



Bayesian Learning

Benjamin Guedj, Ph.D.

<https://bguedj.github.io>
Inria Lille - Nord Europe

2017–2018

[benjamin.guedj@inria.fr - Link]

[<https://bguedj.github.io> - Link]

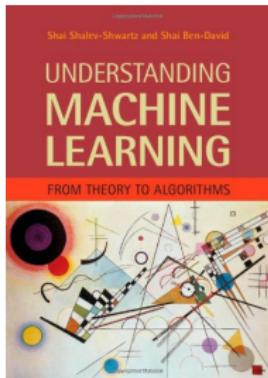
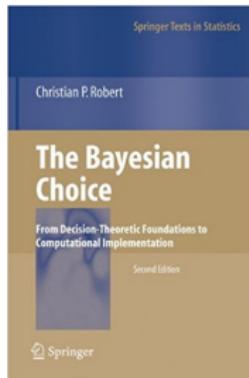
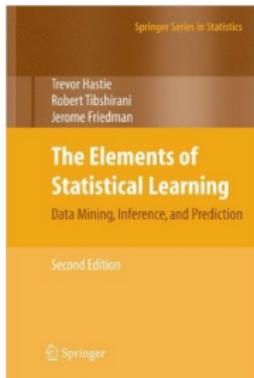
10h de cours (5, 6, 12, 19 et 20 mars 2018)

Projet à rendre. Deadline : **lundi 2 avril 2018 à 23h59.**

[<https://bguedj.github.io/teaching/projet.pdf> - Link]

References

- T. Hastie, R. Tibshirani and J. Friedman. *The Elements of Statistical Learning: Data Mining, Inference and Prediction*, Springer, 2009. [Link]
- C. P. Robert. *The Bayesian Choice*, Springer, 2007. [Link]
- S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning: From Theory to Algorithms*, Cambridge University Press, 2014. [Link]

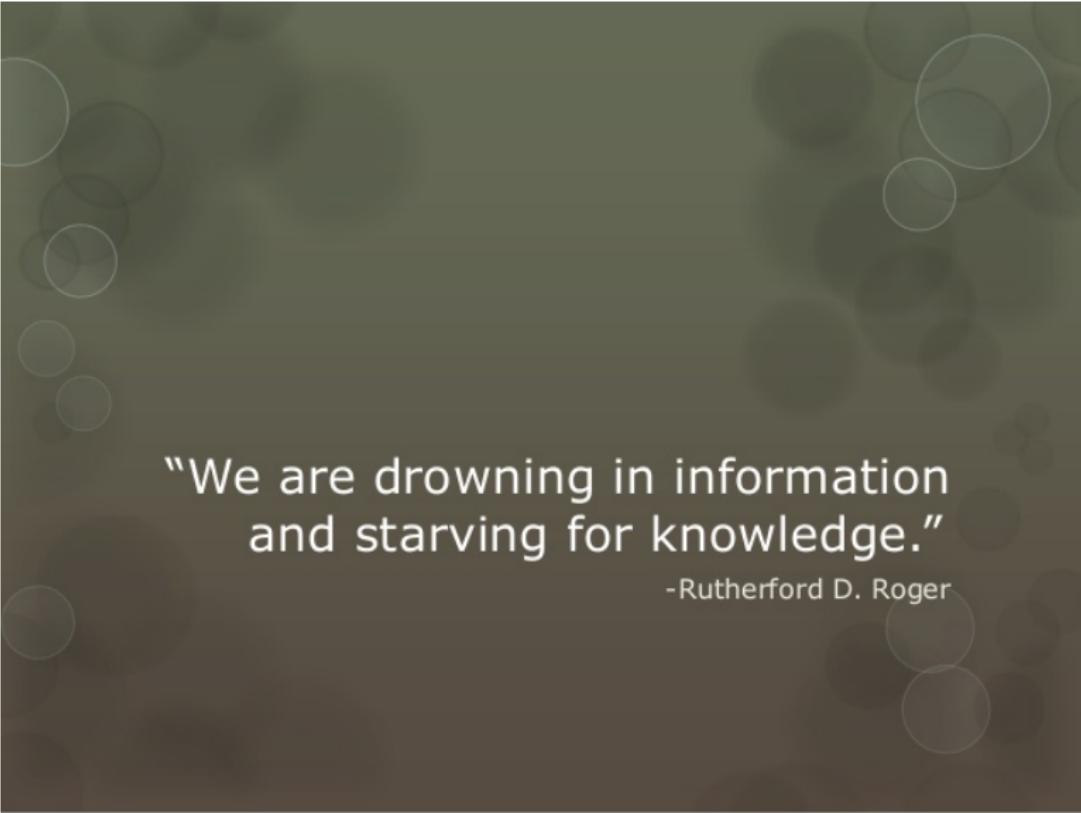


Outline

1. Introduction to statistical / machine learning
2. The Bayesian framework
3. Quasi-Bayesian learning
4. Bayesian learning in practice

The rising of AI

Introduction



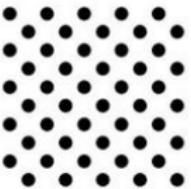
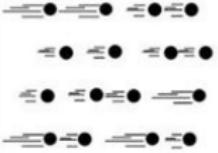
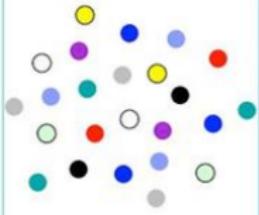
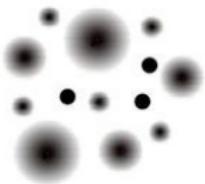
“We are drowning in information
and starving for knowledge.”

