

SPR Day 2

Readings from Bishop PRML

1.2.1 Probability Density functions

1.2.2-1.2.3 Expectations + Bayesian probabilities

2.1 Binary variables

2.1.1 Beta distribution

Gamma functions

Goals

Review ML estimation for Bernoulli data

Review continuous r.v. basic definition

Introduce Beta random variables

PDF formula

What is a Gamma function?

Introduce Beta-Bernoulli model & MAP estimation
interpret α, β parameters as pseudo counts

Next time

Big Idea: Bayesian predictive distribution

Binary Random Variables & the Bernoulli distribution ^①

Consider observing outcome of a random coin toss

We are not sure if coin is fair or not.

Let X be a random variable indicating outcome
Sample space?

1 "heads"

0 "tails"

Probability Mass function?

$$P(X=1) = \mu$$

$$\text{where } 0 \leq \mu \leq 1$$

$$P(X=0) = 1 - \mu$$

The variable μ here is a "parameter", not a "random variable" (yet)

We will often write the PMF as a conditional

$$P(X=1/\mu), \quad \text{just to remind ourselves that } \mu \text{ is required for that calculation}$$

The name for this kind of distribution is Bernoulli

$$\begin{aligned} P(X=x/\mu) &= \text{BernPMF}(X=x/\mu) = \mu^x (1-\mu)^{1-x} \\ &= \begin{cases} \mu & \text{if } x=1 \\ 1-\mu & \text{if } x=0 \end{cases} \end{aligned}$$

Expectation

Given a discrete random variable X ,
we might have a function of X in mind for a task
(e.g. the cost of the fruit selected,
the points I get for each dice roll,
etc)

Let $f(x)$ map from sample space of X
to a real value or
vector of real values.

We want to know the "average" value of $f(x)$
weighting each outcome by its probability.

$$E[f(x)] = \sum_{x \in X} P(X=x) f(x)$$

When each value x is represented as a real value
or vector,
we can compute the mean of random variable X

$$E[X] = \sum_{x \in X} P(X=x) x$$

Expectations of Bernoulli random vars.

$$X \sim \text{Bern}(\mu)$$

$$E[X] = \sum_{x \in \{0,1\}} x P(X=x)$$

$$= 1 \cdot \mu + 0 \cdot (1-\mu)$$

$$= \mu$$

$$E[X^2] = \sum_{x \in \{0,1\}} x^2 P(X=x)$$

$$= 1^2 \cdot \mu + 0^2 \cdot (1-\mu)$$

$$= \mu$$

Recall variance is measure of how much, on average, a random variable's value will differ from its mean.

X is R.V.

$$\text{var}[X] = E[(X - E[X])^2] \quad \text{by definition}$$

$$= E[X^2 - 2X E[X] + E[X]^2]$$

$$= E[X^2] - 2E[X]E[X] + E[X]^2$$

$$= E[X^2] - E[X]^2$$

Thus, variance of a Bernoulli:

$$\text{var}[X] = E[X^2] - E[X]^2 = \mu - \mu^2 = \mu(1-\mu)$$

$X \sim \text{Bern}(\mu)$

Independence vs Conditional Independence

(4)

Two random variables X and Y are "independent" if the distribution of Y given knowledge of X is the same as the marginal distribution of Y

$$\Pr(Y=y | X=x) = \Pr(Y=y) \quad \text{for all } y, x$$

That is, the distribution of Y never depends on the value of X

The joint distribution of two indep. r.v.s looks like a product of marginals

$$\begin{aligned} \Pr(X=x, Y=y) &= \Pr(Y=y | X=x) \Pr(X=x) && \text{product rule} \\ &= \Pr(Y=y) \Pr(X=x) && \text{defn of independence} \end{aligned}$$

