# Introduction to Bayesian methods
# Lecture 17

David Sontag

New York University
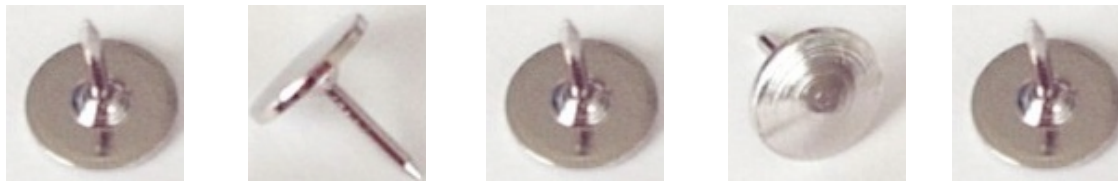
# Bayesian learning

- Bayesian learning uses **probability** to *model* data and *quantify uncertainty* of predictions
  - Facilitates incorporation of prior knowledge
  - Gives optimal predictions
    - Allows for decision-theoretic reasoning

# Your first consulting job

- A billionaire from the suburbs of Manhattan asks you a question:
  - He says: I have thumbtack, if I flip it, what's the probability it will fall with the nail up?
  - You say: Please flip it a few times:

  

  - You say: The probability is:
    - P(heads) = 3/5
  - **He says: Why???**
  - You say: Because…

# Random Variables

- A random variable is some aspect of the world about which we (may) have uncertainty
  - R = Is it raining?
  - D = How long will it take to drive to work?
  - L = Where am I?

- We denote random variables with capital letters

- Random variables have domains
  - R in {true, false}   (sometimes write as {+r, ¬r})
  - D in [0, ∞)
  - L in possible locations, maybe {(0,0), (0,1), …}

# Probability Distributions

- Discrete random variables have distributions

$$P(T)$$

| T | P |
|------|-----|
| warm | 0.5 |
| cold | 0.5 |

$$P(W)$$

| W | P |
|--------|-----|
| sun | 0.6 |
| rain | 0.1 |
| fog | 0.3 |
| meteor | 0.0 |

- A discrete distribution is a TABLE of probabilities of values
- The probability of a state (lower case) is a single number

$$P(W = rain) = 0.1 \qquad P(rain) = 0.1$$

- Must have:

$$\forall x \, P(x) \geq 0 \qquad \sum_x P(x) = 1$$

# Joint Distributions

- A *joint distribution* over a set of random variables: $X_1, X_2, \ldots X_n$ specifies a real number for each assignment:

$$P(X_1 = x_1, X_2 = x_2, \ldots X_n = x_n)$$

$$P(x_1, x_2, \ldots x_n)$$

$P(T, W)$

| T | W | P |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

  - How many assignments if *n* variables with domain sizes *d*?

  - Must obey:

$$P(x_1, x_2, \ldots x_n) \geq 0$$

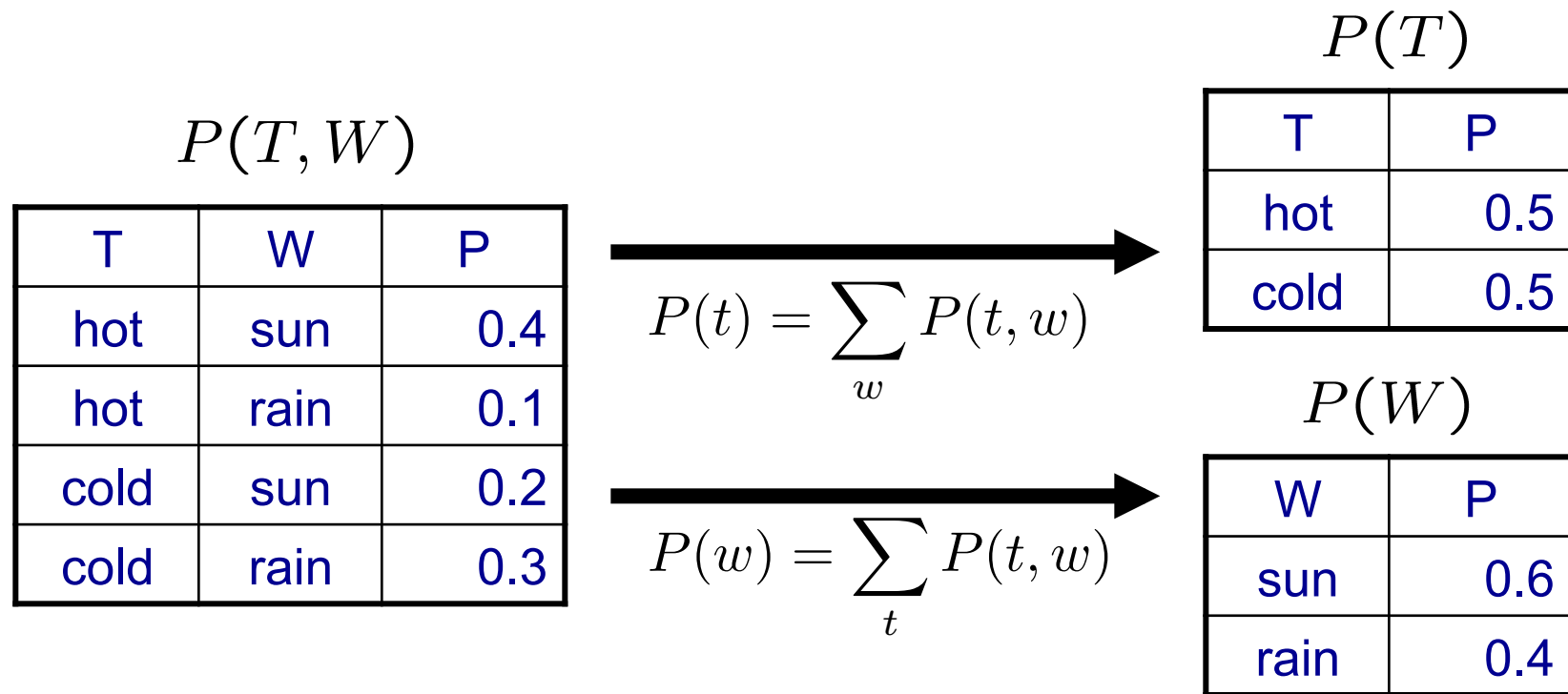$$\sum_{(x_1, x_2, \ldots x_n)} P(x_1, x_2, \ldots x_n) = 1$$

- For all but the smallest distributions, impractical to write out or estimate
  - Instead, we make additional assumptions about the distribution

# Marginal Distributions

- Marginal distributions are sub-tables which eliminate variables
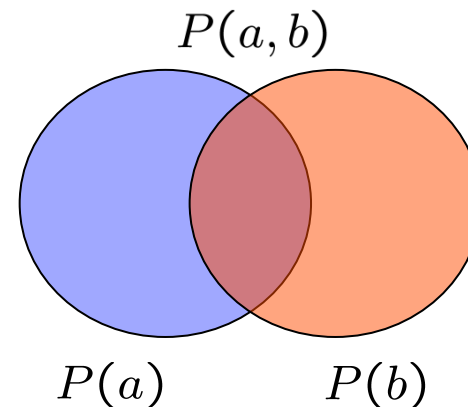- Marginalization (summing out): Combine collapsed rows by adding

$P(T, W)$

| T | W | P |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

$$P(t) = \sum_{w} P(t, w)$$

$P(T)$

| T | P |
|------|-----|
| hot | 0.5 |
| cold | 0.5 |

$$P(w) = \sum_{t} P(t, w)$$

$P(W)$

| W | P |
|------|-----|
| sun | 0.6 |
| rain | 0.4 |

$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2)$$

# Conditional Probabilities

- A simple relation between joint and conditional probabilities
  - In fact, this is taken as the *definition* of a conditional probability

$$P(a|b) = \frac{P(a,b)}{P(b)}$$



$$P(T,W)$$

| T | W | P |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

$$P(W = r | T = c) = ???$$

# Conditional Distributions

- Conditional distributions are probability distributions over some variables given fixed values of others

Conditional Distributions

Joint Distribution

$P(W|T = hot)$

| W | P |
|------|-----|
| sun | 0.8 |
| rain | 0.2 |

$P(W|T = cold)$

| W | P |
|------|-----|
| sun | 0.4 |
| rain | 0.6 |

$P(W|T)$

$P(T, W)$

| T | W | P |
|------|------|-----|
| hot | sun | 0.4 |
| hot | rain | 0.1 |
| cold | sun | 0.2 |
| cold | rain | 0.3 |

# The Product Rule

- Sometimes have conditional distributions but want the joint

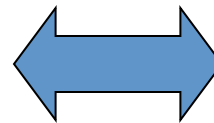$$P(x|y) = \frac{P(x,y)}{P(y)} \quad \Longleftrightarrow \quad P(x,y) = P(x|y)P(y)$$