-Rutherford D. Roger

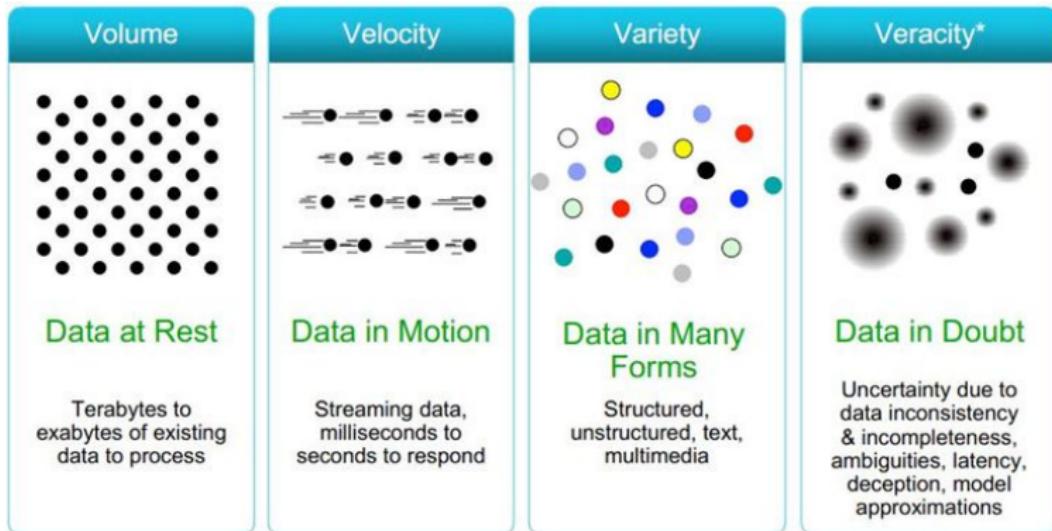
It is vital to remember
that information - in the
sense of raw data - is not
knowledge, that
knowledge is not wisdom,
and that wisdom is not
foresight. But information
is the first essential step
to all of these.

Arthur C Clarke

Big Data 4 V's

Volume	Velocity	Variety	Veracity*
			
Data at Rest Terabytes to exabytes of existing data to process	Data in Motion Streaming data, milliseconds to seconds to respond	Data in Many Forms Structured, unstructured, text, multimedia	Data in Doubt Uncertainty due to data inconsistency & incompleteness, ambiguities, latency, deception, model approximations

Big Data 4 V's



→ Value (\$)

Data Scientists: 100,000 jobs by 2020. Demand is expected to exceed supply by 50 to 60% (McKinsey, 2015)

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day
- ▶ Large Hadron Collider: 40 terabytes/second

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day
- ▶ Large Hadron Collider: 40 terabytes/second
- ▶ Google Street View: 20 petabytes

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day
- ▶ Large Hadron Collider: 40 terabytes/second
- ▶ Google Street View: 20 petabytes
- ▶ AT&T network: 30 petabytes/day

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day
- ▶ Large Hadron Collider: 40 terabytes/second
- ▶ Google Street View: 20 petabytes
- ▶ AT&T network: 30 petabytes/day
- ▶ Human brain may store about 2.5 petabytes of binary data

Volume / Velocity

i	3	6	9	12	15	18	21	24	
10^i	kilo	mega	giga	tera	peta	exa	zeta	yotta	bytes

- ▶ Google: 24 petabytes/day
- ▶ Facebook: 10 terabytes/day
- ▶ Twitter: 7 terabytes/day
- ▶ Large Hadron Collider: 40 terabytes/second
- ▶ Google Street View: 20 petabytes
- ▶ AT&T network: 30 petabytes/day
- ▶ Human brain may store about 2.5 petabytes of binary data
- ▶ ...

Variety / Veracity

The slide illustrates various data types and structures:

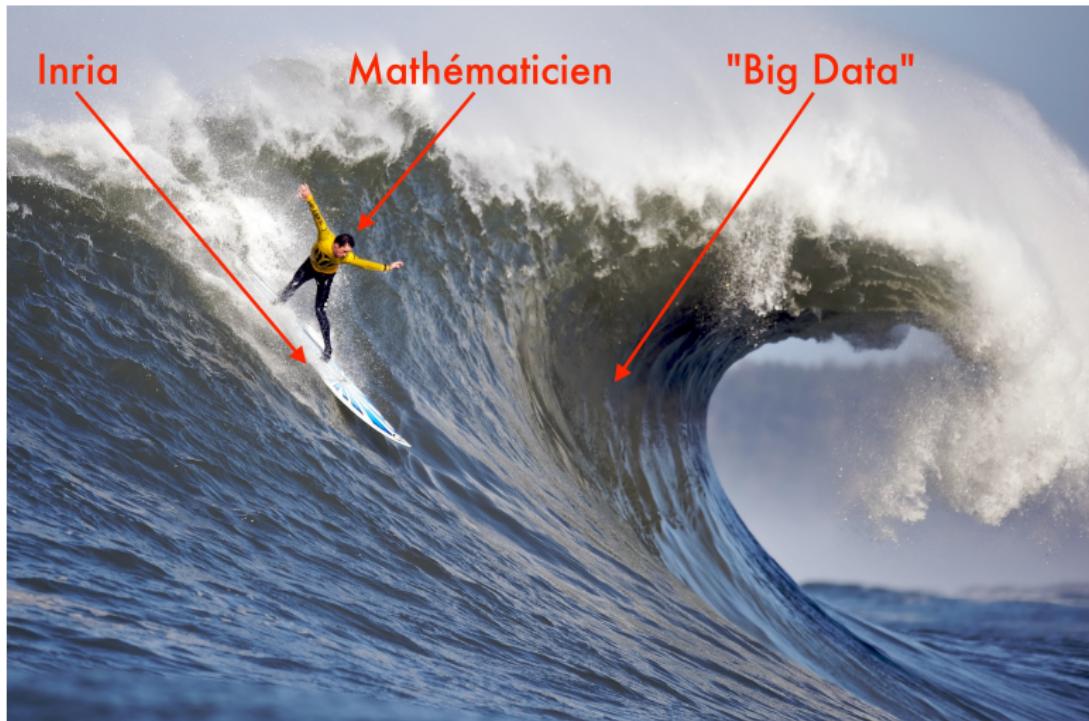
- categorical**: Marital status **married** (Homer Simpson)
- rank**: Drink preference **beer > soda > water**
- integer**: Children **3**
- ordinal**: Intelligence **low**
- missing**: Size (m) **?**
- continuous**: Weight (kg) **119.5**