Conditional Independence

We say two rand vars X and Y are "conditionally independent" given a third random variable Z iff

$$\Pr(Y=y | X=x, Z=z) = \Pr(Y=y | Z=z) \quad \text{for all } x, y, z \text{ values}$$

again, this means the joint distr. looks like

$$\Pr(X=x, Y=y, Z=z) = \Pr(Y=y | Z=z) \Pr(X=x | Z=z) \Pr(Z=z)$$

Models for (independent) coin flips

Suppose we coin toss N times, the outcome can be represented as

$$X_1=x_1, X_2=x_2, \dots, X_N=x_N$$

In general, how many parameters are needed to describe this joint distribution?

ans: $2^N - 1$

Joint proba mass function requires 2^N table entries
but sum must be one so last entry is determined

Now, let us assume each toss is independent

$$P(X_1=x_1, \dots, X_N=x_n) = \prod_{n=1}^N P(X_n=x_n)$$

How many parameters?
 N

Further assume
identically
distributed

$$\rightarrow = \prod_{n=1}^N \text{BernPMF}(X_n=x_n/\mu)$$

How many params?
 1

Insight: We always observe N values

Assumptions like independence & identical distribution help us reduce number of parameters we need to estimate/learn

Estimation of μ from data

Given N observations, $X_1 = x_1, \dots, X_N = x_N$, of binary outcomes,
 and assuming model $\prod_{n=1}^N \text{Bern}(X_n = x_n | \mu)$,
 what is value of μ ?

What principle should we use?

A good one: maximize likelihood

"find value of μ that makes our observed data most likely / most probable under the model"

$$\hat{\mu} = \arg \max_{\mu \in [0, 1]} \prod_{n=1}^N \text{Bern}(X_n = x_n | \mu)$$

Issue:

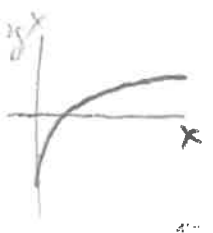
$\prod_{n=1}^N \text{Bern}(X_n / \mu)$ is a product of N numbers between 0 and 1

What can happen in a computer?
 underflow!

$0.9 \cdot 0.9 \cdot 0.9 \cdot \dots = 0.0$ using finite precision
 but should be > 0.0

Correction:

for any objective function $f(x)$, consider $\log f(x)$ as objective instead
 $\mathbb{R} \rightarrow \mathbb{R}'$



$\log(a)$ is a monotonic function, so

$$\arg \max_x f(x) = \arg \max_x \log f(x)$$

Maximizing log likelihood for ML estimation

$$\hat{\mu} = \operatorname{argmax}_{\mu \in [0,1]} \sum_{n=1}^N \log \text{Bern PMF}(X_n = x_n / \mu)$$

$$= \operatorname{argmax}_{\mu \in [0,1]} \sum_{n=1}^N x_n \log \mu + (1-x_n) \log(1-\mu)$$

$$= \operatorname{argmax}_{\mu \in [0,1]} S(x) \log \mu + r(x) \log(1-\mu)$$

$$\text{where } S(x) = \sum_n x_n \geq 0 \\ r(x) = N - S(x) \geq 0$$

This is a constrained optimization problem,

b/c valid solutions μ must be within $[0, 1]$

In general, 2 ways to solve

(1) act like there is no constraint, hope that answer works

(2) method of Lagrange multipliers (next class)

$$\operatorname{argmax}_{\mu} f(\mu), \quad f(\mu) = S \log \mu + r \log(1-\mu)$$

$$\text{set } \frac{\partial f}{\partial \mu} = 0 \text{ and solve for } \mu$$

then, verify second derivative is negative so μ^* is a maximum.

$$\frac{\partial f}{\partial \mu} = \frac{S}{\mu} - \frac{r}{1-\mu} = 0$$

$$\frac{S}{\mu} = \frac{r}{1-\mu}$$

$$S - S\mu = r\mu$$

$$S = (S+r)\mu$$

$$\mu^* = \frac{S}{S+r} = \frac{\sum_n x_n}{N}$$

does this meet our constraint?
yes!

