

Probability Review

Reminders: HW due Thurs. will follow this, but
 Readings: Murphy §2, §3 ↗ adopt standard probability/sets
 Bishop §2 notation

Note: I will update these notes
 as course progresses to maintain
 probability you will need to
 know

For extra resources filling in gaps in background

- "Probability Essentials" Jacod + Protter
 - Short book w/ self-contained chapters
 - + exercises - good for self-study
- "Probability: theory + examples" Durrett
 - Measure theoretic + rigorous coverage
 - for those interested in serious research on probability + ML

» Questions on HW? (collognow)

Def Probability space is a triple $(\Omega, \mathcal{F}(\Omega), P)$
 consisting of:

- Ω - a sample space
 all events which may occur
 ex Flipping two coins
 $\Omega = \{\text{H}_1, \text{H}_2, \text{T}_1, \text{T}_2\}$

- \mathcal{F} - an event space, typically taken as the "σ-algebra" consisting of countably finite intersections & unions of elements from Ω

ex the event that both coins are heads $H_1 \cap H_2$

- P - a probability measure

$$0 \leq P(A) \leq 1, A \in \mathcal{F}$$

↳ frequentist perspective

Given $N \rightarrow \infty$ identical experiments,
 the proportion where the event happens

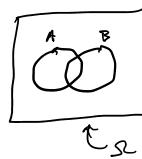
↳ Bayesian perspective

A "modeling function" describing how likely an event is to occur

- Probability of two events

def Given $A, B \in \mathcal{F}$
 AND $P(A, B) := P(A \cap B)$

OR $P(A \cup B)$



def set operations + def's

- Set subtraction $A \setminus B$ or $A - B = A \cap B^c$
- deMorgan's laws $(\bigcup A_i)^c = \bigcap A_i^c$
 $(\bigcap A_i)^c = \bigcup A_i^c$
- Distributive law $A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$
 $A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$
- Cartesian product $A \times B = \{(x, y) \mid x \in A, y \in B\}$
- cardinality $\#A = \text{number of elements in } A$
- inclusion/exclusion principle $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
 $\dots \text{or more generally} \dots$

- inclusion/exclusion principle

$$\begin{aligned} P(A \cup B) &= P(A) + P(B) - P(A \cap B) \\ P(\bigcup_{i=1}^n A_i) &= \sum_{i=1}^n P(A_i) - \sum_{i < j} P(A_i \cap A_j) \\ &\quad \dots + (-1)^{n+1} P(A_1 \cap \dots \cap A_N) \end{aligned}$$

exercise Prove that $A \subseteq B \Rightarrow P(A) \leq P(B)$

$$\begin{aligned} P(B) &= P(A \cup (B \setminus A)) \\ &= P(A) + P(B \setminus A) \quad \leftarrow \text{no intersection} \\ &\geq P(A) \quad \underbrace{\geq 0}_{\geq 0} \end{aligned}$$

Def Conditional Probability & Dependence

- The probability of B conditioned on A
 $P(B|A) := \frac{P(A, B)}{P(A)}$
- Events A & B are independent if
 $P(A, B) = P(A) P(B)$
 i.e. $P(B|A) = P(B)$ (Knowledge of A tells you nothing about B)
 → we write $A \perp B$ in shorthand
- Events A_1, \dots, A_N are conditionally independent on B
 if $P(A_1, \dots, A_N | B) = \prod_{i=1}^N P(A_i | B)$

Random Variables

A function which takes a value w/ a given probability

def a RV taking a finite set of values
 $X = \{x_1, \dots, x_N\}$ is a discrete RV

sample space
 We denote $p(x_i) = P(X = x_i)$ as the probability mass function. noting

$$\begin{aligned} \textcircled{1} \quad 0 &\leq p(x_i) \leq 1 \\ \textcircled{2} \quad \sum_{x \in X} p(x) &= 1 \end{aligned}$$

def a RV X taking values over \mathbb{R} is a continuous RV. We describe prob. over intervals:

• $P(a \leq X \leq b) = P(X \leq b) - P(X \leq a)$

CDF/PDF • we define $F_X(x) = P(X \leq x)$ as cumulative distribution function (CDF)

• If F is differentiable we can work w/ the probability density function (PDF)

$$f_X(x) = \frac{d}{dx} F_X(x)$$

Allowing us to work w/ either

$$P(a \leq X \leq b) = \int_a^b f(x) dx = F(b) - F(a)$$

Types of RVs

Joint distribution $f(x_1, \dots, x_N) = P(X_1 = x_1, \dots, X_N = x_N)$

Marginal distribution $f(\Sigma = x) = \sum_y P(\Sigma = x, I = y)$
 aka rule of total probability

conditional distribution $f(I = y | \Sigma = x) = \frac{f(\Sigma = x, I = y)}{f(\Sigma = x)}$

conditional distribution $f(Y=y | \Sigma=x) = \frac{f(\Sigma=x, Y=y)}{f(\Sigma=x)}$

product rule $f(x, y) = f(x_1) f(y)$

probabilistic chain rule $f(x_1, \dots, x_N) = f(x_1, \dots, x_N | x_1) f(x_1)$
 $\vdots = f(x_1, \dots, x_N | x_1, x_2) f(x_2 | x_1) f(x_1)$
 $= f(x_N | x_1, \dots, x_{N-1}) \dots f(x_2 | x_1) f(x_1)$

marginal independence $x \perp y \Leftrightarrow f(x, y) = f(x) f(y)$ or more generally
 $f(x_1, \dots, x_N) = \prod_{i=1}^N f(x_i)$

conditional independence $x \perp z \perp y \perp z \Leftrightarrow f(x, y | z) = f(x | z) f(y | z)$

