# Point Estimation

Vibhav Gogate

The University of Texas at Dallas

Some slides courtesy of Carlos Guestrin, Chris Bishop, Dan Weld and Luke Zettlemoyer.

# Basics: Expectation and Variance

Random variable $x$ has domain $D(x)$.
Example: $x$ has domain: $\{1, 2, 3, 4\}$
The distribution $P$ is defined over $D(x)$.

$$\mathbb{E}_P[x] = \sum_{x \in D(x)} x P(x)$$

$$\text{var}_P[x] = \sum_{x \in D(x)} (x - \mathbb{E}_P[x])^2 P(x)$$

# Binary Variables (1)

- Coin flipping: heads=1, tails=0

$$p(x = 1 | \mu) = \mu$$

- Bernoulli Distribution

$$\begin{aligned} \mathrm{Bern}(x | \mu) &= \mu^x (1 - \mu)^{1-x} \\ \mathbb{E}[x] &= \mu \\ \mathrm{var}[x] &= \mu(1 - \mu) \end{aligned}$$

# Binary Variables (2)

- N coin flips:

$$p(m \text{ heads}|N, \mu)$$

- Binomial Distribution

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

$$\mathbb{E}[m] \equiv \sum_{m=0}^{N} m \text{Bin}(m|N, \mu) = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^{N} (m - \mathbb{E}[m])^2 \text{Bin}(m|N, \mu) = N\mu(1 - \mu)$$

# Your first consulting job
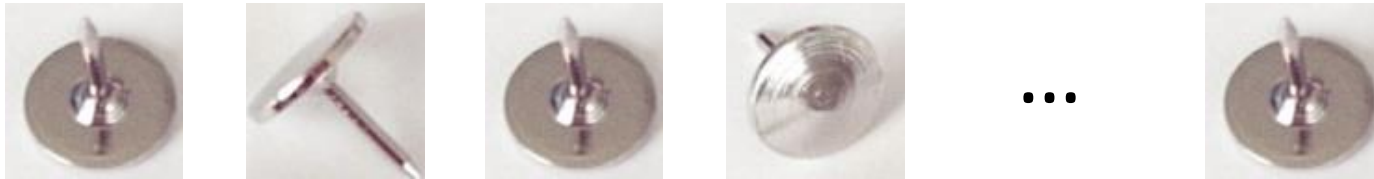
Billionaire in Dallas asks:

- 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) = $\theta$,  P(Tails) = 1-$\theta$



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

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

# Your first parameter learning algorithm

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

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

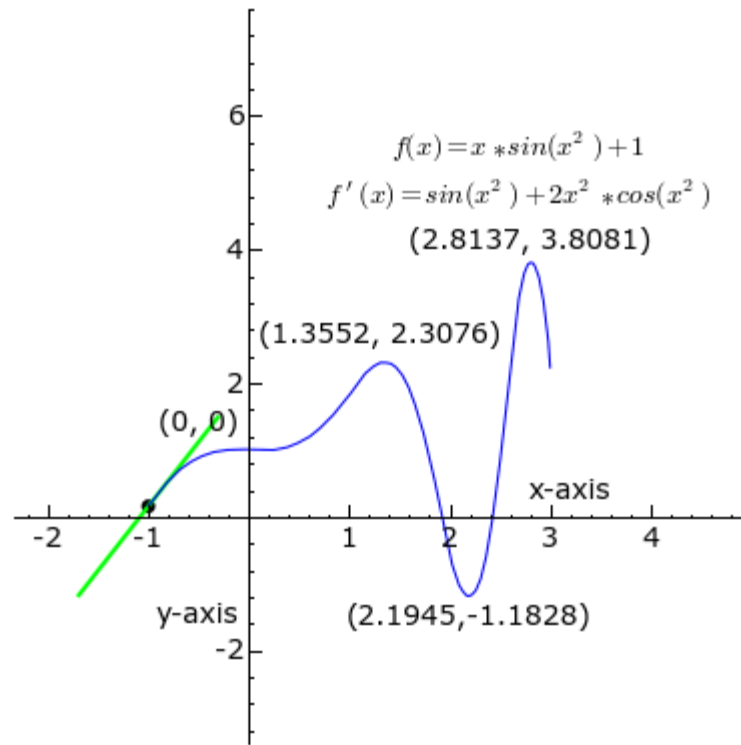- Set derivative to zero, and solve!

