Machine Learning 2020

Volker Roth

Department of Mathematics & Computer Science University of Basel

26th February 2020

Volker Roth (University of Basel)

Machine Learning 2020

A B M A B M

Section 1

Probabilities

Э

<ロト < 回 > < 回 > < 回 > < 回 >

Outline

- Probability basics
- Some important probability distributions
- Some important statistical concepts



Probability theory vs Statistics

Definition (Probability Theory)

A branch of mathematics concerned with the analysis of random phenomena.

 $\mathsf{General} \Rightarrow \mathsf{Specific}$

Definition (Statistics)

The science of collecting, analyzing, presenting, and interpreting data.

 $\mathsf{Specific} \Rightarrow \mathsf{General}$

- Machine learning is closely related to (inferential) statistics.
- Learning algorithms are often probabilistic algorithms.

Probabilities

Definition (Probability Space)

A probability space is the triple

 (Ω, S, P)

where

- Ω is the sample/outcome space, ω ∈ Ω is a sample point/atomic event.
 Example: 6 possible rolls of a die: Ω = {1, 2, 3, 4, 5, 6}
- S is a collection of events to which we are willing to assign probabilities. An event a ∈ S is any subset of Ω, e.g., die roll < 4: a = {1,2,3}
- P is a mapping from events in S to \mathbb{R} that satisfies the probability axioms.

< ロ > < 同 > < 三 > < 三 >

Axioms of Probability

- $P(a) \ge 0 \ \forall a \in S$: probabilities are not negative,
- $P(\Omega) = 1$: "trivial" event has maximal possible prob 1,
- a, b ∈ S and a ∩ b = { } ⇒ P(a ∪ b) = P(a) + P(b): probability of two mutually disjoint events is the sum of their probabilities.
 Example:

P(die roll < 4) = P(1) + P(2) + P(3) = 1/6 + 1/6 + 1/6 = 1/2.

Random Variables

Definition (Random Variable)

A **random variable** X is a function from the sample points to some range, e.g., the reals

 $X: \Omega \to \mathbb{R},$

or booleans

 $X: \Omega \rightarrow \{$ true,false $\}.$

Real random variables are characterized by their distribution function.

Definition (Cumulative Distribution Function)

Let $X:\Omega
ightarrow \mathbb{R}$ be a real valued random variable. We define

 $F_X(x) = P(X \le x).$

This is the probability of the event $\{\omega \in \Omega : X(\omega) \le x\}$

< ロ > < 同 > < 回 > < 回 > < 回 > <

Probability and Propositions

- **Proposition**: event (set of sample points) where the **proposition is** true.
- Given Boolean random variables A and B:
 - event $a = \text{set of atomic events where } A(\omega) = \text{true}$
 - event $\neg a =$ set of atomic events where $A(\omega) =$ false
 - event $a \wedge b$ = atomic events where $A(\omega)$ = true and $B(\omega)$ = true
- With Boolean variables,

event = propositional logic model

e.g., A =true, B =false, or $a \land \neg b$.

• Proposition = disjunction of events in which it is true e.g., $(a \lor b) = (\neg a \land b) \lor (a \land \neg b) \lor (a \land b)$ $\implies P(a \lor b) = P(\neg a \land b) + P(a \land \neg b) + P(a \land b)$

くロト (得) (ヨト (ヨト) ヨ

Syntax for Propositions

• Boolean random variables

e.g., *Cavity* (do I have a cavity?)

Cavity = true is a proposition, also written cavity

- Discrete random variables (finite or infinite)
 e.g., Weather is one of (sunny, rain, cloudy, snow)
 Weather = rain is a proposition
 Values must be exhaustive and mutually exclusive
- **Continuous** random variables (*bounded* or *unbounded*) e.g., *Temp* = 21.6; also allow, e.g., *Temp* < 22.0.

Probability distribution

- Unconditional probabilities of propositions e.g., *P*(*Weather* = *sunny*) = 0.72. Bayesian interpretation: correspond to belief prior to arrival of any (new) evidence
- **Probability distribution** gives values for all possible assignments: $P(Weather) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$ (normalized, sums to 1)
- Joint probability distribution for a set of RVs gives the probability of every atomic event on those RVs (i.e., every sample point) $P(Weather, Cavity) = a 4 \times 2$ matrix of values:

Weather =	sunny	rain	cloudy	snow
<i>Cavity</i> = <i>true</i>	0.144	0.02	0.016	0.02
<i>Cavity</i> = <i>false</i>	0.576	0.08	0.064	0.08

くロト く得ト くヨト くヨト

Probability for continuous variables

Suppose X describes some uncertain continuous quantity. What is the probability that $a < X \le b$?

