# Probability Review and Statistical Estimation

## Instructor: Alan Ritter

Many slides from Tom Mitchell, Pedro Domingos

# Random Variables

- Informally, A is a <u>random variable</u> if
  - A denotes something about which we are uncertain
  - perhaps the outcome of a randomized experiment

- Examples
  A = True if a randomly drawn person from our class is female
  A = The hometown of a randomly drawn person from our class
  A = True if two randomly drawn persons from our class have same birthday

- Define P(A) as "the fraction of possible worlds in which A is true" or "the fraction of times A holds, in repeated runs of the random experiment"
  - the set of possible worlds is called the sample space, S
  - A random variable A is a function defined over S
    $$A: S \rightarrow \{0,1\}$$

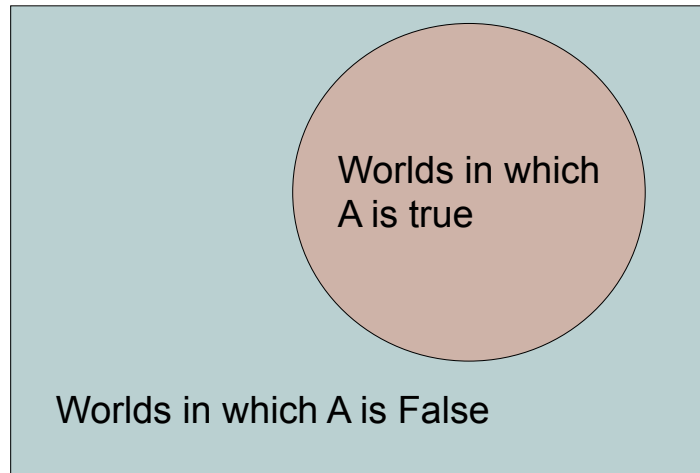# A little formalism

More formally, we have

- a <u>sample space</u> S (e.g., set of students in our class)
    - aka the set of possible worlds

- a <u>random variable</u> is a function defined over the sample space
    - Gender: S → { m, f }
    - Height: S → Reals
- an <u>event</u> is a subset of S
    - e.g., the subset of S for which Gender=f
    - e.g., the subset of S for which (Gender=m) AND (eyeColor=blue)
- we're often interested in probabilities of specific events
- and of specific events conditioned on other specific events

# Visualizing A

Sample space of all possible worlds

Its area is 1

Worlds in which A is true

Worlds in which A is False

P(A) = Area of reddish oval

# The Axioms of Probability

- 0 <= P(A) <= 1
- P(True) = 1
- P(False) = 0
- P(A or B) = P(A) + P(B) - P(A and B)
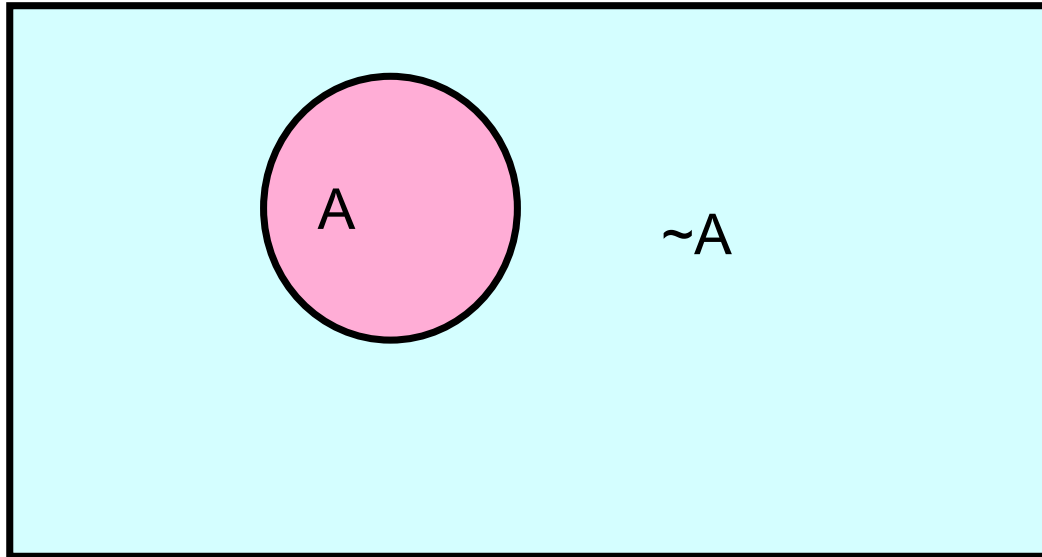
[di Finetti 1931]:

when gambling based on "uncertainty formalism A" you can be exploited by an opponent

iff

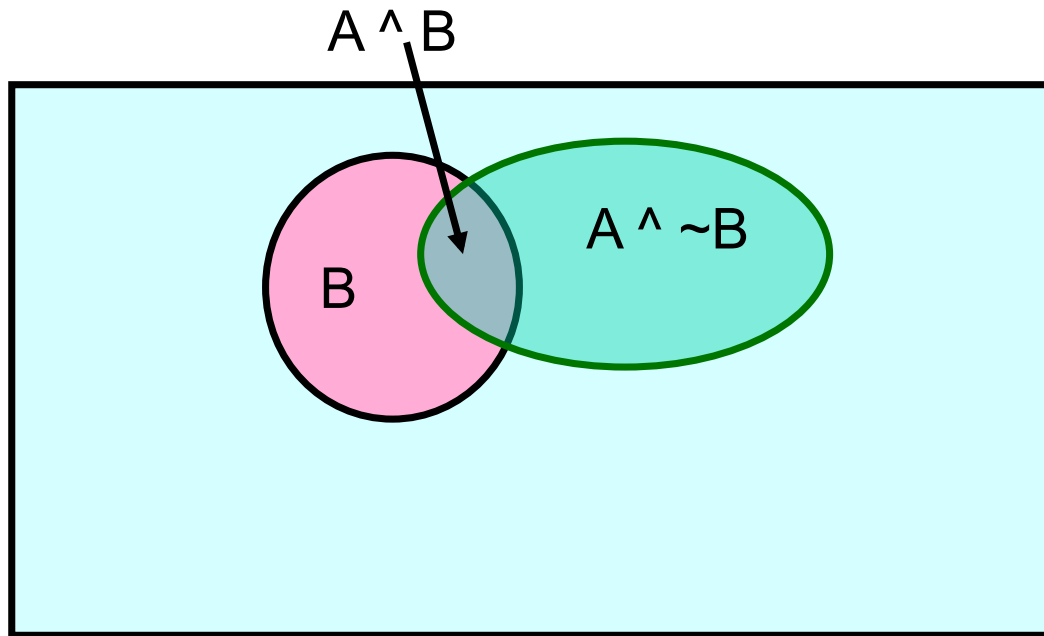your uncertainty formalism A violates these axioms
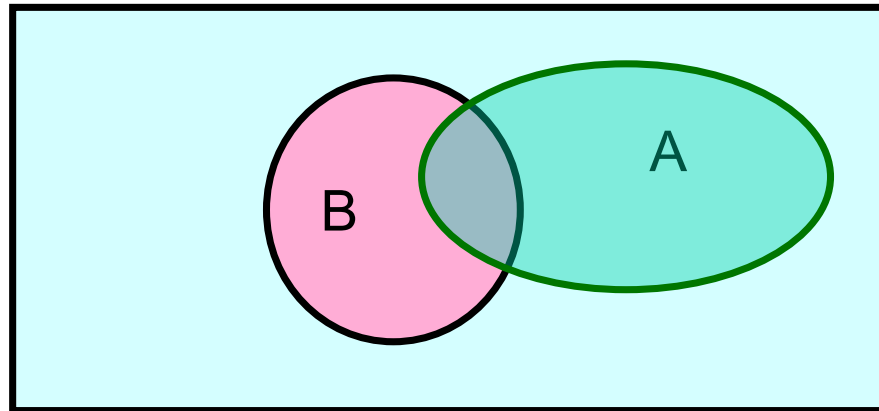
# Elementary Probability in Pictures