Expectations, Variances + other Moments

def $E[\Sigma] := \int x f_\Sigma(x) dx$ or $\sum_{x \in \Sigma} x p(x)$ expectation
 for cont. RV for discrete RV

linearity $E[\alpha \Sigma + \beta] = \alpha E[\Sigma] + \beta$

indep. $E_{x_1, \dots, x_N} \left[\prod_{i=1}^N \Sigma_i \right] = \prod_{i=1}^N E[x_i]$ if x_i are indep

def $\text{Var}(\Sigma) := E_x \left[(x - \underbrace{E_x(x)}_{\mu})^2 \right] = \sigma^2$

equivalently $\approx E_x [x^2] - \mu^2$

so $E[x^2] = \mu^2 + \sigma^2$

linearity $\text{Var}[\alpha \Sigma + \beta] = \alpha^2 \text{Var}[\Sigma]$

Law of total expectation $E_x[x] = E_y [E_x[x | Y]]$

Law of total variance $\text{Var}[\Sigma] = E [\text{Var}[x | Y]] + \text{Var}[E[x | Y]]$

Bayes rule From definition of conditional distribution

$$P(x, y) = P(x | y) P(y) = P(y | x) P(x)$$

$$\Rightarrow P(y | x) = \frac{P(x | y) P(y)}{P(x)}$$

\uparrow posterior distribution \uparrow likelihood \uparrow marginal likelihood
 prior knowledge of y w/o access to x

think of x as hidden data
 we can't directly infer

we can
infer

variables

Marginal

w/o access
to x

Probability Distributions

Lots to cover → we'll focus on two
that form the cornerstone of classification
and regression

Bernoulli Probability of coin landing heads
+ given by $0 \leq \theta \leq 1$

binomial dist $Y \sim \text{Ber}(\theta)$

$$P(Y) = \begin{cases} 1-\theta & Y=0 \\ \theta & Y=1 \end{cases}$$

$$= \theta^Y (1-\theta)^{1-Y}$$

Given N experiments, how many are heads?

$$S = \sum_{i=1}^N Y_i$$

$S \sim \text{Bin}(N, \theta)$

$$P(S) = \binom{N}{S} \theta^S (1-\theta)^{N-S}$$

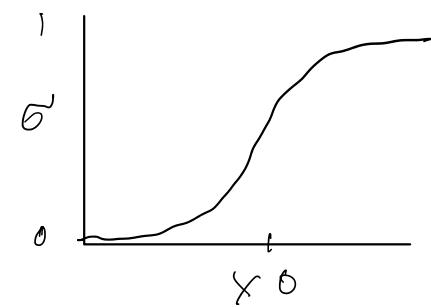
$\binom{N}{S}$ is "N choose K" $= \frac{N!}{(N-S)! S!}$

Sigmoid
Softmax

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \frac{d}{dx} \sigma(x)$$

$$\sigma_+(x) = \log(1 + e^x)$$



Binary
Logistic
Regression

Given set of data
we would like to perform
binary classification

$\sigma \rightarrow \infty$ covers
heaviside

$$P(y|x, \theta) = \text{Ber}(y | \sigma(w^T x + b))$$

$$p(y|x, \theta) = \text{Ber}(y | \text{O}^*(w^T x + b))$$

* Differentiable functions may be
chained together and fit w/
gradient descent HW1

Normal
Distributions

$$Y \sim N(y; \mu, \sigma^2)$$

iff

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\left(\frac{y-\mu}{\sigma}\right)^2\right)$$

$$F(y) = \int_{-\infty}^y f(y') dy'$$

$$\mathbb{E}(Y) = \int_{-\infty}^{\infty} y f(y) dy$$

$$\text{var}(Y) = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2$$

$$\mathbb{E}[Y^2] = \sigma^2 + \mu^2$$

Hetero/homo
skedastic
Regression

$$p(y|x, \theta) = N(y | f_n(x; \theta), f_o(x; \theta))$$

$$f_n \in \mathbb{R}$$

$f_o \in \mathbb{R}^+$ ← need to enforce
homo skedastic → $\sigma^2 = g(x)$
hetero skedastic → $\sigma^2 = f(x)$

$$p(y|x, \theta) = N(y | w_n^T x + b_n, \sigma^2 (w_o^T x))$$

Other Distributions

Make sure you are familiar
w/ other dist (§ 2.7 in Murphy)