- define events $A = (X \le a), B = (X \le b), W = (a < X \le b).$
- $B = A \lor W$, A and W are **mutually exclusive** $\rightsquigarrow p(B) = p(A) + p(W) \rightsquigarrow p(W) = p(B) - p(A)$.
- Define the cumulative distribution function (cdf) as $F(q) := p(X \le q) \rightsquigarrow p(a < X \le b) = F(b) F(a).$
- Assume that F is absolutely continuous: define **probability density** function (pdf) $p(x) := \frac{d}{dx}F(x)$.
- Given a pdf, the probability of a continuous variable being in a finite interval is: P(a < X ≤ b) = ∫_a^b p(x) dx.
- As the size of the interval gets smaller, we can write $P(x < X \le x + dx) \approx p(x) dx$.
- We require p(x) ≥ 0, but it is possible for p(x) > 1, so long as the density integrates to 1.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

Probability for continuous variables



Left: cdf for the standard normal, $\mathcal{N}(0,1)$. Right: corresponding pdf.

- Shaded regions each contain $\alpha/2$ of the probability mass \rightsquigarrow nonshaded region contains 1α .
- Left cutoff point is $\Phi^{-1}(\alpha/2)$, Φ is cdf of standard Gaussian.
- By symmetry, the right cutoff point is $\Phi^{-1}(1 \alpha/2) = -\Phi^{-1}(\alpha/2)$.
- If $\alpha = 0.05$, the central interval is 95%, left cutoff is -1.96, right cutoff is 1.96.

Probability for continuous variables

Example: uniform distribution:

$$\operatorname{Unif}(a,b) = \frac{1}{b-a}\mathbb{I}(a \le x \le b).$$



p(X = 20.5) = 0.125 really means

$$\lim_{dx \to 0} P(20.5 \le X \le 20.5 + dx)/dx = 0.125$$

Mean and Variance

- Most familiar property of a distribution: mean, or expected value, denoted by μ.
- Discrete RVs:

$$E[X] = \sum_{x \in \mathcal{X}} xp(x),$$

Continuous RVs:

$$E[X] = \int_{\mathcal{X}} xp(x) \, dx.$$

If this integral is not finite, the mean is not defined.

- The variance is a measure of the spread of a distribution: $var[X] = E[(X - \mu)^2] = E[X^2] - \mu^2 =: \sigma^2.$
- The square root $\sqrt{\operatorname{var}[X]}$ is the standard deviation.

Common discrete distributions: Binomial and Bernoulli

- Toss a coin *n* times. Let $X \in \{0, ..., n\}$ be the number of heads.
- If the probability of heads is θ, then we say the RV X has a binomial distribution, X ~ Bin(n, θ):

$$\operatorname{Bin}(k|n,\theta) = \binom{n}{k} \theta^k (1-\theta)^{n-k}.$$

Special case for n = 1: Bernoulli distribution.
 Let X ∈ {0,1} → binary random variable. Let θ be the probability of success. We write X ~ Ber(θ).

$$\mathsf{Ber}(x| heta) = heta^{\mathbb{I}(x=1)}(1- heta)^{\mathbb{I}(x=0)}.$$

In other words,

$$\mathsf{Ber}(x| heta) = egin{cases} heta, & ext{if } x = 1 \ 1 - heta, & ext{if } x = 0. \end{cases}$$

Common discrete distributions: Multinomial

- Tossing a K-sided die \rightsquigarrow can use the multinomial distribution.
- Let $X = (X_1, X_2, ..., X_K)$ be a **random vector**. Let x_j be the number of times side j of the die occurs.

$$\mathsf{Mu}(\boldsymbol{x}|n,\boldsymbol{\theta}) = \binom{n}{x_1\cdots x_K} \prod_{j=1}^K \theta_j^{x_j},$$

where θ_j is the probability that side j shows up, and

$$\binom{n}{x_1\cdots x_K} = \frac{n!}{x_1!x_2!\cdots x_K!}$$

is the **multinomial coefficient** (the number of ways to divide a set of size $n = \sum_{k=1}^{K} x_k$ into subsets with sizes x_1 up to x_K).

Common discrete distributions: Multinoulli

- Special case for n = 1: Mutinoulli distribution.
- Rolling a *K*-sided dice once, so *x* will be a vector of 0s and 1s, in which only one bit can be turned on.
- If the dice shows up as face k, then the k'th bit will be on
 → think of x as being a scalar categorical random variable with K states, and x is its dummy encoding.
- Example: K = 3, encode the states 1, 2 and 3 as (1,0,0), (0,1,0), and (0,0,1).
- Also called a one-hot encoding, since we imagine that only one of the K "wires" is "hot" or on.

$$\mathsf{Mu}(\boldsymbol{x}|1, \boldsymbol{ heta}) = \prod_{j=1}^{K} \theta_{j}^{\mathbb{I}(x_{j}=1)}.$$

Common discrete distributions: Empirical

Given a set of data, D = {x₁,..., x_N}, define the empirical distribution, a.k.a. empirical measure:

$$p_{emp}(A) = \frac{1}{N} \sum_{i=1}^{N} \delta_{x_i}(A),$$

where

$$\delta_x(A) = egin{cases} 0, & ext{if } x
ot\in A \ 1, & ext{if } x \in A \end{cases}$$

• In general, we can associate weights with each sample:

$$p(x) = \sum_{i=1}^{N} w_i \delta_{x_i}(x)$$

where we require $0 \le wi \le 1$ and $\sum_{i=1}^{N} w_i = 1$.