- P(~A) + P(A) = 1

# Elementary Probability in Pictures

- P(A) = P(A ^ B) + P(A ^ ~B)

# Definition of Conditional Probability

$$P(A|B) \ = \ \frac{P(A \wedge B)}{P(B)}$$
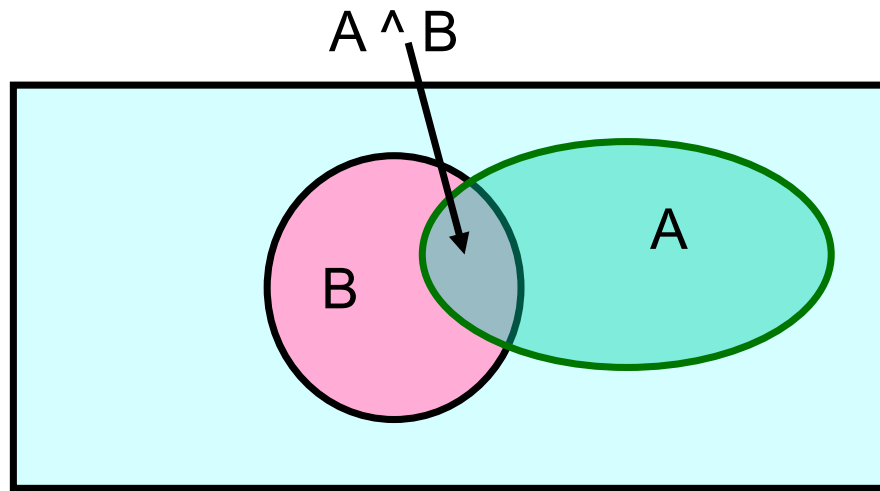
# Definition of Conditional Probability

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

## Corollary: The Chain Rule

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

# Bayes Rule

- let's write 2 expressions for P(A ^ B)

$$P(A|B) = \frac{P(B|A) * P(A)}{P(B)}$$ Bayes' rule

we call P(A) the "prior"

and P(A|B) the "posterior"

**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London,* **53:370-418**

…by no means merely a curious speculation in the doctrine of chances, but necessary to be solved in order to a sure foundation for all our reasonings concerning past facts, and what is likely to be hereafter…. necessary to be considered by any that would give a clear account of the strength of *analogical* or *inductive reasoning…*

# Other Forms of Bayes Rule

$$P(A|B) = \frac{P(B \mid A)P(A)}{P(B \mid A)P(A) + P(B \mid \sim A)P(\sim A)}$$

$$P(A|B \wedge X) = \frac{P(B \mid A \wedge X)P(A \wedge X)}{P(B \wedge X)}$$

# Applying Bayes Rule

$$P(A \mid B) = \frac{P(B \mid A)P(A)}{P(B \mid A)P(A) + P(B \mid \sim A)P(\sim A)}$$

A = you have the flu,   B = you just coughed

Assume:
P(A) = 0.05
P(B|A) = 0.80
P(B| ~A) = 0.2

what is P(flu | cough) = P(A|B)?

what does all this have to do with function approximation?

# The Joint Distribution

Recipe for making a joint distribution of M variables:

# The Joint Distribution

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have $2^M$ rows).

| A | B | C |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 0 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |
| 1 | 1 | 1 |

# The Joint Distribution

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have $2^M$ rows).
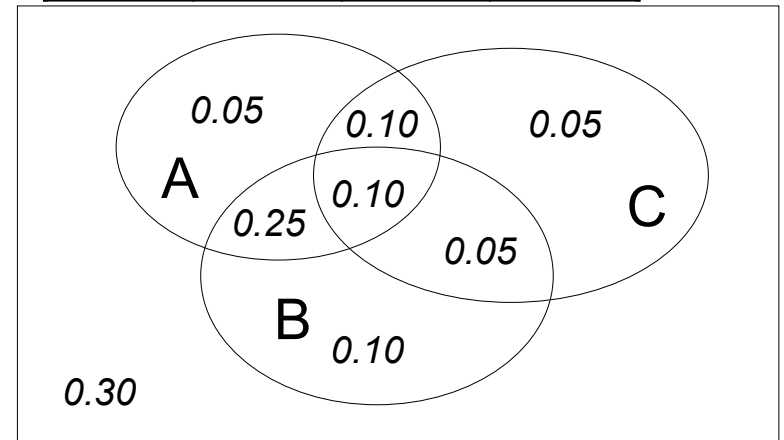2. For each combination of values, say how probable it is.

| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | 0.10 |

# The Joint Distribution

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have $2^M$ rows).
2. For each combination of values, say how probable it is.
3. If you subscribe to the axioms of probability, those numbers must sum to 1.

| A | B | C | Prob |
|---|---|---|------|
| 0 | 0 | 0 | 0.30 |
| 0 | 0 | 1 | 0.05 |
| 0 | 1 | 0 | 0.10 |
| 0 | 1 | 1 | 0.05 |
| 1 | 0 | 0 | 0.05 |
| 1 | 0 | 1 | 0.10 |
| 1 | 1 | 0 | 0.25 |
| 1 | 1 | 1 | 0.10 |

# Using the Joint Distribution

| gender | hours_worked | wealth | | |
|--------|--------------|--------|--------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

One you have the JD you can ask for the probability of **any** logical expression involving these variables

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Using the Joint

| gender | hours_worked | wealth | | |
|--------|-------------|--------|-----------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

P(Poor Male) = 0.4654

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Using the Joint

| gender | hours_worked | wealth | | |
|--------|--------------|--------|---|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

P(Poor) = 0.7604

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

[A. Moore]

# Inference with the Joint

| gender | hours_worked | wealth | | |
|--------|--------------|--------|--------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

$$P(E_1 \mid E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\displaystyle\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\displaystyle\sum_{\text{rows matching } E_2} P(\text{row})}$$

P(Male | Poor) = 0.4654 / 0.7604 = 0.612

[A. Moore]

# Learning and the Joint Distribution



| gender | hours_worked | wealth | | |
|--------|--------------|--------|--------|---|
| Female | v0:40.5- | poor | 0.253122 | |
| | | rich | 0.0245895 | |
| | v1:40.5+ | poor | 0.0421768 | |
| | | rich | 0.0116293 | |
| Male | v0:40.5- | poor | 0.331313 | |
| | | rich | 0.0971295 | |
| | v1:40.5+ | poor | 0.134106 | |
| | | rich | 0.105933 | |

Suppose we want to learn the function f: <G, H> → W

Equivalently, P(W | G, H)

Solution: learn joint distribution from data, calculate P(W | G, H)

e.g., P(W=rich | G = female, H = 40.5- ) =

[A. Moore]

sounds like the solution to
learning F: X $\rightarrow$ Y,
or P(Y | X).


Are we done?

sounds like the solution to
learning F: X $\rightarrow$ Y,
or P(Y | X).

Main problem: learning P(Y|X)
can require more data than we have

consider learning Joint Dist. with 100 attributes
# of rows in this table?
# of people on earth?
fraction of rows with 0 training examples?

# What to do?

1. Be smart about how we estimate probabilities from sparse data
    – maximum likelihood estimates
    – maximum a posteriori estimates

2. Be smart about how to represent joint distributions
    – Bayes networks, graphical models

# 1. Be smart about how we estimate probabilities

# Estimating Probability of Heads

X=1    X=0

- I show you the above coin $X$, and hire you to estimate the probability that it will turn up heads $(X = 1)$ or tails $(X = 0)$

- You flip it repeatedly, observing
  - it turns up heads $\alpha_1$ times
  - it turns up tails $\alpha_0$ times

- Your estimate for $P(X = 1)$ is....?