Problems w/ ML estimation

Suppose 3 coins all heads

$$\hat{\mu}^{ML} = 1.0$$

do we believe that?

ML is vulnerable to overfitting on small training sets

Advantages of ML estimation

Can prove several general properties (we won't prove, just good to know)

(1) ML estimates are consistent

Suppose I generate data from coin flip model with μ^{true}

then, you estimate μ^{ML} from data

can prove that as $N \rightarrow \infty$, $\mu^{ML} = \mu^{true}$

(2) ML estimates are equivariant to different parametrizations

$$\text{model A: } \prod_n \text{Bern}(X_n = x_n | \mu) \quad \mu \in [0, 1]$$

$$\text{B: } \prod_n \text{Bern}(X_n = x_n | s^2) \quad s \in [0, 1] \quad \mu(s) = s^2$$

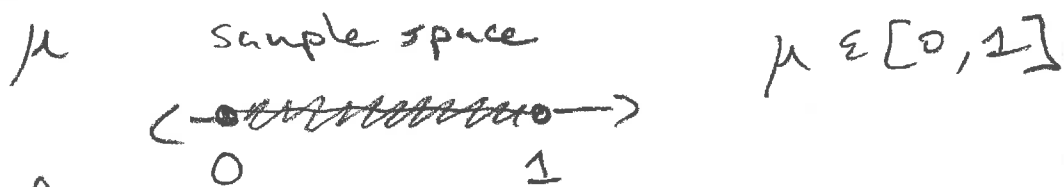
$$\text{C: } \prod_n \text{Bern}(X_n = x_n | \sigma(r)) \quad r \in \mathbb{R}, \quad \sigma(r) = \frac{e^r}{1+e^r}$$

As long as you can map parameters to one-another,
 $\mu^{ML} = \text{map}(s^{ML}) = \text{map}(r^{ML})$

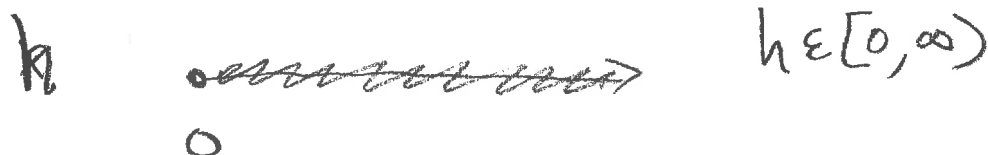
Continuous Random Variables

Examples

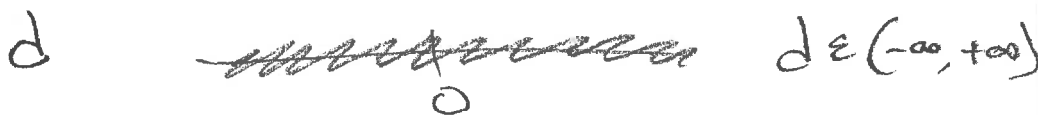
Probability coin ends up heads



Height of avg. student



Distance of random location from Medford



Sample spaces are continuous intervals of real line

Can talk coherently about ^{the probability of} interval events

what's probability that height is above 75 cm

$$P(H \geq 75)$$

Need different language for specific events

should not say what is probability height equals 75

instead, think about limits

$$\text{pdf}(H=75) = P(H \in [75, 75+\delta]) \text{ as } \delta \rightarrow 0$$

↑
probability density function of random variable H

Two rules for a PDF function for continuous r.v. X

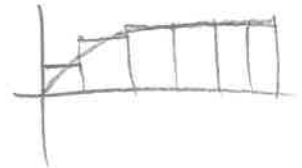
$$\text{pdf}(x) \geq 0 \quad \text{for all outcomes } x \in \mathcal{X}$$

$$\int_{\mathcal{X}} \text{pdf}(x) dx = 1$$

Think of the Riemannian sum of this integral

$$\lim_{\Delta x \rightarrow 0} \sum_{n=1}^N \text{pdf}(x_n) \Delta x_n = 1$$

↑ ↑ ↑
density of volume mass
each of
bin each bin



$$\text{density} = \frac{\text{mass}}{\text{volume}}$$

Can density be larger than one? Yes, if mass of bin is small enough

Joint probabilities also obey these rules
 $X \in \mathbb{R}, Y \in \mathbb{R}$

$$\text{pdf}(x, y) \geq 0 \quad \text{and} \quad \iint \text{pdf}(x, y) dx dy = 1$$