- We can think of this as a histogram, with "spikes" at the data points x_i , where w_i determines the height of spike *i*.
- This distribution assigns 0 probability to any point not in the data set.

Common continuous distributions: Normal

• The pdf of the normal distribution is

$$p(x \mid \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}},$$

where μ is the mean, σ^2 is the variance. The inverse variance is sometimes called **precision**.

• The cdf of the standard normal distribution is the integral

$$\Phi(x)=\frac{1}{\sqrt{2\pi}}\int_{-\infty}^{x}e^{-t^2/2}\,dt.$$

It has no closed form expression.

• The cdf is sometimes expressed in terms of the error function

$$erf(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt,$$

as follows:

$$\Phi(x) = \frac{1}{2} \left[1 + erf\left(\frac{x}{\sqrt{2}}\right) \right].$$

Common continuous distributions: Normal

- If σ tends to zero, p(x) tends to zero at any x ≠ μ, but grows without limit if x = μ, while its integral remains equal to 1.
- Can be defined as a generalized function: **Dirac's delta function** δ translated by the mean: $p(x) = \delta(x \mu)$, where

$$\delta(x) = \begin{cases} +\infty, & x = 0 \\ 0, & x \neq 0, \end{cases}$$

additionally constrained to satisfy the identity

$$\int_{-\infty}^{\infty} \delta(x) \, dx = 1.$$

• Sifting property: selecting out a single term from a sum or integral: $\int_{-\infty}^{\infty} f(x)\delta(x-z) \, dx = f(z)$

since the integrand is only non-zero if x - z = 0.

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Central Limit Theorem

- Under certain (fairly common) conditions, the sum of many random variables will have an **approximately normal distribution**.
- Let X_1, \ldots, X_n be i.i.d. RVs with the same (arbitrary) distribution, zero mean, and variance σ^2 .
- Let

$$Z = \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^{n} X_i \right)$$

• Then, as *n* increases, the probability distribution of *Z* will tend to the normal distribution with zero mean and variance σ^2 .

Common continuous distributions: Beta

- The beta distribution is supported on the unit interval [0, 1]
- For $0 \le x \le 1$, and shape parameters $\alpha, \beta > 0$, the pdf is

$$f(x; \alpha, \beta) = rac{1}{\mathrm{B}(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

The beta function, B, is a normalization constant to ensure that the total probability is 1.



Common continuous distributions: Multivariate Normal

• The multivariate normal distribution of a k-dimensional random vector $\mathbf{X} = (X_1, \ldots, X_k)^t$ can be written as: $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with k-dimensional mean vector

$$\boldsymbol{\mu} = \mathsf{E}[\mathbf{X}] = [\mathsf{E}[X_1], \mathsf{E}[X_2], \dots, \mathsf{E}[X_k]]^{\mathrm{t}}$$

and $k \times k$ covariance matrix

$$\Sigma =: \mathsf{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^{\mathrm{t}}] = [\mathsf{Cov}[X_i, X_j]; 1 \le i, j \le k],$$

where

$$\operatorname{Cov}[X_i, X_j] = \mathsf{E}[(X_i - \mu_i)(X_j - \mu_j)].$$

- The inverse of the covariance matrix is the precision matrix $Q = \Sigma^{-1}$.
- The pdf of the multivariate normal distribution is

$$p(x_1,\ldots,x_k) = rac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-rac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{t}} \Sigma^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)$$

The 2D Normal distribution



Affine transformations: If $\mathbf{y} = \mathbf{c} + B\mathbf{x}$ is an affine transformation of $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, then $\mathbf{Y} \sim \mathcal{N}(\mathbf{c} + B\boldsymbol{\mu}, B\boldsymbol{\Sigma}B^{t})$

Volker Roth (University of Basel)

Machine Learning 2020

The 2D Gaussian distribution

2D Gaussian:
$$p(\mathbf{x}|\boldsymbol{\mu} = \mathbf{0}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{2\pi|\boldsymbol{\Sigma}|}} \exp(-\frac{1}{2}\mathbf{x}^t \boldsymbol{\Sigma}^{-1} \mathbf{x})$$

Covariance

(also written "*co*-variance") is a measure of how much **two** random variables vary together:

- +1: perfect linear coherence,
- -1: perfect negative linear coherence,
- 0: no linear coherence.



