# Point Estimation

# Your first consulting job

Billionaire in Dallas asks:

- 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:
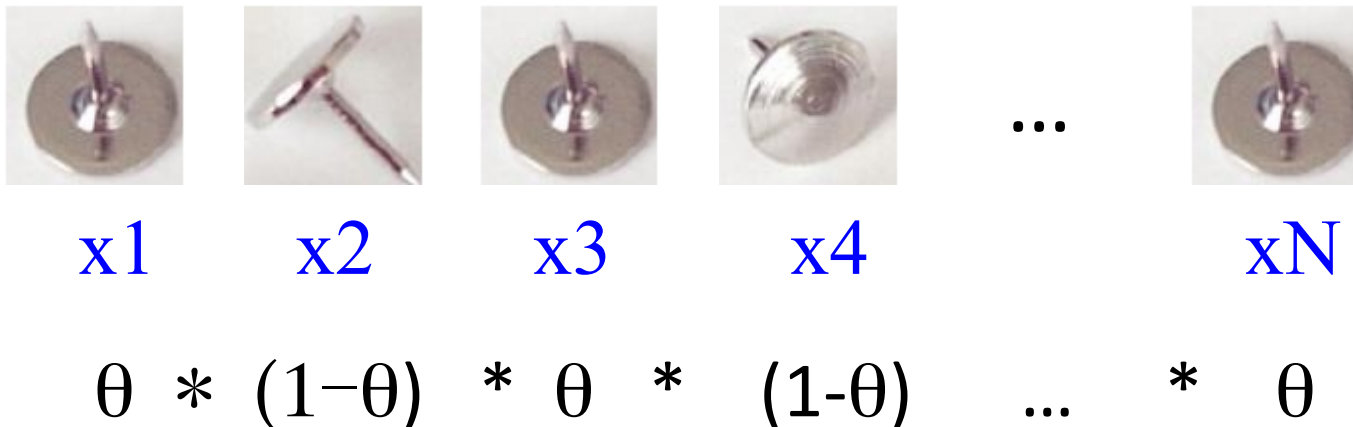
# Your first consulting job

Billionaire in Dallas asks:

– 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(H) = 3/5

–**He says: Why???**
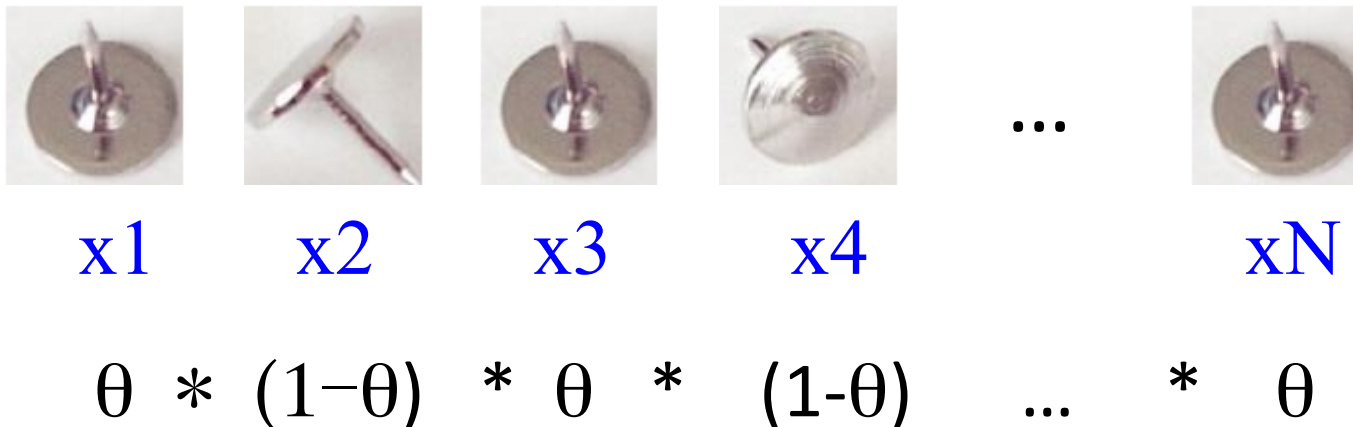
– You say: Because…

# Thumbtack – Binomial Distribution

- P(Heads) = θ,  P(Tails) = 1-θ



x1     x2     x3     x4          xN

θ  *  (1−θ)  *  θ  *  (1-θ)    ...    *    θ

- Flips are *i.i.d.*:
  - Independent events
  - Identically distributed according to Binomial distribution