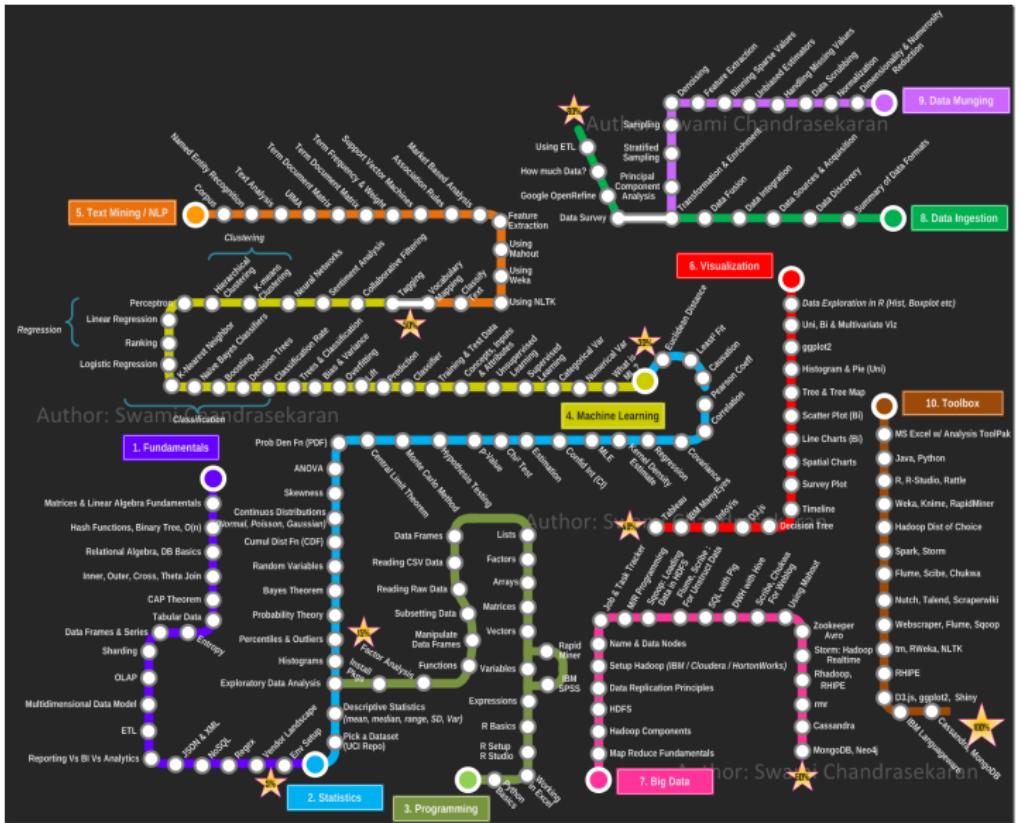
functional: Drink consumption graph showing consumption over time (8h, 12h, 16h, 20h, 24h).

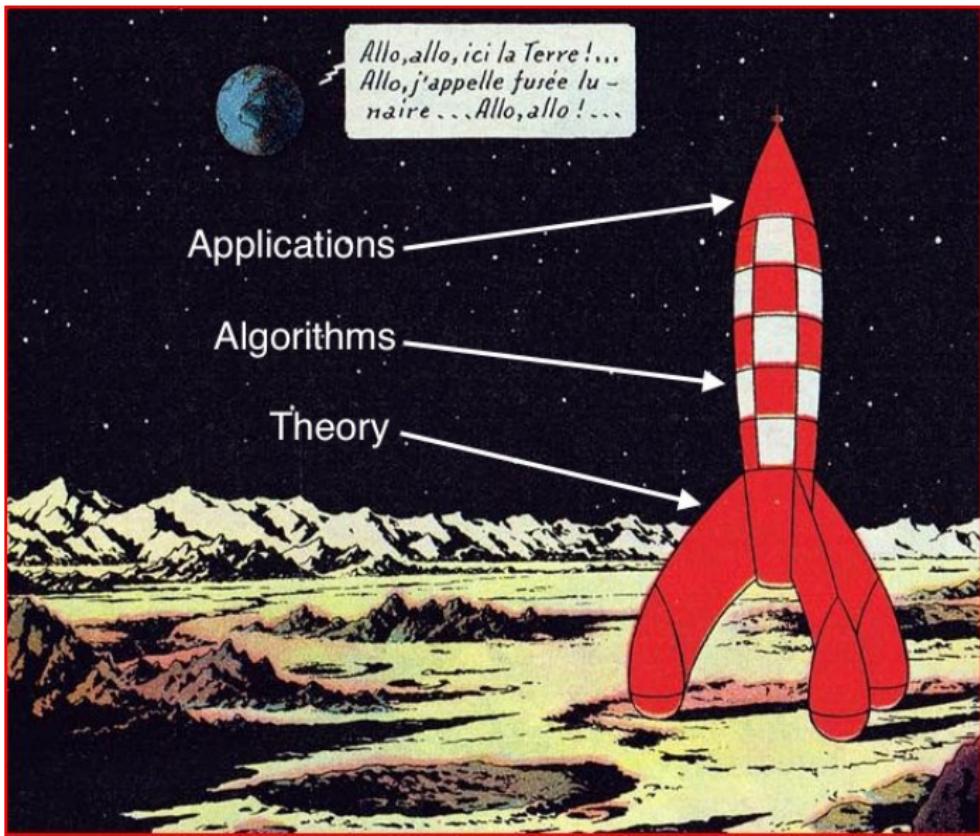
graph: Family graph showing relationships between Homer, Marge, Bart, Lisa, and Maggie.