Sum Rule

$$\text{pdf}(x) = \int \text{pdf}(x, y) dy$$

Product Rule

$$\begin{aligned} \text{pdf}(x, y) &= \text{pdf}(x|y) \text{pdf}(y) \\ &= \text{pdf}(y|x) \text{pdf}(x) \end{aligned}$$

Beta distribution

Random variable μ

Sample space $\mu \in [0, 1]$

PDF function parameters $a > 0$
 $b > 0$

Beta pdf $(\mu|a,b)$

Gamma function: see notebook online.

$$\Gamma(a) = \int_0^{\infty} u^{a-1} e^{-u} du$$

$$\Gamma: \mathbb{R} \rightarrow \mathbb{R}$$

scalar \rightarrow scalar
real \rightarrow real

$$\frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}$$

Const. wrt μ

we'll call this
"normalizing"
constant

"interesting part
is a function"

$f(\mu|a,b)$

Remember:

$$\int_0^1 \text{BetaPDF}(\mu|a,b) d\mu = 1$$

$$\int_0^1 c(a,b) f(\mu|a,b) d\mu = 1$$

implies

$$\int_0^1 f(\mu|a,b) d\mu = \frac{1}{c(a,b)}$$

thus,

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

Now, consider a joint model for μ and the coin observations x_1, \dots, x_N

$$P(\underbrace{x_1, x_2, \dots, x_N}_{\text{discrete}}, \underbrace{\mu}_{\text{continuous}}) = \text{Beta PDF}(\mu | a, b) \prod_{n=1}^N \text{Bern}(x_n | \mu)$$

We observe x_1, \dots, x_N , how to estimate μ ?

several ways:

- (1) find μ that maximizes $P(\mu | x_1, \dots, x_N)$
called maximum a-posteriori (MAP)
- (2) take a more Bayesian approach (next time)

MAP estimator

$$\begin{aligned} \mu^* &= \arg \max_{\mu \in [0, 1]} P(\mu | x_1, \dots, x_N) \\ &= \arg \max_{\mu \in [0, 1]} \frac{P(\mu) P(x_1, \dots, x_N | \mu)}{\text{const wrt } \mu} \end{aligned}$$

$$= \arg \max_{\mu \in [0, 1]} \log \text{Beta PDF}(\mu | a, b) + \sum_n \log \text{Bern PMF}(x_n | \mu)$$

$$= \arg \max_{\mu \in [0, 1]} (a-1) \log \mu + (b-1) \log(1-\mu) + S(x) \log \mu + r(x) \log(1-\mu)$$

MAP estimator (cont'd)

$$\mu^* = \arg \max_{\mu \in [0,1]} \left(\overset{s'}{s(x)+a-1} \right) \log \mu + \left(\underset{r'}{r(x)+b-1} \right) \log (1-\mu)$$

We've solved this before!

$$\mu^* = \frac{s'}{s'+r'} \quad \text{when } \begin{matrix} s' \geq 0 \\ r' \geq 0 \end{matrix}$$

$$\mu^* = \frac{s(x)+a-1}{\underbrace{N+a+b-2}_w = s(x)+r(x)}$$

requires

$$a \geq 1$$

$$b \geq 1$$

otherwise

MAP does not exist

Interpretation

a acts like a "pseudocount" for heads

b acts like a "pseudocount" for tails