# Estimating θ = P(X=1)

X=1    X=0

Test A:

  100 flips: 51 Heads (X=1), 49 Tails (X=0)

Test B:

  3 flips:  2 Heads (X=1), 1 Tails (X=0)

# Estimating θ = P(X=1)

X=1    X=0

Case C: (online learning)

- keep flipping, want single learning algorithm that gives reasonable estimate after each flip

# Principles for Estimating Probabilities

Principle 1 (maximum likelihood):

- choose parameters θ that maximize **P(data | θ)**

- e.g.,
$$\hat{\theta}^{MLE} = \frac{\alpha_1}{\alpha_1 + \alpha_0}$$

Principle 2 (maximum a posteriori prob.):

- choose parameters θ that maximize **P(θ | data)**

- e.g.

$$\hat{\theta}^{MAP} = \frac{\alpha_1 + \#\text{hallucinated\_1s}}{(\alpha_1 + \#\text{hallucinated\_1s}) + (\alpha_0 + \#\text{hallucinated\_0s})}$$

# Maximum Likelihood Estimation

P(X=1) = θ          P(X=0) = (1-θ)

X=1     X=0

Data D:

Flips produce data D with $\alpha_1$ heads, $\alpha_0$ tails
- flips are independent, identically distributed 1's and 0's (Bernoulli)
- $\alpha_1$ and $\alpha_0$ are counts that sum these outcomes (Binomial)

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1}(1-\theta)^{\alpha_0}$$

# Maximum Likelihood Estimate for $\Theta$

$$\hat{\theta} = \arg\max_{\theta} \ln P(\mathcal{D} \mid \theta)$$

$$= \arg\max_{\theta} \ln \theta^{\alpha_H}(1-\theta)^{\alpha_T}$$

- Set derivative to zero: $\dfrac{d}{d\theta} \ln P(\mathcal{D} \mid \theta) = 0$

[C. Guestrin]

$$\hat{\theta} = \arg\max_{\theta} \ \ln P(D|\theta)$$

$$= \arg\max_{\theta} \ \ln \left[ \theta^{\alpha_1}(1-\theta)^{\alpha_0} \right]$$

- Set derivative to zero: $\frac{d}{d\theta} \ln P(\mathcal{D}\mid\theta) = 0$

hint: $\frac{\partial \ln\theta}{\partial\theta} = \frac{1}{\theta}$

# Summary:
## Maximum Likelihood Estimate

X=1    X=0

$P(X=1) = \theta$
$P(X=0) = 1-\theta$
(Bernoulli)

- Each flip yields boolean value for $X$

$$X \sim \text{Bernoulli}: P(X) = \theta^X (1 - \theta)^{(1-X)}$$

- Data set $D$ of independent, identically distributed (iid) flips produces $\alpha_1$ ones, $\alpha_0$ zeros (Binomial)

$$P(D|\theta) = P(\alpha_1, \alpha_0|\theta) = \theta^{\alpha_1}(1 - \theta)^{\alpha_0}$$

$$\hat{\theta}^{MLE} = \text{argmax}_\theta\, P(D|\theta) = \frac{\alpha_1}{\alpha_1+\alpha_0}$$

# Principles for Estimating Probabilities

Principle 1 (maximum likelihood):