And so on...

My job (allegory)







A foretaste of Learning Theory

{Statistical,Machine} Learning: building automatic procedures to infer general rules from examples.

{Statistical,Machine} Learning: building automatic procedures to infer general rules from examples.

In the (rather not so?) long term: mimic the inductive functioning of the humain brain to develop an artificial intelligence.

{Statistical,Machine} Learning: building automatic procedures to infer general rules from examples.

In the (rather not so?) long term: mimic the inductive functioning of the humain brain to develop an artificial intelligence.

In the Big Data Era, very dynamic field at the crossroads of Computer Science, Optimization and Statistics.

Probabilistic framework: n -sample $\mathcal{D}_n = (\mathbf{X}_i, \mathbf{Y}_i)_{i=1}^n$ of i.i.d. replications of some random variable

$$(\mathbf{X}, \mathbf{Y}) \in \mathcal{X} \times \mathcal{Y}, \quad \dim(\mathcal{X}) = d.$$

Probabilistic framework: n -sample $\mathcal{D}_n = (\mathbf{X}_i, \mathbf{Y}_i)_{i=1}^n$ of i.i.d. replications of some random variable

$$(\mathbf{X}, \mathbf{Y}) \in \mathcal{X} \times \mathcal{Y}, \quad \dim(\mathcal{X}) = d.$$

We want to infer the link between the explanatory variable \mathbf{X} and the response variable \mathbf{Y} , *i.e.*, use \mathcal{D}_n to build up $\hat{\phi}$ such that $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

Probabilistic framework: n -sample $\mathcal{D}_n = (\mathbf{X}_i, \mathbf{Y}_i)_{i=1}^n$ of i.i.d. replications of some random variable

$$(\mathbf{X}, \mathbf{Y}) \in \mathcal{X} \times \mathcal{Y}, \quad \dim(\mathcal{X}) = d.$$

We want to infer the link between the explanatory variable \mathbf{X} and the response variable \mathbf{Y} , *i.e.*, use \mathcal{D}_n to build up $\hat{\phi}$ such that $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

- ▶ Classification: \mathcal{Y} is discrete.
- ▶ Regression: \mathcal{Y} is a continuum.

- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

Big Data Era: easy/cheap to collect massive amounts of data,
hence typically $\mathcal{X} = \mathbb{R}^d$ where d and n may be (extremely) large.

- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

Big Data Era: easy/cheap to collect massive amounts of data,
hence typically $\mathcal{X} = \mathbb{R}^d$ where d and n may be (extremely) large.

- ▶ Good ol' statistics: n and d small

- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

Big Data Era: easy/cheap to collect massive amounts of data,
hence typically $\mathcal{X} = \mathbb{R}^d$ where d and n may be (extremely) large.

- ▶ Good ol' statistics: n and d small
- ▶ Tall data: $n \gg d$

- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

Big Data Era: easy/cheap to collect massive amounts of data,
hence typically $\mathcal{X} = \mathbb{R}^d$ where d and n may be (extremely) large.

- ▶ Good ol' statistics: n and d small
- ▶ Fat data: $d \gg n$
- ▶ Tall data: $n \gg d$

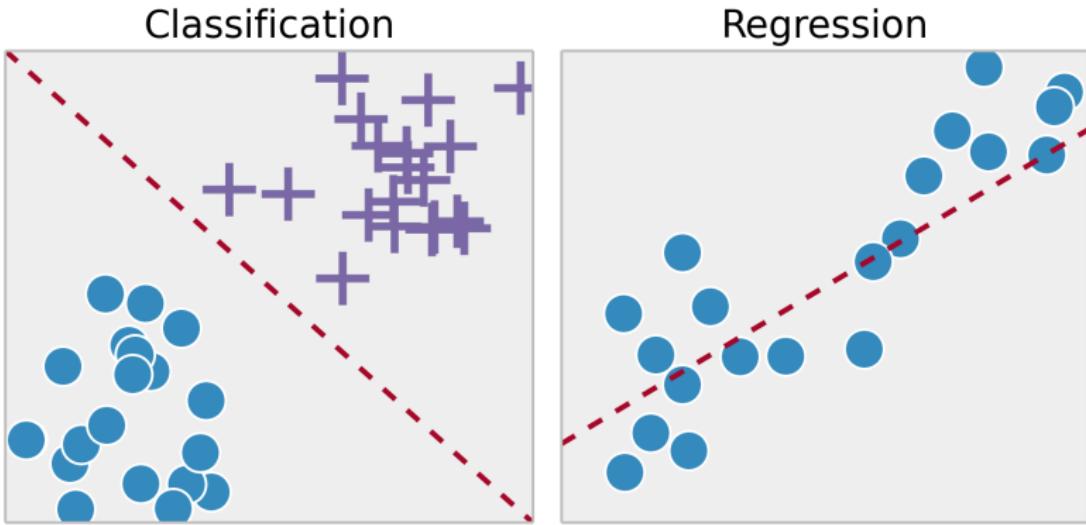
- ▶ Online learning: observations are revealed over time.
- ▶ Batch learning: all observations are revealed at once.

Big Data Era: easy/cheap to collect massive amounts of data,
hence typically $\mathcal{X} = \mathbb{R}^d$ where d and n may be (extremely) large.

- ▶ Good ol' statistics: n and d small
- ▶ Tall data: $n \gg d$
- ▶ Fat data: $d \gg n$
- ▶ Big/massive data: n and d huge

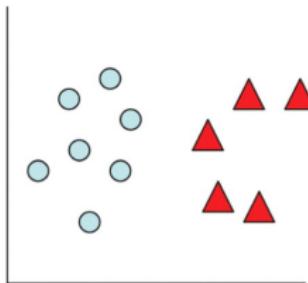
- ▶ Supervised learning: all of the \mathbf{Y}_i s are observed.

- ▶ Supervised learning: all of the \mathbf{Y}_i s are observed.