- Example:

$P(W)$

| W | P |
|------|-----|
| sun | 0.8 |
| rain | 0.2 |

$P(D|W)$

| D | W | P |
|-----|------|-----|
| wet | sun | 0.1 |
| dry | sun | 0.9 |
| wet | rain | 0.7 |
| dry | rain | 0.3 |

$P(D,W)$

| D | W | P |
|-----|------|------|
| wet | sun | 0.08 |
| dry | sun | 0.72 |
| wet | rain | 0.14 |
| dry | rain | 0.06 |

# Bayes' Rule

- Two ways to factor a joint distribution over two variables:

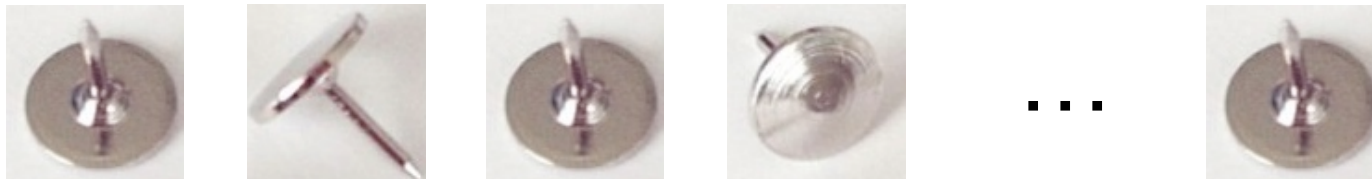$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

- Dividing, we get:

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

- Why is this at all helpful?
  - Let's us build one conditional from its reverse
  - Often one conditional is tricky but the other one is simple
  - Foundation of many practical systems (e.g. ASR, MT)

- In the running for most important ML equation!

# Returning to thumbtack example...

- P(Heads) = θ,  P(Tails) = 1-θ



- Flips are *i.i.d.*:  $D=\{x_i \mid i=1\ldots n\},\ P(D \mid \theta) = \Pi_i P(x_i \mid \theta)$

  – Independent events

  – Identically distributed according to Bernoulli distribution

- Sequence *D* of $\alpha_H$ Heads and $\alpha_T$ Tails

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

Called the "likelihood" of the data under the model

# Maximum Likelihood Estimation

- **Data:** Observed set $D$ of $\alpha_H$ Heads and $\alpha_T$ Tails

- **Hypothesis:** Bernoulli distribution

- **Learning:** finding $\theta$ is an optimization problem
  - What's the objective function?

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

- **MLE:** Choose $\theta$ to maximize probability of $D$

$$\widehat{\theta} = \arg\max_{\theta} \; P(\mathcal{D} \mid \theta)$$

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

# Your first parameter learning algorithm

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

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

- Set derivative to zero, and solve!

$$\frac{d}{d\theta} \ln P(\mathcal{D} \mid \theta) = \frac{d}{d\theta} \left[\ln \theta^{\alpha_H}(1-\theta)^{\alpha_T}\right]$$

$$= \frac{d}{d\theta} \left[\alpha_H \ln \theta + \alpha_T \ln(1-\theta)\right]$$

$$= \alpha_H \frac{d}{d\theta} \ln \theta + \alpha_T \frac{d}{d\theta} \ln(1-\theta)$$

$$= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} = 0 \qquad \boxed{\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}}$$
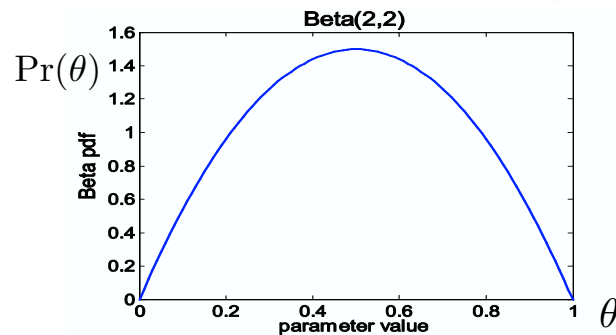
**Data**



$$L(\theta; \mathcal{D}) = \ln P(\mathcal{D}|\theta)$$



$$\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

$L(\theta{:}\mathcal{D})$

$\theta$

# What if I have prior beliefs?

- Billionaire says: Wait, I know that the thumbtack is "close" to 50-50. What can you do for me now?

