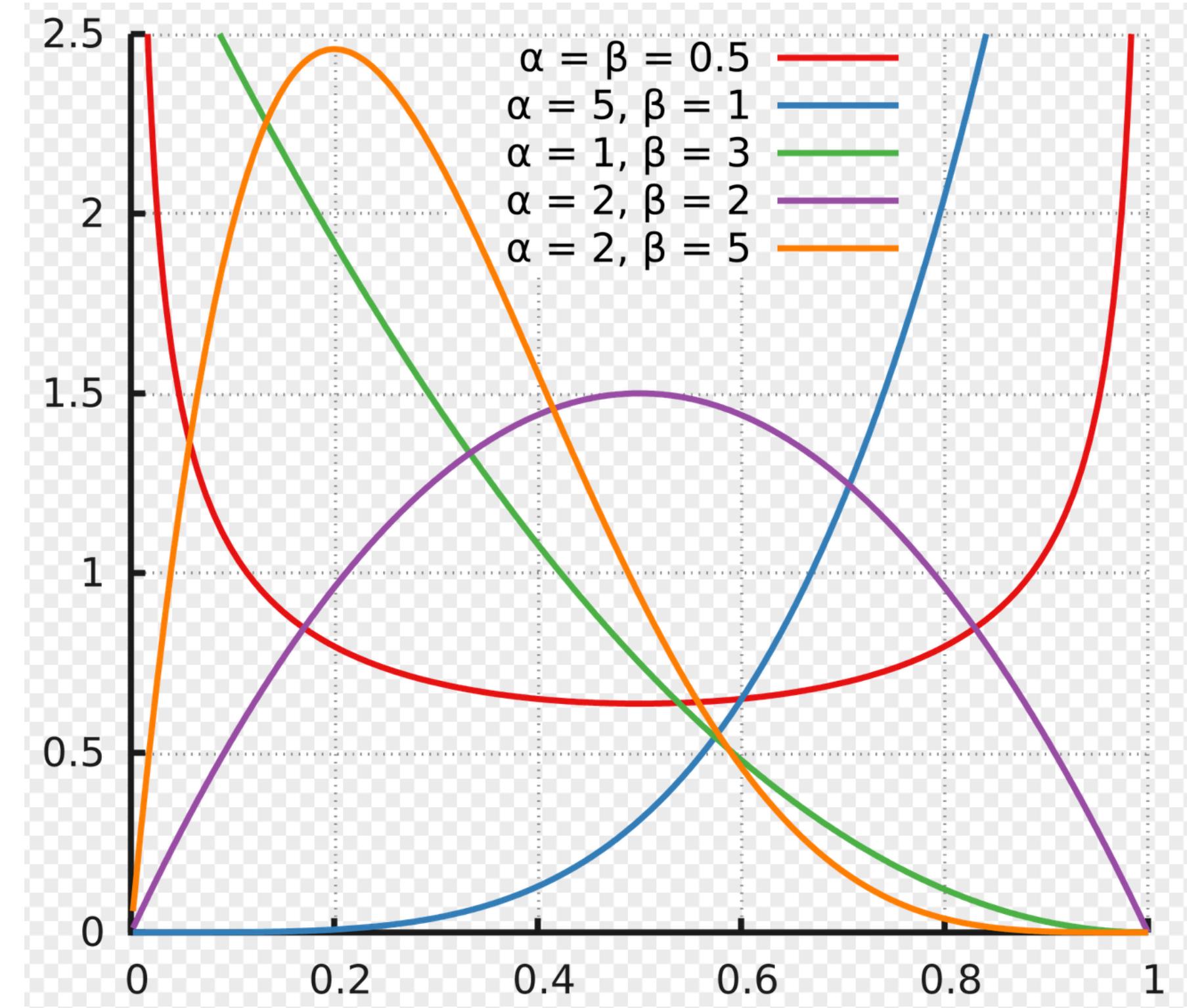
variance : $V(x) = \frac{\alpha\beta}{(\alpha + \beta)^2 \cdot (\alpha + \beta - 1)}$

mode : $Mode[x] = \frac{\alpha - 1}{\alpha + \beta - 2}$

Example

$$B(\theta | \alpha_{posterior}, \beta_{posterior}) — P(\theta | x) = \frac{Ber(x | \theta) \times P(\theta)}{P(x)} \longrightarrow B(\theta | \alpha_{prior}, \beta_{prior})$$

In the context of a Bernoulli distribution, the prior is a beta !



4. Conjugate distributions

Usual distributions : Beta distribution

Gamma distribution

PDF : $B(x | \alpha, \beta) = \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$ with $\alpha, \beta > 0$ and $x \in [0,1]$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha) \cdot \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

mean : $\mathbb{E}[x] = \frac{\alpha}{\alpha + \beta}$

variance : $V(x) = \frac{\alpha\beta}{(\alpha + \beta)^2 \cdot (\alpha + \beta - 1)}$