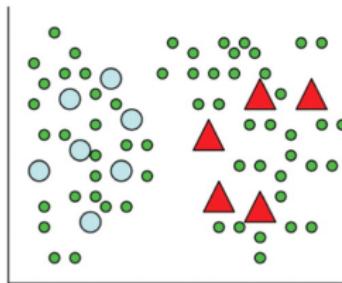


- ▶ Semi-supervised learning: some of the \mathbf{Y}_i s are observed (labeling is expensive or difficult).

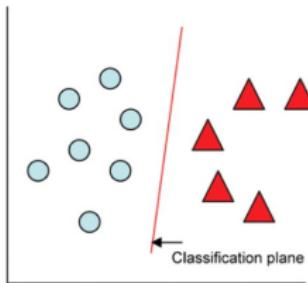
- ▶ Semi-supervised learning: some of the \mathbf{Y}_i s are observed (labeling is expensive or difficult).



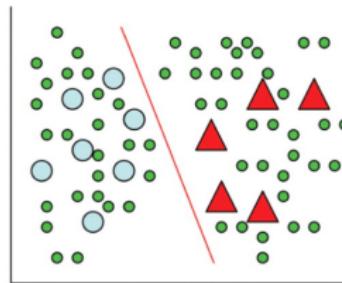
Labeled Data
(a)



Labeled and Unlabeled Data
(b)



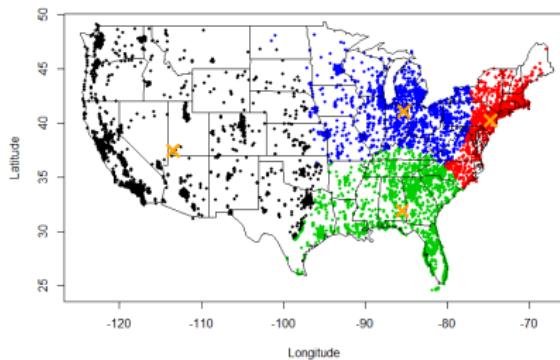
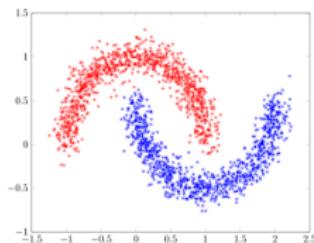
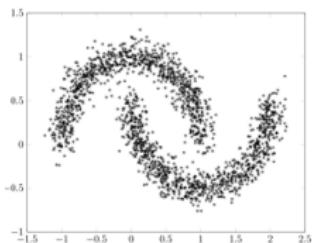
Supervised Learning
(c)



Semi-Supervised Learning
(d)

- ▶ Unsupervised learning: none of the \mathbf{Y}_i 's are observed (detect patterns).

- Unsupervised learning: none of the \mathbf{Y}_i 's are observed (detect patterns).



- ▶ Unsupervised learning: none of the \mathbf{Y}_i 's are observed (detect patterns).

- ▶ Unsupervised learning: none of the \mathbf{Y}_i 's are observed (detect patterns).



- ▶ Reinforcement learning: feedback (possibly adversarial) from the environment (robotics, adversarial environments, training...).

- ▶ Reinforcement learning: feedback (possibly adversarial) from the environment (robotics, adversarial environments, training...).



Loss function:

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

Loss function:

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

$\ell(\hat{\phi}(\mathbf{X}), Y)$ (random) quantifies how a predictor $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

Loss function:

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

$\ell(\hat{\phi}(\mathbf{X}), Y)$ (random) quantifies how a predictor $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

A predictor is any mapping

$$\hat{\phi}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{Y}^{\mathcal{X}}.$$

Loss function:

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

$\ell(\hat{\phi}(\mathbf{X}), Y)$ (random) quantifies how a predictor $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

A predictor is any mapping

$$\hat{\phi}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{Y}^{\mathcal{X}}.$$

Risk

$$R(\hat{\phi}) = \mathbb{E} \left[\ell \left(\hat{\phi}(\mathbf{X}), Y \right) \right],$$

Loss function:

$$\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+.$$

$\ell(\hat{\phi}(\mathbf{X}), Y)$ (random) quantifies how a predictor $\hat{\phi}(\mathbf{X})$ is a "good" approximation of \mathbf{Y} .

A predictor is any mapping

$$\hat{\phi}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{Y}^{\mathcal{X}}.$$

Risk

$$R(\hat{\phi}) = \mathbb{E} \left[\ell \left(\hat{\phi}(\mathbf{X}), Y \right) \right],$$

Empirical risk

$$R_n(\hat{\phi}) = \frac{1}{n} \sum_{i=1}^n \left[\ell \left(\hat{\phi}(\mathbf{X}_i), Y_i \right) \right].$$

- Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.

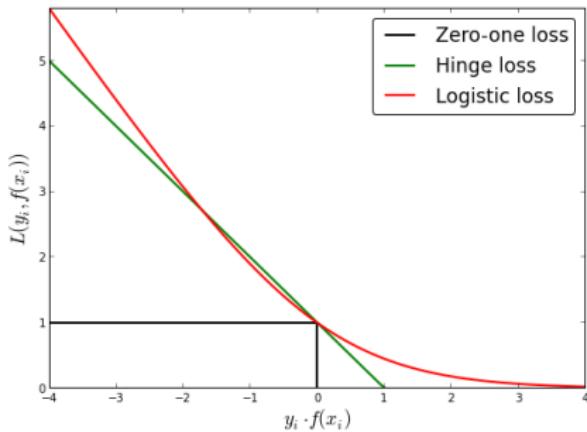
- ▶ Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.
- ▶ Absolute loss (regression): $\ell(a, b) = |a - b|$.

- ▶ Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.
- ▶ Absolute loss (regression): $\ell(a, b) = |a - b|$.
- ▶ 0-1 loss (classification): $\ell(a, b) = \mathbb{1}_{\{a \neq b\}}$.