# Thumbtack – Binomial Distribution

- P(Heads) = θ,  P(Tails) = 1-θ



$$x1 \quad x2 \quad x3 \quad x4 \quad\quad\quad xN$$

$$\theta \ * \ (1-\theta) \ * \ \theta \ * \ (1-\theta) \ \ldots \ * \ \theta$$

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

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

# Maximum Likelihood Estimation

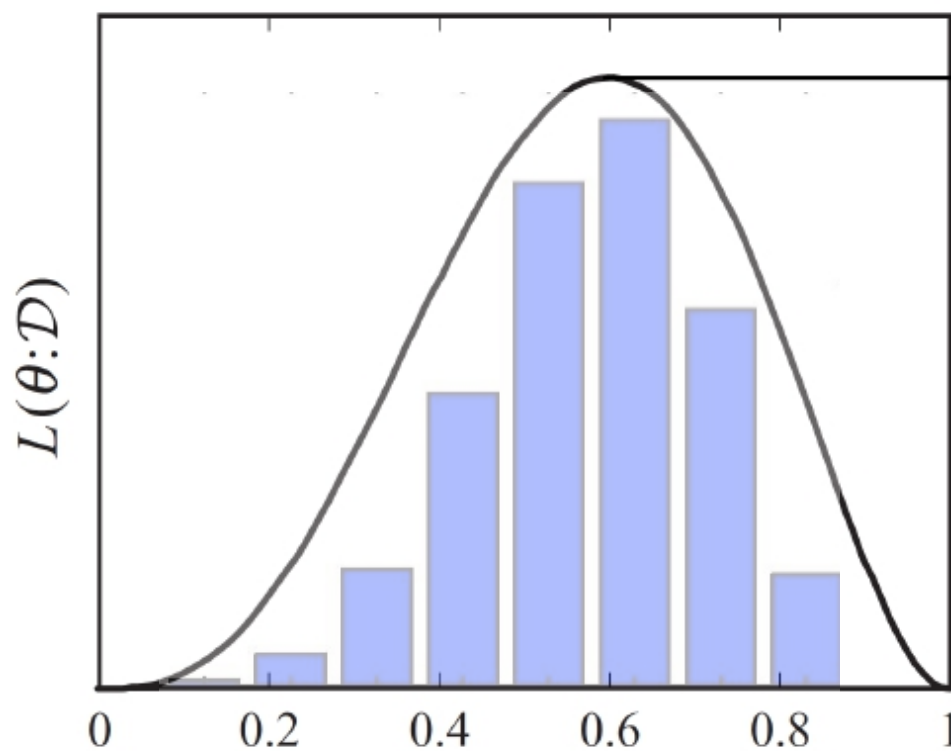- **Data:** Observed set $D$ of $\alpha_H$ Heads and $\alpha_T$ Tails
- **Hypothesis:** Binomial distribution
- **Learning:** finding $\theta$ is an optimization problem
  - What's the objective function?

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

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

$$
\begin{aligned}
\widehat{\theta} &= \arg\max_{\theta} \; P(\mathcal{D} \mid \theta) \\
&= \arg\max_{\theta} \; \ln P(\mathcal{D} \mid \theta)
\end{aligned}
$$