mode : $Mode[x] = \frac{\alpha - 1}{\alpha + \beta - 2}$

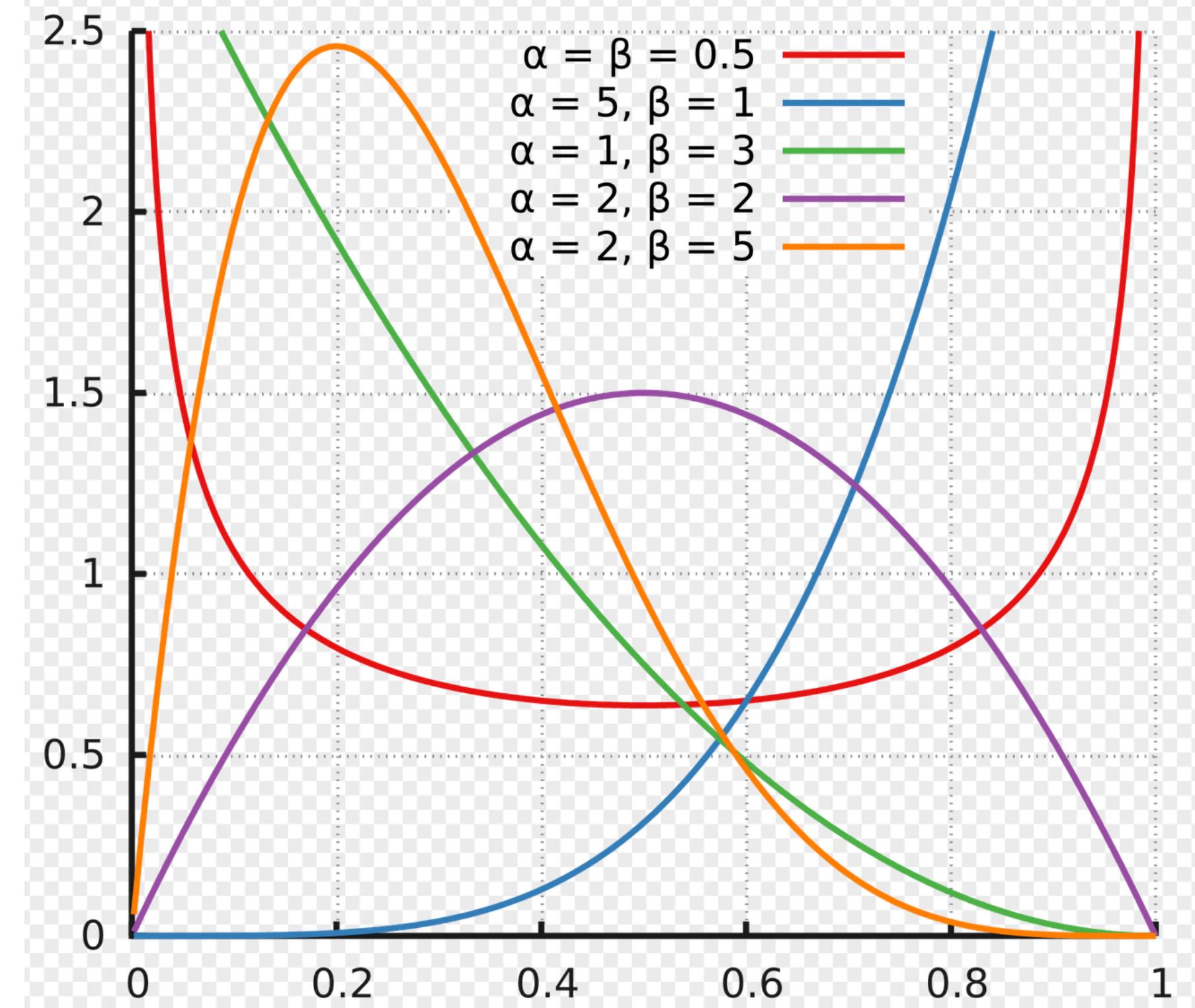
Exercice (left as an exercice, correction in the next lecture)

$$B(\theta | \alpha_{posterior}, \beta_{posterior})$$

$$B(\theta | n_1 + \alpha_{posterior}, n_0 + \beta_{posterior})$$

$$P(\theta | x) = \frac{Ber(x | \theta) \times P(\theta)}{P(x)} \quad \frac{\theta^{n_1} \cdot (1-\theta)^{n_0}}{B(\theta | \alpha_{prior}, \beta_{prior})}$$

In the context of a Bernoulli distribution, the prior is a beta !



4. Conjugate distributions

Limitations

Posterior distribution

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \times P(\theta)}{P(X)}$$

Posterior

Likelihood Prior

Evidence

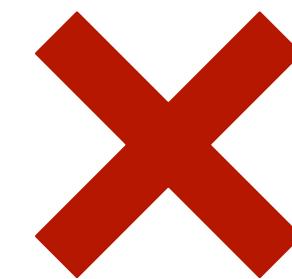
Fixed by model Fixed by us
Fixed by data

Remarks

- We have to **avoid computing** the evidence
- We can choose a **convenient prior** which enable us to compute the posterior :
Conjugate prior



- It computes the **exact posterior**
- Easy for **online learning**



- For some (**complex**) models, the conjugate prior can be **inadequate (improper prior)**
- Can be **unrealistic (non-informative prior)**

!

Road map

Bayesian statistics (02/05/21)



1

Bayesian perspective :

$$P(\theta | X) = \frac{P(X, \theta)}{P(X)} = \frac{P(X | \theta) \cdot P(\theta)}{P(X)}$$

Likelihood Prior distribution
Posterior distribution

θ parameters

X observations

Exemple :
Naive Bayes classifier,
Linear regression,

MAP : $\arg \max_{\theta} P(X | \theta) \cdot P(\theta)$

Conjugate distribution

Pros :
- exact posterior

Cons :
- conjugate prior
maybe inadequate

Latent variable models (17/05/21)

2

Variational Inference (31/05/21)

3

Markov Chain Monte Carlo (07/06/21)

4

Extensions (14/06/21)

5