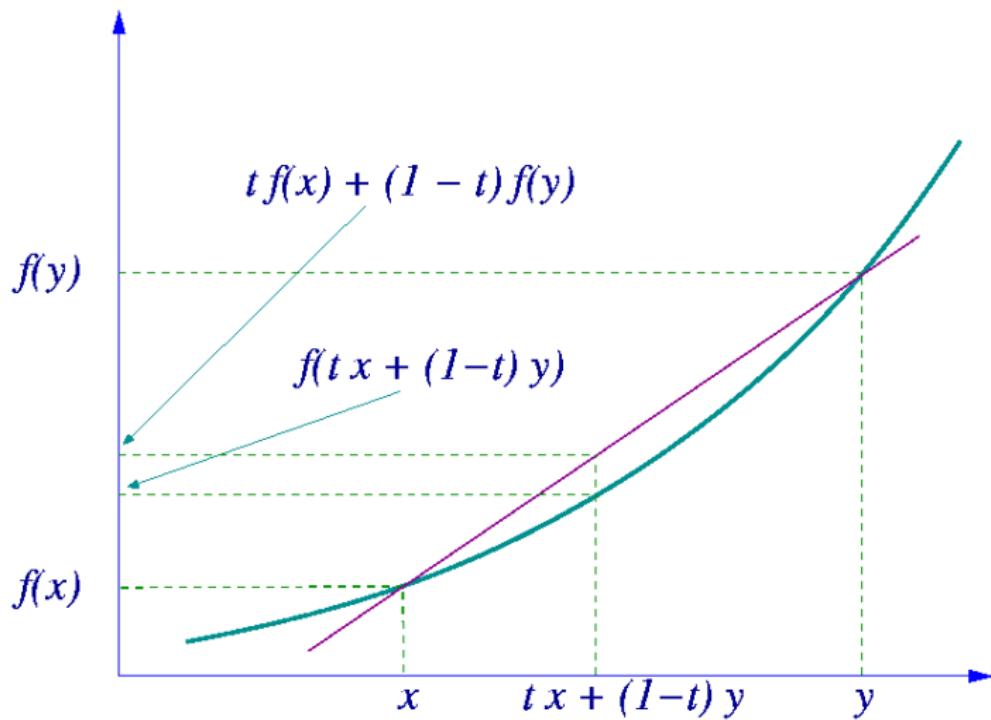
- ▶ Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.
- ▶ Absolute loss (regression): $\ell(a, b) = |a - b|$.
- ▶ 0-1 loss (classification): $\ell(a, b) = \mathbb{1}_{\{a \neq b\}}$.
- ▶ Hinge loss (classification): $\ell(a, b) = \max(0, 1 - ab)$.

- ▶ Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.
- ▶ Absolute loss (regression): $\ell(a, b) = |a - b|$.
- ▶ 0-1 loss (classification): $\ell(a, b) = \mathbb{1}_{\{a \neq b\}}$.
- ▶ Hinge loss (classification): $\ell(a, b) = \max(0, 1 - ab)$.
- ▶ Logistic loss (classification): $\ell(a, b) = \log[1 + \exp(-ab)]$.

- ▶ Quadratic loss (regression): $\ell(a, b) = (a - b)^2$.
- ▶ Absolute loss (regression): $\ell(a, b) = |a - b|$.
- ▶ 0-1 loss (classification): $\ell(a, b) = \mathbb{1}_{\{a \neq b\}}$.
- ▶ Hinge loss (classification): $\ell(a, b) = \max(0, 1 - ab)$.
- ▶ Logistic loss (classification): $\ell(a, b) = \log[1 + \exp(-ab)]$.



Convexity is (often) crucial



Statistical Learning vs. Machine Learning

Same task, different approaches:

Statistical Learning vs. Machine Learning

Same task, different approaches:

- ▶ In machine learning, given some deterministic sequence (\mathbf{x}_i, y_i) , solve

$$\hat{\phi}(\cdot) = \arg \min_m \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) \right\}.$$

Statistical Learning vs. Machine Learning

Same task, different approaches:

- ▶ In machine learning, given some deterministic sequence (\mathbf{x}_i, y_i) , solve

$$\hat{\phi}(\cdot) = \arg \min_m \left\{ \sum_{i=1}^n \ell(y_i, m(\mathbf{x}_i)) \right\}.$$

- ▶ In statistical learning, assume that the Y_i s are realisations of some random variable Y (given \mathbf{X}) with distribution P . Solve

$$\hat{\phi}(\cdot) = \arg \max_m \left\{ \sum_{i=1}^n \log dP(Y_i; m(\mathbf{X}_i)) \right\}.$$

SL vs. ML in the simple parametric case

SL vs. ML in the simple parametric case

- ▶ In machine learning, given some deterministic sequence (\mathbf{x}_i, y_i) , solve

$$\hat{\phi}(\cdot) = \arg \min_{\theta} \left\{ \sum_{i=1}^n \ell(y_i, \langle \theta, \mathbf{x}_i \rangle) \right\}.$$

SL vs. ML in the simple parametric case

- ▶ In machine learning, given some deterministic sequence (\mathbf{x}_i, y_i) , solve

$$\hat{\phi}(\cdot) = \arg \min_{\theta} \left\{ \sum_{i=1}^n \ell(y_i, \langle \theta, \mathbf{x}_i \rangle) \right\}.$$

- ▶ In statistical learning, assume that the Y_i s are realisations of some random variable Y (given \mathbf{X}) with distribution P . Solve

$$\hat{\phi}(\cdot) = \arg \max_{\theta} \left\{ \sum_{i=1}^n \log dP(Y_i | \mathbf{x}_i, \theta) \right\}.$$

*All models are wrong
but some are useful*



George E.P. Box



If the only tool you have is a hammer, you tend to see every problem as a nail.

(Abraham Maslow)



A primer on probability distributions

All words are hyperlinks.

