

Probability Review

Reminders: HW due Thurs.  $\checkmark$  will follow other, but adopt standard probability set notation  
 Readings: Murphy §2, §3  
 Bishop §2

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? (extra material)

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

- $\Omega$  - 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 "sigma-algebra" consisting of countably finite intersections & unions of elements from  $\Omega$   
ex the event that both sides are heads  $\{\text{H}, \text{A H}\}$
- $P$  - 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

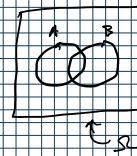
↳ Bayesian perspective

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



def set operations + def's

- Set intersection  $A \cap B$  or  $A - B = A \cap B^c$
- DeMorgan's Laws  $(\cup A_i)^c = \cap A_i^c$   
 $(\cap A_i)^c = \cup A_i^c$
- 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  $A \times B = \{(x, y) | x \in A, y \in B\}$
- cardinality  $\#A = \text{number of elements in } A$
- inclusion/exclusion principle  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$   
or more generally

- inclusion/exclusion principle

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

or more generally

$$P(\bigcup_{i=1}^n A_i) = \sum_{i=1}^n P(A_i) - \sum_{1 \leq i < j \leq n} P(A_i \cap A_j) + \dots + (-1)^{n+1} P(A_1 \cap A_2 \cap \dots \cap A_n)$$

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

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

### Def Conditional Probability & Independence

- 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 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  $p(x) = P(X=x)$  as the probability mass function. writing  
 (1)  $0 \leq p(x) \leq 1$   
 (2)  $\sum_{x \in X} p(x) = 1$

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)

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

Notation convention  
we drop the  $X$  subscript

allowing us to work w/ either

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

### Sort of RVs

Joint distribution  $f(x_1, \dots, x_n) = P(X_1=x_1, \dots, X_n=x_n)$

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

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

conditional distribution  $f(\Sigma = \mathbf{z} | \Sigma = \mathbf{x}) = \frac{f(\Sigma = \mathbf{z}, \Sigma = \mathbf{x})}{f(\Sigma = \mathbf{x})}$

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

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

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

Expectations, Variances + other moments

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

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

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

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

equivalently  $M = E_x [x^2] - M^2$

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

linearly  $\text{var}[\alpha \Sigma + \beta] = \alpha^2 \text{var}[\Sigma]$

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

Law of total variance  
 $\text{var}[\Sigma] = E_y [\text{var}[\Sigma | Y]] + \text{var}[E_x[\Sigma | Y]]$

Bayes rule from definition of conditional distribution

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

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

↑                      ↓                      ↑  
 likelihood      posterior distribution      marginal likelihood

think of  $x$  as middle data  
 we care about  $y$   
 w/o access to  $x$

prior knowledge of  $y$   
 w/o access to  $x$

we can  
see

most

w/o access  
to  $x$

## Probability Distributions

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

Binomial

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?

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

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

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

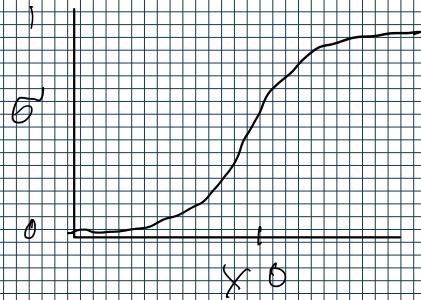
$$\binom{N}{S} \text{ is "N choose S"} = \frac{N!}{(N-S)! S!}$$

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

$\lambda \rightarrow \infty$  covers  
heaviside

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

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

\* Differentiable  
Chained  
gradient descent  
functions may be  
fogged together and fit w/  
(HW)

Normal

Distributions

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

i.e.

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

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

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

$$= \mu$$

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

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

Hetero/Homo  
Skedastic  
Regression

$$p(y|x, \theta) = N(y | f_\mu(x; \theta), f_\sigma(x; \theta))$$

$$f_\mu \in \mathbb{R}$$

$$f_\sigma \in \mathbb{R}^+ \quad \begin{matrix} \leftarrow \text{need to enforce} \\ \text{to get valid} \end{matrix}$$

$$\text{homo skedastic} \rightarrow \sigma^2 = f(x)$$

$$\text{hetero skedastic} \rightarrow \sigma^2 = f(x)$$

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

Other Distributions

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