**Data**



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

Set the derivative to 0 and solve

# 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!

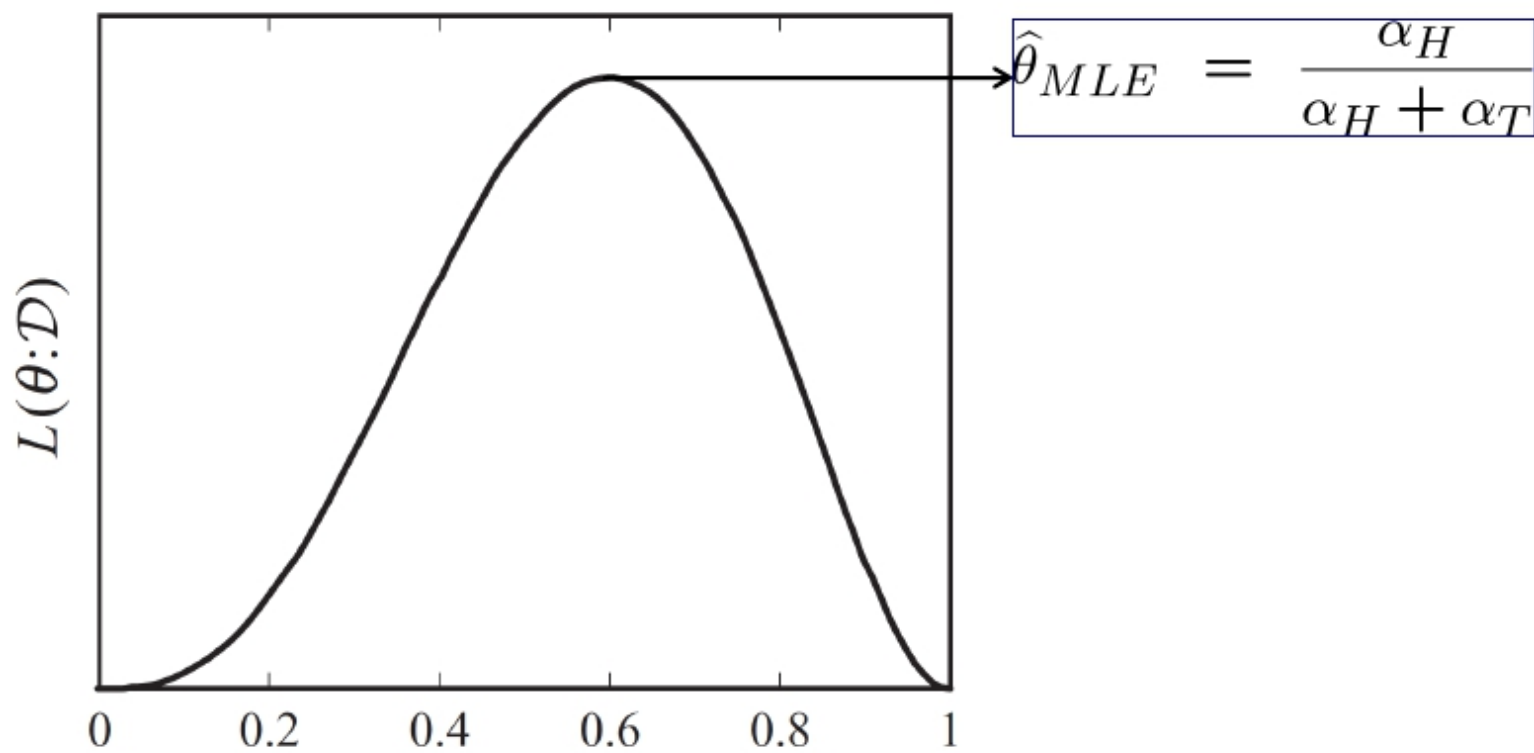$$\frac{d}{d\theta} \ln P(\mathcal{D} \mid \theta) = \frac{d}{d\theta} [\ln \theta^{\alpha_H}(1-\theta)^{\alpha_T}]$$

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

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

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

**Data**



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

# But, how many flips do I need?

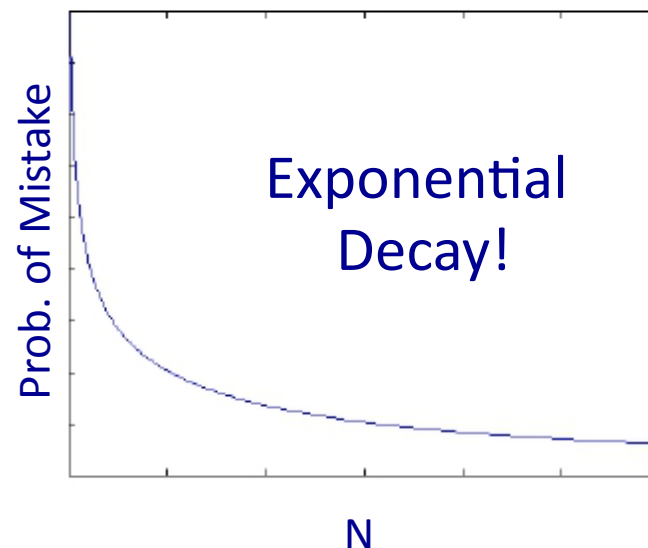$$\hat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}$$

- Billionaire says: I flipped 3 heads and 2 tails.
- You say: $\theta$ = 3/5, I can prove it!
- He says: What if I flipped 30 heads and 20 tails?
- You say: Same answer, I can prove it!
- **He says: What's better?**
- You say: Umm... The more the merrier???
- He says: Is this why I am paying you the big bucks???
- You say: I will give you a theoretical bound.