Common continuous distributions: Dirichlet

- The Dirichlet distribution of order K ≥ 2 with parameters α₁,..., α_K > 0 is a multivariate generalization of the beta distribution.
- $\bullet~$ Its pdf on \mathbb{R}^{K-1} is

$$f(x_1,\ldots,x_K;\alpha_1,\ldots,\alpha_K)=\frac{1}{\mathrm{B}(\boldsymbol{\alpha})}\prod_{i=1}^K x_i^{\alpha_i-1},$$

where $\{x_k\}_{k=1}^{k=K}$ belong to the standard K-1 simplex: $\sum_{i=1}^{K} x_i = 1 \text{ and } x_i \ge 0 \forall i \in [1, K]$

• The normalizing constant is the multivariate beta function.

Common continuous distributions: Dirichlet



wikimedia.org/w/index.php?curid=49908662

Machine Learning 2020

Conditional probability

• Conditional or posterior probabilities

e.g., P(cavity|toothache) = 0.8
i.e., given that toothache is all I know
NOT "if toothache then 80% chance of cavity"

- Notation for conditional distributions:
 P(Cavity|Toothache) = 2-element vector of 2-elem. vectors.
- If we know more, e.g., cavity is also given, then we have P(cavity|toothache, cavity) = 1
 Note: the less specific belief remains valid after more evidence arrives, but is not always useful.
- New evidence may be **irrelevant**, allowing **simplification**: P(cavity|toothache, die roll = 3) = P(cavity|toothache) = 0.8

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Conditional probability

• Definition of conditional probability:

$$P(a|b) = rac{P(a \wedge b)}{P(b)}$$
 if $P(b) \neq 0$

Product rule gives an alternative formulation: $P(a \land b) = P(a|b)P(b) = P(b|a)P(a)$

- A general version holds for whole **distributions**, e.g., P(Weather, Cavity) = P(Weather|Cavity)P(Cavity)
- **Chain rule** is derived by successive application of product rule: $P(X_1,...,X_n) = P(X_1,...,X_{n-1}) P(X_n|X_1,...,X_{n-1})$ $= P(X_1,...,X_{n-2}) P(X_{n-1}|X_1,...,X_{n-2}) P(X_n|X_1,...,X_{n-1})$ = ... $= \prod_{i=1}^n P(X_i|X_1,...,X_{i-1})$

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Inference by enumeration

Start with the joint distribution:

	toothache		⊐ toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

For any proposition ϕ , sum the atomic events where it is true: P(toothache) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2

Inference by enumeration

Start with the joint distribution:

	toothache		⊐ toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

For any proposition ϕ , sum the atomic events where it is true: $P(cavity \lor toothache) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$

Inference by enumeration

Start with the joint distribution:

	toothache		⊐ toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

Can also compute conditional probabilities:

$$P(\neg cavity | toothache) = \frac{P(\neg cavity \land toothache)}{P(toothache)}$$
$$= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} = 0.4$$

Normalization

	toothache		⊐ too	⊐ toothache	
	catch	\neg catch	catch	\neg catch	
cavity	.108	.012	.072	.008	
\neg cavity	.016	.064	.144	.576	

Denominator can be viewed as a normalization constant α

 $\mathbf{P}(Cavity | toothache) = \alpha \, \mathbf{P}(Cavity, toothache)$

- $= \alpha [\mathbf{P}(Cavity, toothache, catch) + \mathbf{P}(Cavity, toothache, \neg catch)]$
- $= \alpha [\langle 0.108, 0.016 \rangle + \langle 0.012, 0.064 \rangle]$
- $= \alpha \left< 0.12, 0.08 \right> = \left< 0.6, 0.4 \right>$

General idea: compute distribution on query variable by fixing **evidence variables** and summing over **hidden variables**.

Machine Learning 2020

Inference by enumeration, contd.

Let X be all the variables. Typically, we want the posterior joint distribution of the **query variables** Y given specific values e for the **evidence variables** E

Let the **hidden variables** be $\mathbf{H} = \mathbf{X} - \mathbf{Y} - \mathbf{E}$

Then the required summation of joint entries is done by **summing out the hidden variables**:

$$P(\mathbf{Y}|\mathbf{E}=\mathbf{e}) = \alpha P(\mathbf{Y}, \mathbf{E}=\mathbf{e}) = \alpha \sum_{\mathbf{h}} P(\mathbf{Y}, \mathbf{E}=\mathbf{e}, \mathbf{H}=\mathbf{h})$$

Joint probability $p(x) = p(x_1, ..., x_n) \rightsquigarrow$ number of states: $\prod_{i=1}^{n} |arity(x_i)|.$ Obvious problems:

1) Worst-case time complexity $O(d^n)$ where d is the largest arity

- 2) Space complexity $O(d^n)$ to store the joint distribution
- 3) How to find the numbers for $O(d^n)$ entries???

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Inference in Jointly Gaussian Distributions: Marginalization

$$\mathbf{x} \sim \mathcal{N}(\mathbf{\mu}, \mathbf{\Sigma})$$
. Let $\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$ and $\mathbf{\Sigma} = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$.
Then $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{\mu}_1, \Sigma_{11})$ and $\mathbf{x}_2 \sim \mathcal{N}(\mathbf{\mu}_2, \Sigma_{22})$.