- **You say: I can learn it the Bayesian way…**

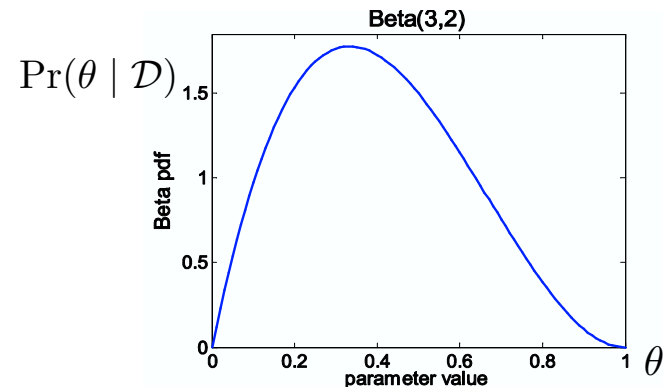- Rather than estimating a single $\theta$, we obtain a distribution over possible values of $\theta$

In the beginning

After observations



Observe flips
e.g.: {tails, tails}

# Bayesian Learning

- **Use Bayes' rule!**

Data Likelihood

Prior



Posterior



$$P(\theta \mid \mathcal{D}) = \frac{P(\mathcal{D} \mid \theta)P(\theta)}{P(\mathcal{D})}$$

Normalization

- Or equivalently: $P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$

- For *uniform* priors, this reduces to maximum likelihood estimation!

$$P(\theta) \propto 1 \qquad P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)$$

# Bayesian Learning for Thumbtacks

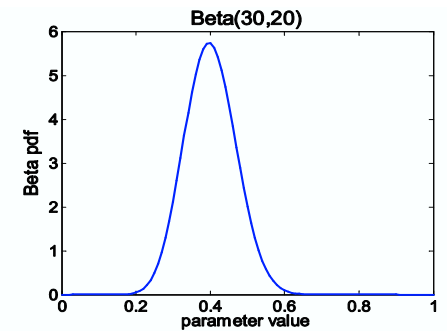$$P(\theta \mid \mathcal{D}) \propto P(\mathcal{D} \mid \theta)P(\theta)$$
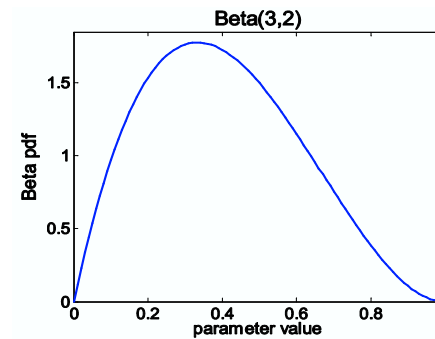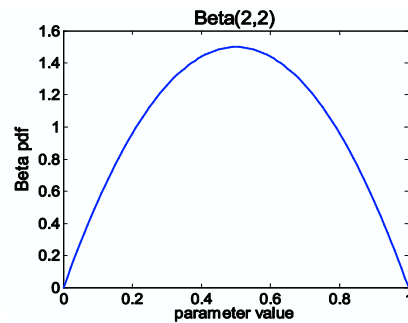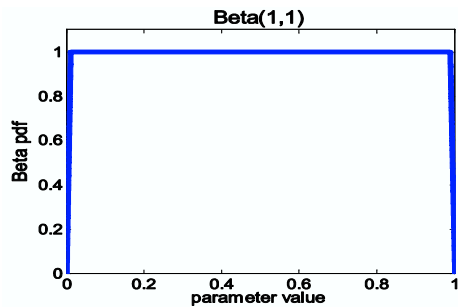
Likelihood: $P(\mathcal{D} \mid \theta) = \theta^{\alpha_H}(1 - \theta)^{\alpha_T}$

- What should the prior be?
  - Represent expert knowledge
  - Simple posterior form

- For binary variables, commonly used prior is the Beta distribution:

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

# Beta prior distribution – P(θ)

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



- The posterior distribution:

$$P(\theta \mid \mathcal{D}) \quad \propto \quad P(\mathcal{D} \mid \theta)P(\theta)$$

$$\propto \theta^{\alpha_H}(1 - \theta)^{\alpha_T} \ \theta^{\beta_H - 1}(1 - \theta)^{\beta_T - 1}$$

$$= \theta^{\alpha_H + \beta_H - 1}(1 - \theta)^{\alpha_T + \beta_t + 1}$$

$$= Beta(\alpha_H + \beta_H, \alpha_T + \beta_T)$$