# A bound  (from Hoeffding's inequality)

For $N = \alpha_H + \alpha_T$, and $\quad \hat{\theta}_{MLE} = \dfrac{\alpha_H}{\alpha_H + \alpha_T}$

Let $\theta_*$ be the true parameter, for any $\varepsilon > 0$:

$$P(\,|\,\hat{\theta} - \theta^*\,| \geq \epsilon) \;\leq\; 2e^{-2N\epsilon^2}$$



Prob. of Mistake

Exponential Decay!

N

# PAC Learning

- **PAC:** Probably Approximately Correct
- **Billionaire says:** I want to know the thumbtack $\theta$, within $\varepsilon$ = 0.1, with probability at least 1-$\delta$ = 0.95.
- **How many flips?** Or, how big do I set $N$?

$$P(|\ \widehat{\theta} - \theta^* | \geq \epsilon) \ \leq \ 2e^{-2N\epsilon^2}$$

$$\delta \geq 2e^{-2N\epsilon^2} \geq P(\text{mistake})$$

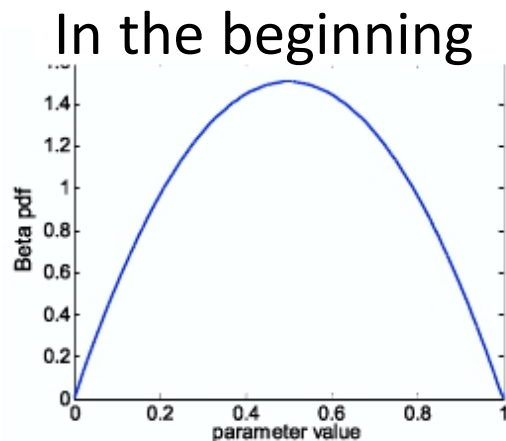$$\ln \delta \geq \ln 2 - 2N\epsilon^2$$

$$N \geq \frac{\ln(2/\delta)}{2\epsilon^2}$$

*Interesting! Lets look at some numbers!*

$\varepsilon$ = 0.1, $\delta$=0.05

$$N \geq \frac{\ln(2/0.05)}{2 \times 0.1^2} \approx \frac{3.8}{0.02} = 190$$

# 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...**

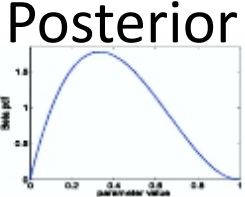- Rather than estimating a single $\theta$, we obtain a distribution over possible values of $\theta$

In the beginning

After observations

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

# Bayesian Learning

Use Bayes rule:

Data Likelihood

Prior



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

Posterior



Normalization

Or equivalently:

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

# Bayesian Learning for Thumbtacks

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

Likelihood function is Binomial:

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

- What about prior?
  - Represent expert knowledge
  - Simple posterior form
- Conjugate priors:
  - Closed-form representation of posterior
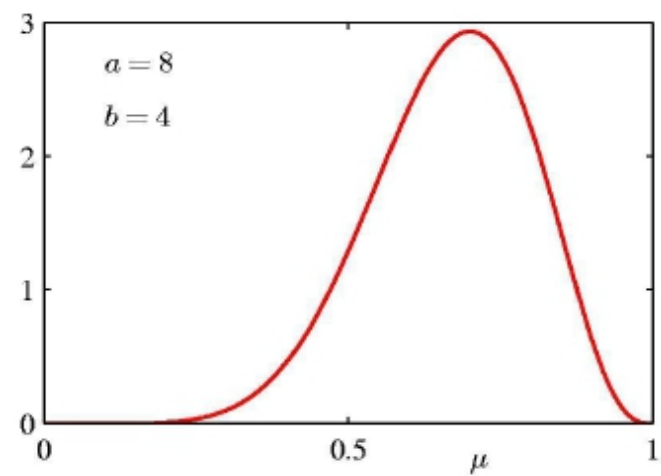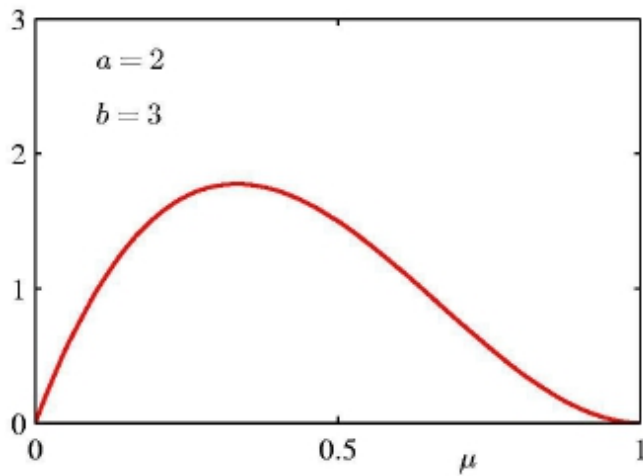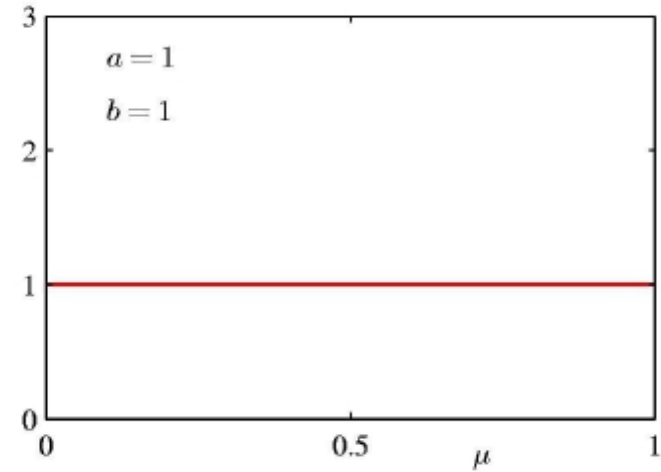  - **For Binomial, conjugate prior is Beta distribution**
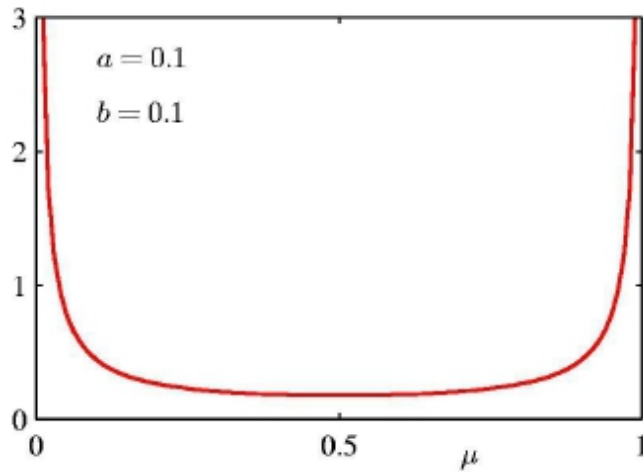
# Beta Distribution

- Distribution over $\mu \in [0, 1]$.

$$B(a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$$

$$
\begin{aligned}
\text{Beta}(\mu|a, b) &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1}(1-\mu)^{b-1} \\
\mathbb{E}[\mu] &= \frac{a}{a+b} \\
\text{var}[\mu] &= \frac{ab}{(a+b)^2(a+b+1)}
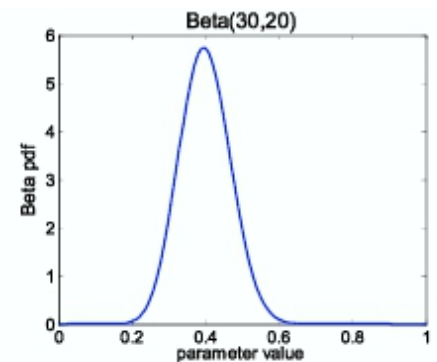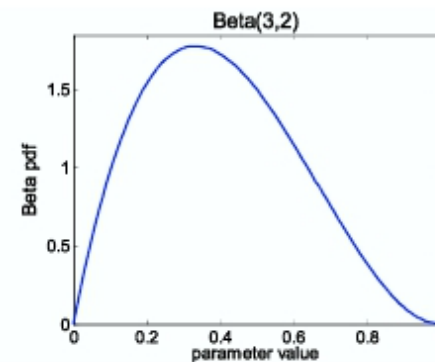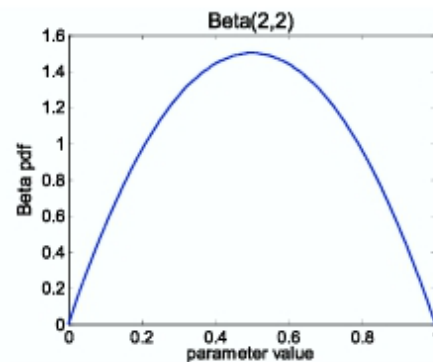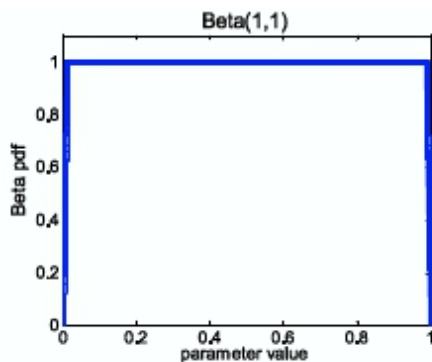\end{aligned}
$$