Marginals of Gaussians are again Gaussian!

Volker Roth (University of Basel)

Machine Learning 2020

Inference in Jointly Gaussian Distributions



wikimedia.org/w/index.php?curid=25235145

Inference in Jointly Gaussian Distributions

$$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
. Let $\mathbf{x} = \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}$ and $\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$.
Then $\mathbf{x}_2 | \mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{22} - \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12})$.



Conditionals of Gaussians are again Gaussian!

Volker Roth (University of Basel)

Machine Learning 2020

Independence

A and B are **independent** iff $\mathbf{P}(A|B) = \mathbf{P}(A)$ or $\mathbf{P}(B|A) = \mathbf{P}(B)$ or $\mathbf{P}(A, B) = \mathbf{P}(A)\mathbf{P}(B)$



$$\begin{split} & \mathsf{P}(\textit{Toothache, Catch, Cavity, Weather}) \\ &= \mathsf{P}(\textit{Toothache, Catch, Cavity})\mathsf{P}(\textit{Weather}) \\ &\rightsquigarrow 4 \cdot 8 = 32 \text{ entries reduced to } 4 + 8 = 12. \\ & \mathsf{Absolute independence powerful but rare...} \\ & \mathsf{Dentistry is a large field with hundreds of variables,} \\ & \mathsf{none of which are independent. What to do?} \end{split}$$

くロト く得ト くほト くほう

Conditional independence

P(Toothache, Cavity, Catch) has $2^3 - 1 = 7$ independent entries

If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:

(1) P(catch|toothache, cavity) = P(catch|cavity)

The same independence holds if I haven't got a cavity: (2) $P(catch|toothache, \neg cavity) = P(catch|\neg cavity)$

Catch is **conditionally independent** of Toothache given Cavity: P(Catch|Toothache, Cavity) = P(Catch|Cavity)

Equivalent statements:

$$\begin{split} & \mathsf{P}(\mathit{Toothache}|\mathit{Catch},\mathit{Cavity}) = \mathsf{P}(\mathit{Toothache}|\mathit{Cavity}) \\ & \mathsf{P}(\mathit{Toothache},\mathit{Catch}|\mathit{Cavity}) = \mathsf{P}(\mathit{Toothache}|\mathit{Cavity}) \mathsf{P}(\mathit{Catch}|\mathit{Cavity}) \end{split}$$

イロト 不得 トイヨト イヨト 二日

Conditional independence contd.

Write out full joint distribution using chain rule:

P(*Toothache*, *Catch*, *Cavity*)

- $= \mathbf{P}(\mathit{Toothache}|\mathit{Catch},\mathit{Cavity})\mathbf{P}(\mathit{Catch},\mathit{Cavity})$
- $= {\sf P}(\mathit{Toothache}|\mathit{Catch},\mathit{Cavity}) {\sf P}(\mathit{Catch}|\mathit{Cavity}) {\sf P}(\mathit{Cavity})$
- $= \mathsf{P}(\mathit{\textit{Toothache}}|\mathit{\textit{Cavity}})\mathsf{P}(\mathit{\textit{Catch}}|\mathit{\textit{Cavity}})\mathsf{P}(\mathit{\textit{Cavity}})$
- I.e., only 2 + 2 + 1 = 5 independent numbers.

Often, conditional independence reduces the size of the representation of the joint distribution from **exponential** in n to **linear** in n.

Conditional independence is our most basic and robust form of knowledge about uncertain environments.

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Bayes Rule

Bayes Rule $P(B|A) = \frac{P(A|B)P(B)}{P(A)}$

Proof.

$$P(A|B)P(B) = P(A,B) = P(B|A)P(A)$$



イロト イヨト イヨト

Bayes Rule (cont'd)

• Useful for assessing diagnostic probability from causal probability:

$$P(Cause | Effect) = \overbrace{\frac{P(Effect | Cause) P(Cause)}{P(Effect)}}^{Prob(symptoms) Prevalence}$$

E.g., let *M* be meningitis (acute inflammation of the protective membranes covering the brain and spinal cord),
 S be stiff neck. Assume the doctor knows that the prevalence of meningitis is 1/50,000, that the prior probability of a stiff neck is 0.01, and that the symptom stiff neck occurs with a probability of 0.7.

$$P(m|s) = rac{P(s|m)P(m)}{P(s)} = rac{0.7 imes 1/50000}{0.01} = 0.0014.$$

• Note: the posterior probability of meningitis is still very small (1 in 700 patients)!

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

Bayes rule (cont'd)

Question: Why should it be easier to estimate the conditional probabilities in the causal direction P(Effect|Cause), as compared to the diagnostic direction, P(Cause|Effect)?