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

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

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

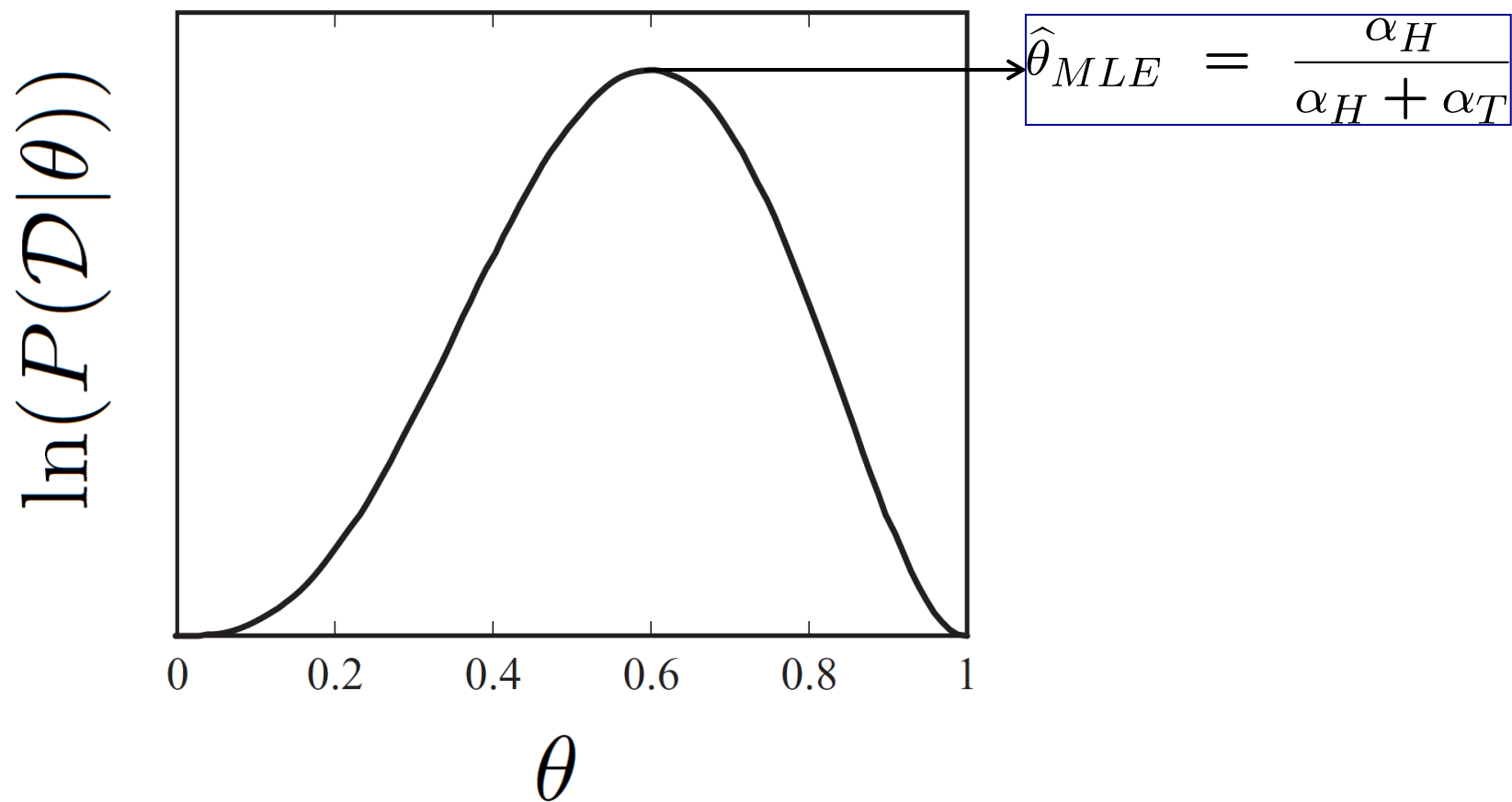$$= \frac{\alpha_H}{\theta} - \frac{\alpha_T}{1-\theta} = 0 \qquad \boxed{\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T}}$$