- choose parameters θ that maximize P(data | θ)
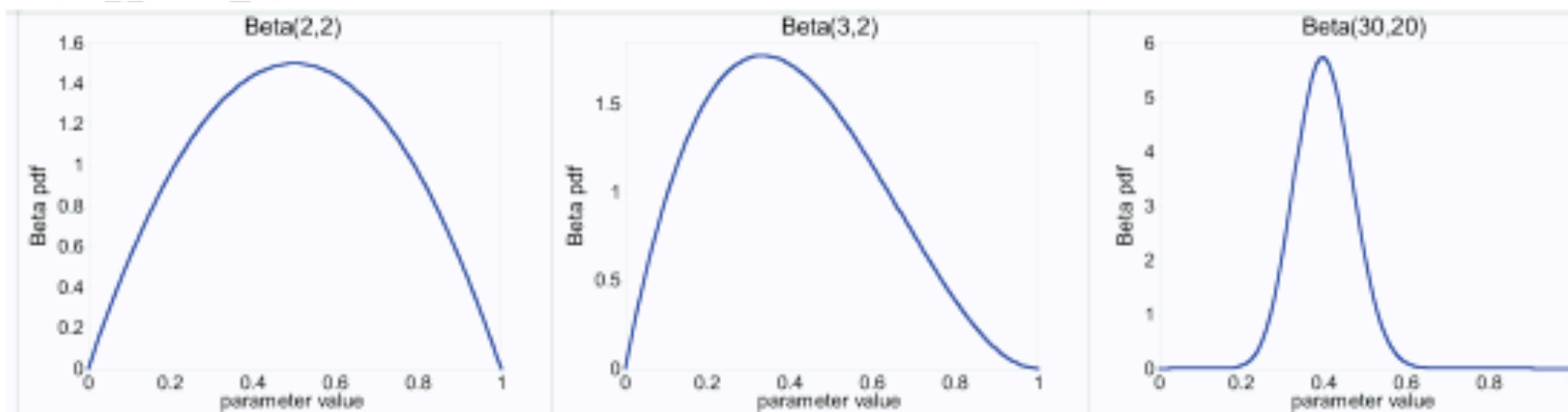
Principle 2 (maximum a posteriori prob.):

- choose parameters θ that maximize

$$P(\theta \mid data) = \frac{P(data \mid \theta)\ P(\theta)}{P(data)}$$

# Beta prior distribution – P(θ)

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1-\theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

- **Likelihood function:** $P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1-\theta)^{\alpha_T}$
- **Posterior:** $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

# Beta prior distribution – P(θ)

$$P(\theta) = \frac{\theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$



[C. Guestrin]

Eg. 1 Coin flip problem

Likelihood is ~ Binomial

$$P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1-\theta)^{\alpha_T}$$

If prior is Beta distribution,

$$P(\theta) = \frac{\theta^{\beta_H-1}(1-\theta)^{\beta_T-1}}{B(\beta_H, \beta_T)} \sim Beta(\beta_H, \beta_T)$$

Then posterior is Beta distribution

$$P(\theta|D) \sim Beta(\alpha_H + \beta_H, \alpha_H + \beta_H)$$

and MAP estimate is therefore

$$\hat{\theta}^{MAP} = \frac{\alpha_H + \beta_H - 1}{(\alpha_H + \beta_H - 1) + (\alpha_T + \beta_T - 1)}$$

**Eg. 2** Dice roll problem (6 outcomes instead of 2)

Likelihood is ~ Multinomial($\theta = \{\theta_1, \theta_2, \ldots, \theta_k\}$)

$$P(\mathcal{D} \mid \theta) = \theta_1^{\alpha_1} \theta_2^{\alpha_2} \ldots \theta_k^{\alpha_k}$$

If prior is Dirichlet distribution,

$$P(\theta) = \frac{\theta_1^{\beta_1-1} \; \theta_2^{\beta_2-1} \ldots \theta_k^{\beta_k-1}}{B(\beta_1, \ldots, \beta_k)} \sim \mathrm{Dirichlet}(\beta_1, \ldots, \beta_k)$$

Then posterior is Dirichlet distribution

$$P(\theta|D) \sim \mathrm{Dirichlet}(\beta_1 + \alpha_1, \ldots, \beta_k + \alpha_k)$$

and MAP estimate is therefore

$$\hat{\theta}_i^{MAP} = \frac{\alpha_i + \beta_i - 1}{\sum_{j=1}^{k}(\alpha_j + \beta_j - 1)}$$

# Some terminology

- Likelihood function:  P(data | θ)
- Prior: P(θ)
- Posterior: P(θ | data)

- Conjugate prior: P(θ) is the conjugate prior for likelihood function P(data | θ) if the forms of P(θ) and P(θ | data) are the same.

# You should know

- Probability basics
    - random variables, conditional probs, …
    - Bayes rule
    - Joint probability distributions
    - calculating probabilities from the joint distribution
- Estimating parameters from data
    - maximum likelihood estimates
    - maximum a posteriori estimates
    - distributions – binomial, Beta, Dirichlet, …
    - conjugate priors