# Beta Distribution

# Posterior Distribution

- Prior: $Beta(\beta_H, \beta_T)$

- Data: $\alpha_H$ heads and $\alpha_T$ tails

- Posterior distribution:

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

# Bayesian Posterior Inference



Beta(30,20)

- Posterior distribution:
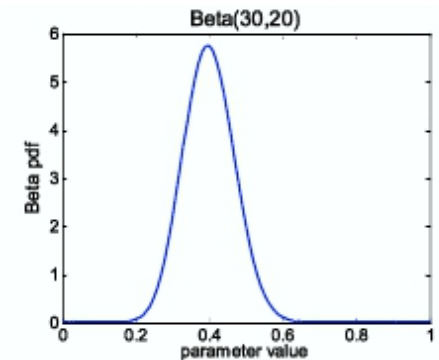
$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- Bayesian inference:
  - No longer single parameter
  - For any specific $f$, the function of interest
  - Compute the expected value of $f$

$$E[f(\theta)] = \int_0^1 f(\theta)P(\theta \mid \mathcal{D})d\theta$$

  - Integral is often hard to compute
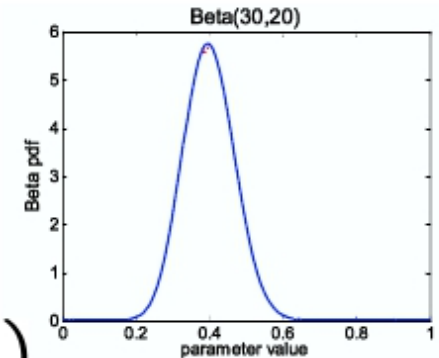
# MAP: Maximum a Posteriori Approximation


Beta(30,20)

$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

$$E[f(\theta)] = \int_0^1 f(\theta)P(\theta \mid \mathcal{D})d\theta$$

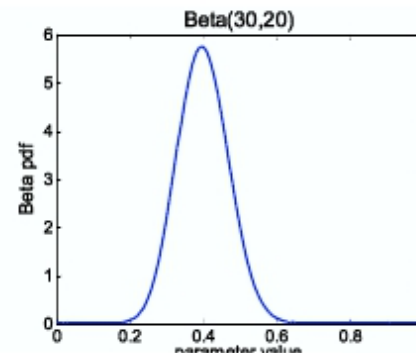- As more data is observed, Beta is more certain

- MAP: use most likely parameter to approximate the expectation

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

$$E[f(\theta)] \approx f(\widehat{\theta})$$

# MAP for Beta distribution


Beta(30,20)

$$P(\theta \mid \mathcal{D}) = \frac{\theta^{\beta_H + \alpha_H - 1}(1 - \theta)^{\beta_T + \alpha_T - 1}}{B(\beta_H + \alpha_H, \beta_T + \alpha_T)} \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

MAP: use most likely parameter:

$$\widehat{\theta} = \arg \max_{\theta} P(\theta \mid \mathcal{D}) = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

Beta prior equivalent to extra thumbtack flips

As $N \rightarrow \infty$, prior is "forgotten"

**But, for small sample size, prior is important!**

# Estimating Parameters

- Maximum Likelihood Estimate (MLE): choose
  θ that maximizes probability of observed data $\mathcal{D}$

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

- Maximum a Posteriori (MAP) estimate:
  choose θ that is most probable given prior
  probability and the data

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

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