$f(x) = x * sin(x^2) + 1$

$f'(x) = sin(x^2) + 2x^2 * cos(x^2)$

(2.8137, 3.8081)

(1.3552, 2.3076)

(0, 0)

x-axis

y-axis

(2.1945, -1.1828)

At each point, the derivative is the slope of a line that is tangent to the curve. Note: derivative is **positive where green**, **negative where red**, and **zero where black**.

**Source: Wikipedia.com**

**Data**





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

# But, how many flips do I need?

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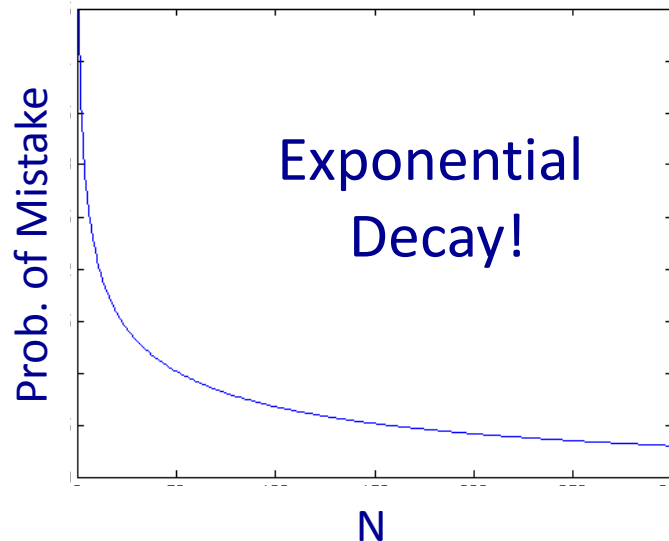
- Billionaire says: I flipped 3 heads and 2 tails.
- You say: $\theta$ = 3/5, I can prove it!
- He says: What if I flipped 30 heads and 20 tails?
- You say: Same answer, I can prove it!
- **He says: What's better?**
- You say: Umm… The more the merrier???
- He says: Is this why I am paying you the big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

# A bound (from Hoeffding's inequality)

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

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

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

Prob. of Mistake

Exponential Decay!

N

# PAC Learning

- PAC: Probably Approximate Correct
- Billionaire says: I want to know the thumbtack $\theta$, within $\varepsilon = 0.1$, with probability of mistake, $\delta <= 0.05$.
- How many flips? Or, how big do I set $N$?

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

$P(\text{mistake})$ is less than or equal to $2e^{-2N\epsilon^2} \leq \delta$

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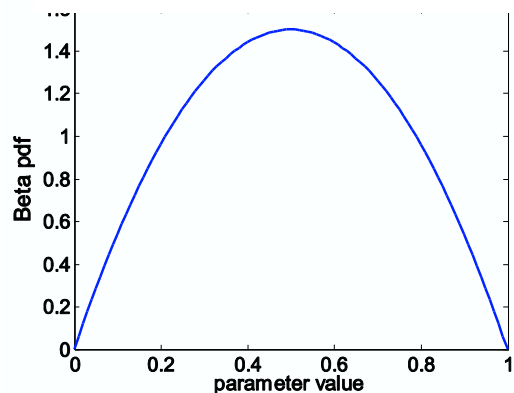
$$N \geq \frac{\ln(2/\delta)}{2\epsilon^2}$$