There are two possible answers (in a medical setting):

- We might have access to a collection of health records for patients having meningitis. This collections will provide us with estimates of P(s|m). For directly estimating P(m|s) we would need a database of all cases of the very unspecific symptom.
- Diagnostic knowledge might be **more fragile** than causal knowledge. Assume a doctor has directly estimated P(m|s). If there is a sudden epidemic of meningitis, P(m) will go up, but this doctor will have no idea how to update P(m|s). The other doctor who uses Bayes rule knows that P(m|s) should go up proportionately with p(m). Note that causal information P(s|m) is **unaffected by the epidemic** (it simply reflects the way how meningitis works)!

・ロト ・ 一 ト ・ ヨ ト ・ ヨ ト ・ ヨ

Bayes' Rule and conditional independence

 $P(Cavity | toothache \land catch)$

- $= \alpha P(toothache \land catch|Cavity)P(Cavity)$
- $= \alpha P(toothache|Cavity)P(catch|Cavity)P(Cavity)$

This is an example of a naive Bayes model:

 $P(Cause, Effect_1, \dots, Effect_n) = P(Cause) \prod_i P(Effect_i | Cause)$



Total number of parameters is **linear** in *n*

Example: Wumpus World

- The wumpus is a beast that eats anyone who enters the room.

- Some rooms contain bottomless pits that will trap anyone entering the room (except for the wumpus, which is too big to fall in!)
- The only positive aspect is the possibility of finding ${\boldsymbol{gold}}...$



Squares adjacent to wumpus are **smelly**. Squares adjacent to pit are **breezy**. **Glitter** if and only if gold is in the same square. **Shooting** kills the wumpus if you are facing it. Shooting uses up the only **arrow**. **Grabbing** picks up the gold if in the same square. **Releasing** drops the gold in the same square.

Goal: Get gold back to start without entering pit or wumpus square

Wumpus World

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
^{1,2} B OK	2,2	3,2	4,2
1,1	^{2,1} B	3,1	4,1
OK	OK		

 $P_{ij} = true \text{ iff } [i, j] \text{ contains a pit}$ $B_{ij} = true \text{ iff } [i, j] \text{ is breezy}$ Include only $B_{1,1}, B_{1,2}, B_{2,1}$ in the probability model , A_{ij}

Volker Roth (University of Basel)

Machine Learning 2020

3

Specifying the probability model

The full joint distribution is $P(P_{1,1}, \ldots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$ Apply product rule: $P(B_{1,1}, B_{1,2}, B_{2,1} | P_{1,1}, \ldots, P_{4,4})P(P_{1,1}, \ldots, P_{4,4})$ (Do it this way to get P(Effect | Cause).) First term: 1 if pits are adjacent to breezes, 0 otherwise Second term: pits are placed randomly, probability 0.2 per square:

$$\mathbf{P}(P_{1,1},\ldots,P_{4,4}) = \prod_{i,j=1,1}^{4,4} \mathbf{P}(P_{i,j}) = 0.2^n \times 0.8^{16-n}$$

for *n* pits.

Observations and query

We know the following facts:

$$\begin{split} b &= \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1} \\ known &= \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1} \\ \text{Query is } \mathbf{P}(P_{1,3}|known,b) \\ \text{Define } Unknown &= P_{ij} \text{s other than } P_{1,3} \text{ and } Known \\ \text{For inference by enumeration, we have} \end{split}$$

 $\mathbf{P}(P_{1,3}|known, b) = \alpha \sum_{unknown} \mathbf{P}(P_{1,3}, unknown, known, b)$

Grows exponentially with number of squares!

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Using conditional independence

Basic insight: observations are conditionally independent of other hidden squares given neighbouring hidden squares



Define $Unknown = Fringe \cup Other$ $P(b|P_{1,3}, Known, Unknown) = P(b|P_{1,3}, Known, Fringe)$ Manipulate query into a form where we can use this!

Using conditional independence contd.

Volker

$$\begin{aligned} \mathbf{P}(P_{1,3}|known, b) &= \alpha \sum_{unknown} \mathbf{P}(P_{1,3}, unknown, known, b) \\ &= \alpha \sum_{unknown} \mathbf{P}(b|P_{1,3}, known, unknown) \mathbf{P}(P_{1,3}, known, unknown) \\ &= \alpha \sum_{unknown} \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe, other) \mathbf{P}(P_{1,3}, known, fringe, other) \\ &= \alpha \sum_{fringe} \sum_{other} \mathbf{P}(b|known, P_{1,3}, fringe) \mathbf{P}(P_{1,3}, known, fringe, other) \\ &= \alpha \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) \sum_{other} \mathbf{P}(P_{1,3}, known, fringe, other) \\ &= \alpha \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) \sum_{other} \mathbf{P}(P_{1,3}) P(known) P(fringe) P(other) \\ &= \alpha P(known) \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{fringe} \mathbf{P}(b|known, P_{1,3}, fringe) P(fringe) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{other} P(other) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{other} P(other) \\ &= \alpha'$$

Using conditional independence contd.