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]
- ▶ [Binomial - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]
- ▶ [Binomial - Link]
- ▶ [Bernoulli - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]
- ▶ [Binomial - Link]
- ▶ [Bernoulli - Link]
- ▶ [Gamma - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]
- ▶ [Binomial - Link]
- ▶ [Bernoulli - Link]
- ▶ [Gamma - Link]
- ▶ [Student - Link]

A primer on probability distributions

All words are hyperlinks.

- ▶ [Normal - Link]
- ▶ [Inverse Gaussian (a.k.a Wald) - Link]
- ▶ [Beta - Link]
- ▶ [Poisson - Link]
- ▶ [Binomial - Link]
- ▶ [Bernoulli - Link]
- ▶ [Gamma - Link]
- ▶ [Student - Link]
- ▶ ...

The Bayesian paradigm

Introductory example

Consider observations $\mathbf{x} = (x_1, \dots, x_n)$ generated from a probability distribution with density $f(\cdot|\theta)$.

Introductory example

Consider observations $\mathbf{x} = (x_1, \dots, x_n)$ generated from a probability distribution with density $f(\cdot|\theta)$.

The associated likelihood is the inverted density:

$$\mathcal{L}(\theta|\mathbf{x}) = f(\mathbf{x}|\theta).$$

Example $f(\cdot|\theta) = \mathcal{N}(\theta, 1)$.

Bayes' Theorem

Bayes' Theorem

Inversion of probabilities a.k.a actualisation principle.

Bayes' Theorem

Inversion of probabilities a.k.a actualisation principle.

If A and B are events such that $\mathbb{P}(B) \neq 0$,

$$\begin{aligned}\mathbb{P}(A|B) &= \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B|A)\mathbb{P}(A) + \mathbb{P}(B|A^c)\mathbb{P}(A^c)} \\ &= \frac{\mathbb{P}(B|A)\mathbb{P}(A)}{\mathbb{P}(B)}.\end{aligned}$$

(due to Thomas Bayes, published in 1764)

Who was Thomas Bayes?

Who was Thomas Bayes?



Reverend Thomas Bayes (ca. 1702–1761) Presbyterian minister in Kent from 1731. Election to the Royal Society based on a tract of 1736 where he defended the views and philosophy of Newton. Sole probability paper, "Essay Towards Solving a Problem in the Doctrine of Chances", published posthumously in 1763 and containing the seeds of Bayes' Theorem.

A new paradigm

Bayes introduces a whole new perspective.

A new paradigm

Bayes introduces a whole new perspective.

- ▶ Uncertainty on the parameter θ , modeled through a probability distribution π , called *prior distribution*.

A new paradigm

Bayes introduces a whole new perspective.

- ▶ Uncertainty on the parameter θ , modeled through a probability distribution π , called *prior distribution*.
- ▶ Inference based on the distribution of θ conditional on \mathbf{X} $\pi(\theta|\mathbf{x})$, called *posterior distribution*

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{\int f(\mathbf{x}|\theta)\pi(\theta)d\theta}.$$

A Bayesian model

. . . is made of a parametric (in this course) statistical model defined through its likelihood $f(\mathbf{x}|\theta)$ and a prior distribution on the parameter $\pi(\theta)$.

Consequences

- ▶ Semantic drift from unknown to random

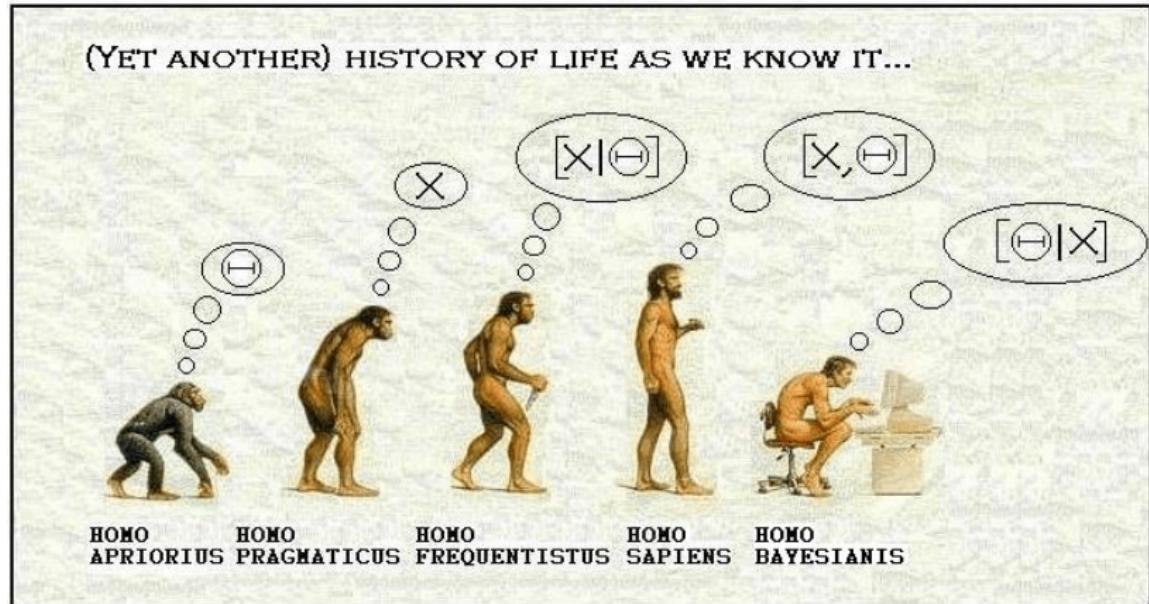
Consequences

- ▶ Semantic drift from unknown to random
- ▶ Actualization of θ by extracting the information contained in the observation x

Consequences

- ▶ Semantic drift from unknown to random
- ▶ Actualization of θ by extracting the information contained in the observation x
- ▶ Allows incorporation of imperfect information in the decision process

The advantages of being a Bayesian



Distributions (1/2)

Given the likelihood $f(\mathbf{x}|\theta)$ and the prior $\pi(\theta)$, several distributions of interest:

Distributions (1/2)

Given the likelihood $f(\mathbf{x}|\theta)$ and the prior $\pi(\theta)$, several distributions of interest:

- ▶ The *joint distribution* of (θ, \mathbf{x})

$$\varphi(\theta, \mathbf{x}) = f(\mathbf{x}|\theta)\pi(\theta).$$

Distributions (1/2)

Given the likelihood $f(\mathbf{x}|\theta)$ and the prior $\pi(\theta)$, several distributions of interest:

- ▶ The *joint distribution* of (θ, \mathbf{x})

$$\varphi(\theta, \mathbf{x}) = f(\mathbf{x}|\theta)\pi(\theta).$$

- ▶ The *marginal distribution* of \mathbf{x}

$$m(\mathbf{x}) = \int \varphi(\theta, \mathbf{x})d\theta = \int f(\mathbf{x}|\theta)\pi(\theta)d\theta.$$

Distributions (2/2)

Distributions (2/2)

- ▶ The *posterior distribution* of θ

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{\int f(\mathbf{x}|\theta)\pi(\theta)d\theta} = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})}.$$

Distributions (2/2)

- ▶ The *posterior distribution* of θ

$$\pi(\theta|\mathbf{x}) = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{\int f(\mathbf{x}|\theta)\pi(\theta)d\theta} = \frac{f(\mathbf{x}|\theta)\pi(\theta)}{m(\mathbf{x})}.$$

- ▶ The *predictive distribution* of y when $y \sim g(\cdot|\theta, \mathbf{x})$

$$g(y|\mathbf{x}) = \int g(y|\theta, \mathbf{x})\pi(\theta|\mathbf{x})d\theta.$$

A comprehensive example normal-normal

Assume that we model $\mathbf{x} \sim \mathcal{N}(\theta, 1)$ and use the prior $\theta \sim \mathcal{N}(a, 10)$.

A comprehensive example normal-normal

Assume that we model $\mathbf{x} \sim \mathcal{N}(\theta, 1)$ and use the prior $\theta \sim \mathcal{N}(a, 10)$.

$$\begin{aligned}\pi(\theta|\mathbf{x}) &\propto f(\mathbf{x}|\theta)\pi(\theta) \\ &\propto \exp\left(-\frac{(\mathbf{x}-\theta)^2}{2} - \frac{(\theta-a)^2}{20}\right) \\ &\propto \exp\left(-\frac{11}{20}\theta^2 + \theta(\mathbf{x} + a/10)\right) \\ &\propto \exp\left(-\frac{11}{20}\left(\theta - \frac{10\mathbf{x} + a}{11}\right)^2\right)\end{aligned}$$

A comprehensive example normal-normal

Assume that we model $\mathbf{x} \sim \mathcal{N}(\theta, 1)$ and use the prior $\theta \sim \mathcal{N}(a, 10)$.

$$\begin{aligned}\pi(\theta|\mathbf{x}) &\propto f(\mathbf{x}|\theta)\pi(\theta) \\ &\propto \exp\left(-\frac{(\mathbf{x}-\theta)^2}{2} - \frac{(\theta-a)^2}{20}\right) \\ &\propto \exp\left(-\frac{11}{20}\theta^2 + \theta(\mathbf{x} + a/10)\right) \\ &\propto \exp\left(-\frac{11}{20}\left(\theta - \frac{10\mathbf{x}+a}{11}\right)^2\right)\end{aligned}$$

which means $\theta|\mathbf{x} \sim \mathcal{N}\left(\frac{10\mathbf{x}+a}{11}, \frac{10}{11}\right)$.

A comprehensive example uniform-binomial

Bayes' very own example: a billiard ball W rolls on a line of length one, with a uniform probability of stopping anywhere: W stops at p .

A comprehensive example uniform-binomial

Bayes' very own example: a billiard ball W rolls on a line of length one, with a uniform probability of stopping anywhere: W stops at p .

A second ball O then rolls n times under the same assumptions. X denotes the number of times the ball O stopped on the left of W .

A comprehensive example uniform-binomial

Bayes' very own example: a billiard ball W rolls on a line of length one, with a uniform probability of stopping anywhere: W stops at p .

A second ball O then rolls n times under the same assumptions. X denotes the number of times the ball O stopped on the left of W .

Bayes' question: given X , what inference can we make on p ?

Mathematical translation

Derive the posterior distribution of p given X , when $p \sim \mathcal{U}(0, 1)$ and $X \sim \mathcal{B}(n, p)$.

Resolution 1/2

Since

$$\mathbb{P}(X = x | p) = \binom{n}{x} p^x (1 - p)^{n-x},$$

Resolution 1/2

Since

$$\mathbb{P}(X = x | p) = \binom{n}{x} p^x (1 - p)^{n-x},$$

$$\mathbb{P}(a < p < b \text{ and } X = x) = \int_a^b \binom{n}{x} p^x (1 - p)^{n-x} dp,$$

Resolution 1/2

Since

$$\mathbb{P}(X = x | p) = \binom{n}{x} p^x (1 - p)^{n-x},$$

$$\mathbb{P}(a < p < b \text{ and } X = x) = \int_a^b \binom{n}{x} p^x (1 - p)^{n-x} dp,$$

and

$$\mathbb{P}(X = x) = \int_0^1 \binom{n}{x} p^x (1 - p)^{n-x} dp,$$

Resolution 2/2

then

$$\begin{aligned}\mathbb{P}(a < p < b | X = x) &= \frac{\int_a^b \binom{n}{x} p^x (1-p)^{n-x} dp}{\int_0^1 \binom{n}{x} p^x (1-p)^{n-x} dp} \\ &= \frac{\int_a^b \binom{n}{x} p^x (1-p)^{n-x} dp}{\mathcal{B}(x+1, n-x+1)},\end{aligned}$$

i.e., $p|x \sim \mathcal{B}(x+1, n-x+1)$.

(Beta distribution)