*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…**

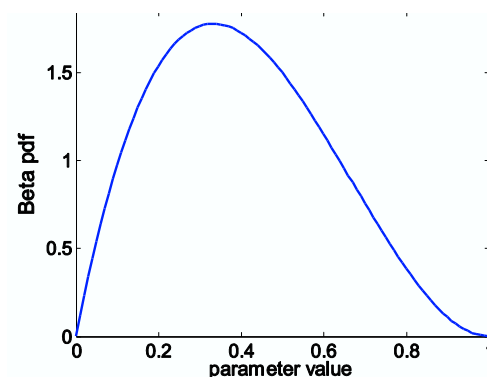- Rather than estimating a single $\theta$, we obtain a distribution over possible values of $\theta$

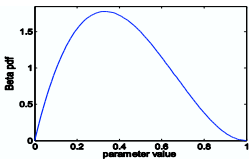In the beginning

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

After observations

# Bayesian Learning
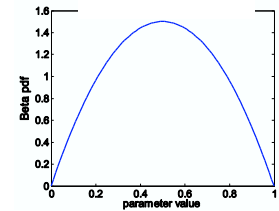
Use Bayes rule:

Data Likelihood

Prior

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

Posterior

Normalization

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

Also, for uniform priors:

→ reduces to MLE objective

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

# Bayesian Learning for Thumbtacks

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

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** $\mu \in [0, 1]$.        $B(a,b) = \dfrac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}$

$$\text{Beta}(\mu|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\mu^{a-1}(1-\mu)^{b-1}$$

$$\mathbb{E}[\mu] = \frac{a}{a+b}$$

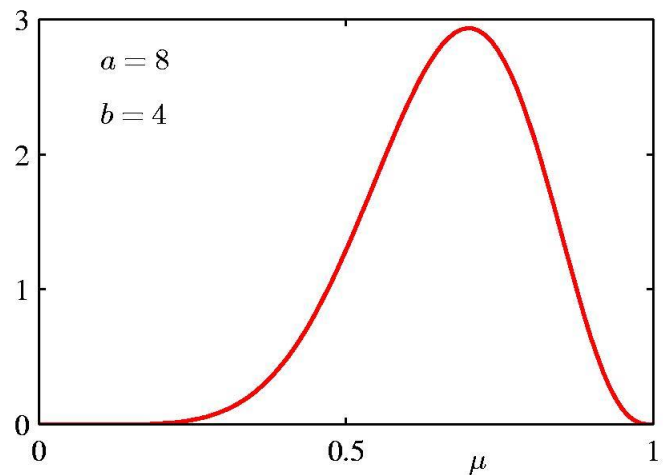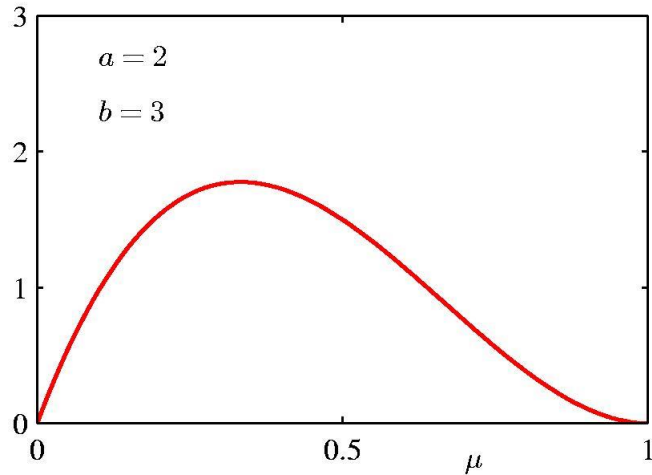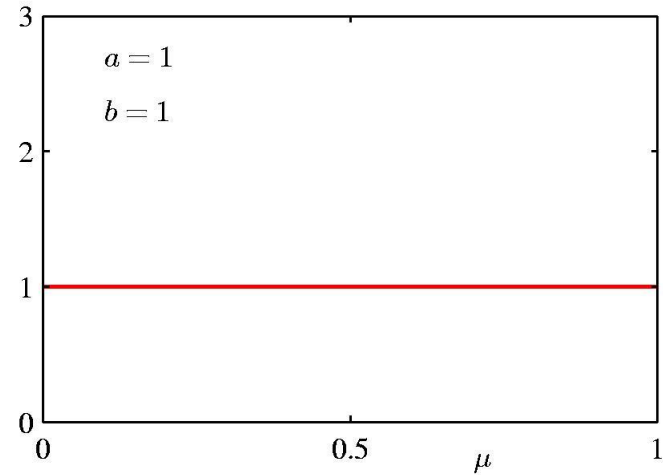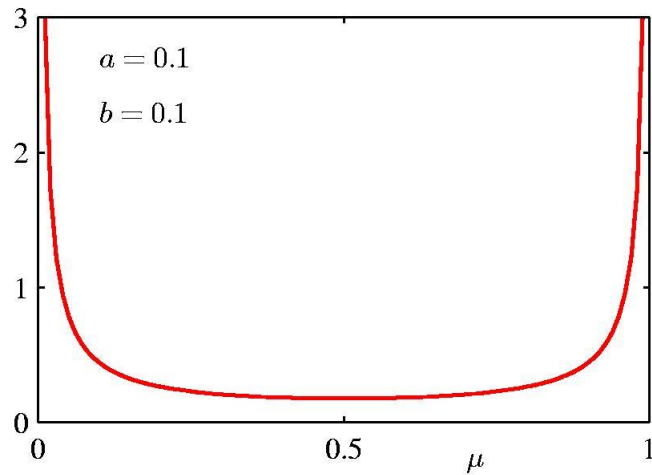$$\text{var}[\mu] = \frac{ab}{(a+b)^2(a+b+1)}$$