 $\begin{aligned} \mathsf{P}(P_{1,3}|\textit{known}, b) &= \alpha' \left< 0.2(0.04 + 0.16 + 0.16), \ 0.8(0.04 + 0.16) \right> \\ &\approx \left< 0.31, 0.69 \right> \end{aligned}$

 $\mathsf{P}(P_{2,2}|known, b) \approx \langle 0.86, 0.14 \rangle$

Summary

- Probability is a rigorous formalism for uncertain knowledge
- Joint probability distribution specifies probability of every atomic event
- Queries can be answered by summing over atomic events
- For nontrivial domains, we must find a way to reduce the joint size
- Independence and conditional independence provide the tools

Subsection 1

Origins of probabilities

E

• • = • • = •

Origins of probabilities I

Historically speaking, probabilities have been regarded in a number of different ways:

- Frequentist position: probabilities come from measurements. The assertion P(cavity) = 0.05 means that 0.05 is the fraction that would be observed in the limit of infinitely many samples. From a finite sample, we can estimate this true fraction and also calculate how accurate this estimates is likely to be.
- **Objectivist view:** probabilities are actual **properties of the universe** An excellent example: quantum phenomena.

A less clear example: coin flipping – the uncertainty is probably due to our uncertainty about the initial conditions of the coin.

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Origins of probabilities II

- **Subjectivist view**: probabilities are an **agent's degrees of belief**, rather than having any external physical significance.
- The **Bayesian view** allows any self-consistent ascription of prior probabilities to propositions, but then insists on proper **Bayesian updating** as evidence arrives.
- For example P(cavity) = 0.05 denotes the degree of belief that a random person has a cavity before we make any actual observation of that person.

Updating in the light of further evidence "person has a toothache":

 $P(cavity|toothache) = \alpha P(toothache|cavity)P(cavity)$

イロト イポト イヨト イヨト 三日

The reference class problem

- Even a strict frequentist position involves subjective analysis.
- **Example:** Say a doctor takes a frequentist approach to diagnosis. She examines a large number of people to establish the probability of whether or not they have heart disease.

To be accurate she tries to measure "similar people" (she knows for example that gender might be important.)

• Taken to an extreme, all people are different and therefore the **reference class is empty**.

This has been a vexing problem in the philosophy of science.

- Assume x₁,..., x_n are drawn i.i.d. from normal N(μ, σ²) with known variance σ². What can be said about μ?
- Frequentist view: no further probabilistic assumptions \rightsquigarrow treat μ as an **unknown constant**.

• The sample mean $\bar{x} = \sum_i x_i/n$ is the observed value of the RV $\bar{X} \sim N(\mu, \bar{\sigma}^2)$, with $\bar{\sigma}^2 = \sigma^2/n$.

Now define the linearly transformed random variable

$$B := rac{\mu - X}{ar{\sigma}} \sim N(0, 1), \ (\ {
m i.e. \ standard \ normal}).$$

Use normal cdf $\Phi(k_c) = P(B < k_c)$ to derive an upper limit for μ :

$$P(B < k_c) = \Phi(k_c) = 1 - c$$

= $P(-\bar{\sigma}B > -\bar{\sigma}k_c)$
= $P(\underbrace{\mu - \bar{\sigma}B}_{\bar{X}} > \mu - \bar{\sigma}k_c)$
= $P(\bar{X} + \bar{\sigma}k_c > \mu).$

- The statement P(μ < X̄ + σ̄k_c) = 1 − c can be interpreted as specifying a hypothetical long run of statements about the constant μ, a portion 1 − c of which is correct. (Note that X̄ is a RV!)
- Plugging in the observed x
 x, the statement μ < x
 x + σ
 k c can be interpreted as one of a long run of such statements about μ.
- Arguments involving probability only via its (hypothetical) long-run frequency interpretation are called frequentist.
- That is, in the frequentist world we define procedures for assessing evidence that are calibrated by how they would perform were they used repeatedly.

イロト イポト イヨト イヨト 三日

- From the Bayesian viewpoint, we treat μ as having a probability distribution both with and without the data. That is, μ is the unobserved value of the random variable M.
- Bayes' theorem: $p_{M|X}(\mu|x) = \alpha \times p_{X|M}(x|\mu)p_M(\mu)$.
- Intuitive idea:
 - \blacktriangleright all relevant information about μ is in the conditional distribution, given the data;
 - this distribution is determined by the elementary formulae of probability theory;
 - remaining problems are solely computational.