$$B(a,b) = \int_0^1 u^{a-1}(1-u)^{b-1}du, \quad \text{a>0, b>0}$$

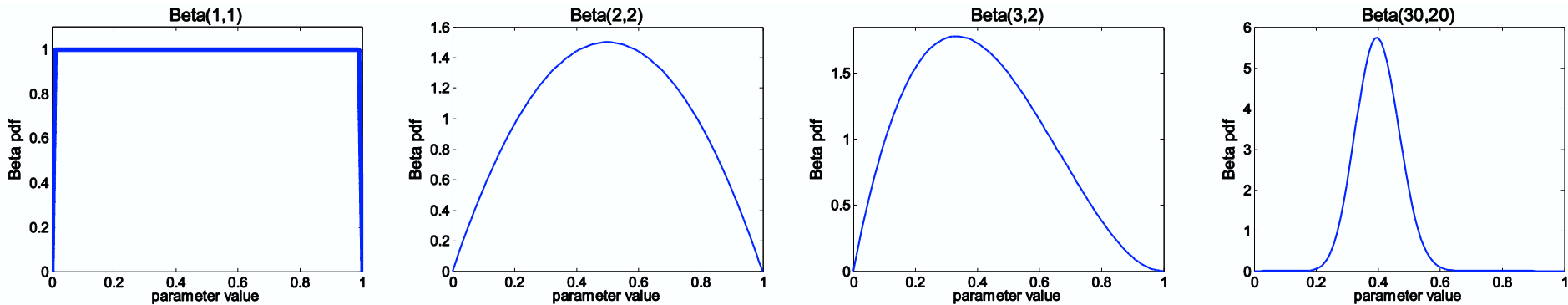$$\Gamma(a) = \int_0^\infty u^{a-1}e^{-a}du$$

# Beta Distribution

$$\text{Beta}(\mu|a, b) \quad = \quad \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\mu^{a-1}(1-\mu)^{b-1}$$