• Example: choose $p(\mu) = N(m, \nu^2) \rightsquigarrow p(\mu|x) = N(\tilde{m}, \tilde{\nu}^2)$ with $\tilde{m} = \frac{\bar{x}/\bar{\sigma}^2 + m/\nu}{1/\bar{\sigma}^2 + 1/\nu^2}, \quad \tilde{\nu}^2 = \frac{1}{1/\bar{\sigma}^2 + 1/\nu^2}$

"Normal likelihood times normal prior gives normal posterior"

イロト イポト イヨト イヨト 三日

- Same reasoning as before: define transfomed $ilde{B}:=rac{\mu- ilde{m}}{ ilde{
 u}}\sim N(0,1)$
- Upper limit for μ : $P(\mu < \tilde{m} + k_c \tilde{\nu}) = 1 c$.
- If the prior variance $\nu^2 \gg \bar{\sigma}^2$ and the prior mean m is not too different from \bar{X} , this limit agrees closely with the one obtained by the frequentist method (because then $\tilde{m} \approx \bar{X}$ and $\tilde{\nu} \approx \bar{\sigma}$).
- This broad parallel between the different types of analysis is in no way specific to the normal distribution.
- See the beautiful book (D.R. Cox, Principles of statistical inference, Cambridge, 2006) for further details.

Subsection 2

Some important statistical concepts

Э

Convergence of random variables

Definition (Convergence in Probability)

Let X_1, X_2, \ldots be random variables. We say that X_n converges in probability to the random variable X as $n \to \infty$, iff, for all $\varepsilon > 0$,

$$P(|X_n - X| > \varepsilon) \to 0$$
, as $n \to \infty$.

We write $X_n \xrightarrow{p} X$ as $n \to \infty$.

Example: Weak law of large numbers

Theorem (Bernoulli's Theorem, Weak law of large numbers)

Let X_1, X_2, \ldots , be a sequence of independent and identically distributed (i.i.d.) random variables, each having mean μ (and standard deviation σ). Let $S_n = X_1 + \ldots + X_n$. Then

$$P(|S_n/n-\mu|>\varepsilon)\to 0$$

as $n
ightarrow \infty$.

- Given sufficiently many observations x_i , the sample mean $\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$ will approach the true mean μ .
- Note that |S_n/n − μ| > ε might happen an infinite number of times, although at infrequent intervals.
- The strong law even says that for any $\varepsilon > 0$ the inequality $|S_n/n \mu| < \varepsilon$ holds for all large enough *n*, but we will not discuss this further...

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

Example: convergence of empirical CDF

Definition (Empirical cumulative distribution function)

Let X_1, X_2, \ldots, X_n be iid real random variables with the common cdf F(t). Then the empirical distribution function is defined as

$$\hat{F}_n(t) = rac{\#(ext{elements}) ext{ in the sample} \leq t}{n} = rac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i \leq t\},$$

where $\mathbf{1}$ {*A*} is the indicator of event *A*.

- For a fixed t, the indicator 1{X_i ≤ t} is a Bernoulli random variable with mean μ = F(t).
- Weak law of large numbers \Rightarrow estimator $\hat{F}_n(t)$ converges in probability to F(t) as $n \to \infty$, for every value of t:

 $\hat{F}_n(t) \xrightarrow{p} F(t).$

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Expectation

Definition (Expectation)

Let X be a random variable with probability density function f_X . The expectation is

$$E[X] := \int_{-\infty}^{\infty} x f_X(x) \, dx.$$

Definition (Sample mean)

Let a sample $x = \{x_1, x_2, \dots, x_n\}$ be given. The sample mean is

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i.$$

The sample mean is an **unbiased estimator** of $\mu = E[X]$.

・ロト ・ 一日 ・ ・ 日 ・ ・ 日 ・

Variance

Definition (Variance)

Let X be a random variable with density function f_X . The **variance** is

$$Var[X] = E[(X - E[X])^2] = E[X^2] - (E[X])^2.$$

The square root $\sqrt{Var[X]}$ is the standard deviation.

Definition (Sample Variance)

Let the sample $x = \{x_1, x_2, ..., x_n\}$ with sample mean \overline{x} be given. The **sample variance** is

$$s^{2} = \frac{1}{n-1} \sum_{i=1}^{n} (x_{i} - \overline{x})^{2}.$$

The sample variance is an **unbiased estimator** of Var[X].

Volker Roth (University of Basel)

・ロト ・ 同ト ・ ヨト ・ ヨト … ヨ

Bias and Variance of an Estimator

- Assume that a statistical model parametrized by θ gives rise to a **probability distribution for observed data**, $P(x|\theta)$.
- Let θ̂ be an estimator of θ based on any observed data x, i.e. θ̂ maps observed x to values that we hope are close to θ.
- $\bullet\,$ The bias of $\hat{\theta}$ is defined to be

$$\mathsf{Bias}[\hat{\theta}] = E_{P(x|\theta)}[\hat{\theta}] - \theta = E_{P(x|\theta)}[\hat{\theta} - \theta],$$

where $E_{P(x|\theta)}[\cdot]$ denotes expected value over the distribution $P(x|\theta)$, i.e. averaging over all possible observations x.

- An estimator is unbiased if its bias is zero for all values of parameter θ.
- The variance of $\hat{\theta}$ is the expected value of the squared sampling deviations: $var(\hat{\theta}) = E[(\hat{\theta} E(\hat{\theta}))^2].$