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

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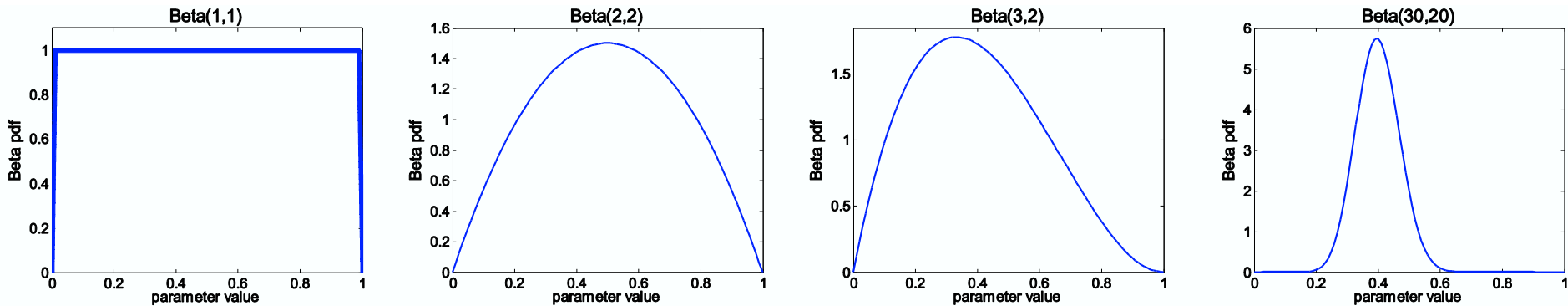
$$= \theta^{\alpha_H + \beta_H - 1}(1 - \theta)^{\alpha_T + \beta_T - 1}$$

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

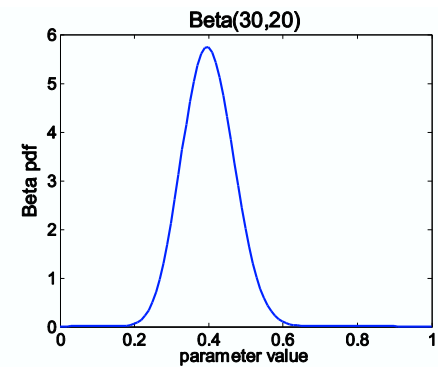# Posterior Distribution

- 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:

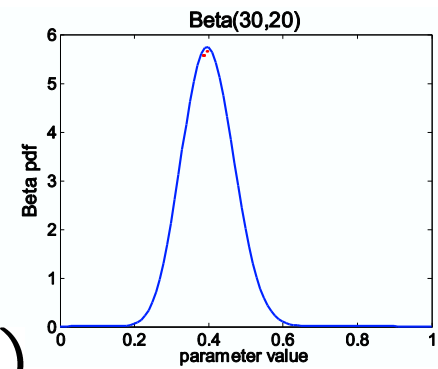$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

- 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)

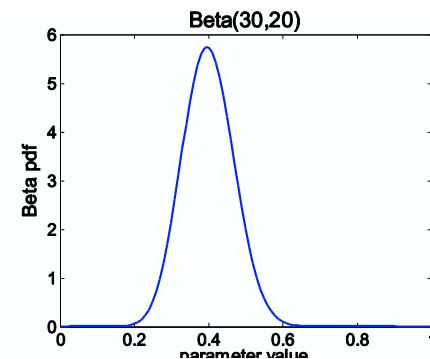$$P(\theta \mid \mathcal{D}) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

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

- 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:

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

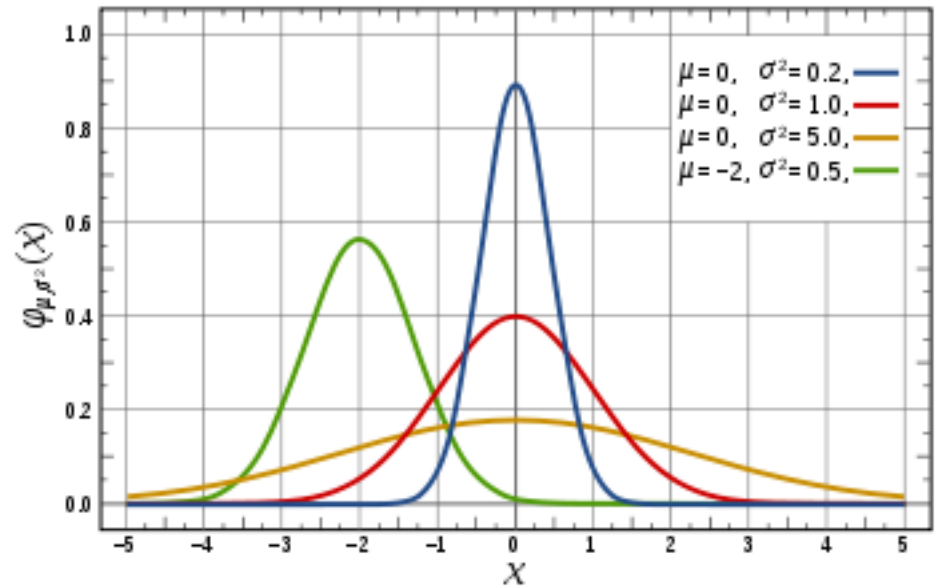Beta prior equivalent to extra thumbtack flips

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

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

# What about continuous variables?

- Billionaire says: If I am measuring a continuous variable, what can you do for me?

- **You say: Let me tell you about Gaussians...**



$$P(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

# Learning a Gaussian

| $X_i = i$ | Exam Score |
|---|---|
| 0 | 85 |
| 1 | 95 |
| 2 | 100 |
| 3 | 12 |
| ... | ... |
| 99 | 89 |

- Collect a bunch of data
  - Hopefully, i.i.d. samples
  - e.g., exam scores

- Learn parameters
  - Mean: $\mu$
  - Variance: $\sigma$

$$P(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$$

# MLE for Gaussian: $P(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}}$

- Prob. of i.i.d. samples $D = \{x_1, \ldots, x_N\}$:

$$P(\mathcal{D} \mid \mu, \sigma) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^N \prod_{i=1}^{N} e^{\frac{-(x_i-\mu)^2}{2\sigma^2}}$$

$$\mu_{MLE}, \sigma_{MLE} = \arg\max_{\mu,\sigma} P(\mathcal{D} \mid \mu, \sigma)$$

- Log-likelihood of data:

$$\ln P(\mathcal{D} \mid \mu, \sigma) = \ln\left[\left(\frac{1}{\sigma\sqrt{2\pi}}\right)^N \prod_{i=1}^{N} e^{\frac{-(x_i-\mu)^2}{2\sigma^2}}\right]$$

$$= -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2}$$

# Your second learning algorithm: MLE for mean of a Gaussian

- What's MLE for mean?

$$\frac{d}{d\mu} \ln P(\mathcal{D} \mid \mu, \sigma) = \frac{d}{d\mu} \left[ -N \ln \sigma\sqrt{2\pi} - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \frac{d}{d\mu} \left[ -N \ln \sigma\sqrt{2\pi} \right] - \sum_{i=1}^{N} \frac{d}{d\mu} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= -\sum_{i=1}^{N} \frac{(x_i - \mu)}{\sigma^2} = 0$$

$$= -\sum_{i=1}^{N} x_i + N\mu = 0$$

$$\widehat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^{N} x_i$$

# MLE for variance

- Again, set derivative to zero:

$$\frac{d}{d\sigma} \ln P(\mathcal{D} \mid \mu, \sigma) = \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} - \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= \frac{d}{d\sigma} \left[ -N \ln \sigma \sqrt{2\pi} \right] - \sum_{i=1}^{N} \frac{d}{d\sigma} \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$= -\frac{N}{\sigma} + \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{\sigma^3} = 0$$

$$\boxed{\widehat{\sigma}^2_{MLE} = \frac{1}{N} \sum_{i=1}^{N} (x_i - \widehat{\mu})^2}$$

# Learning Gaussian parameters

- MLE:

$$\widehat{\mu}_{MLE} = \frac{1}{N}\sum_{i=1}^{N} x_i$$

$$\widehat{\sigma}^2_{MLE} = \frac{1}{N}\sum_{i=1}^{N} (x_i - \widehat{\mu})^2$$

- BTW. MLE for the variance of a Gaussian is **biased**
  - Expected result of estimation is **not** true parameter!
  - Unbiased variance estimator:

$$\widehat{\sigma}^2_{unbiased} = \frac{1}{N-1}\sum_{i=1}^{N} (x_i - \widehat{\mu})^2$$