# Density Estimation
# (Parametric Approach)

# Parametric Approach

- P(Y|X)=P(X|Y)P(Y)/P(X)
- P(Y):easy. Experience or training data.
- P(X|Y):difficult. Too few training data; high-dimensional feature space (computation,storage).
- Parametric Approach
  - Form known
  - Parameters unknown

# Your first consulting job

A billionaire from the suburbs of Seattle asks you a question:

- He says: I have a coin, if I flip it, what's the probability it will fall with the head up?
- You say: Please flip it a few times:



- You say: The probability is: **3/5**
- **He says: Why???**
- You say: Because...

# Bernoulli Distribution

Data, D = 

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

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

Choose $\theta$ that maximizes the probability of observed data

# Maximum Likelihood Estimation

Choose θ that maximizes the probability of observed data

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

MLE of probability of head:

$$\widehat{\theta}_{MLE} = \frac{\alpha_H}{\alpha_H + \alpha_T} \quad \textbf{= 3/5}$$

"Frequency of heads"

# 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: Hmm... The more the merrier???
- He says: Is this why I am paying you the big bucks???

# Simple bound (Hoeffding's inequality)

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

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

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

# PAC Learning

- PAC: Probably Approximate Correct

- Billionaire says: I want to know the coin parameter $\theta$, within $\epsilon$ = 0.1, with probability at least 1-$\delta$ = 0.95. How many flips?
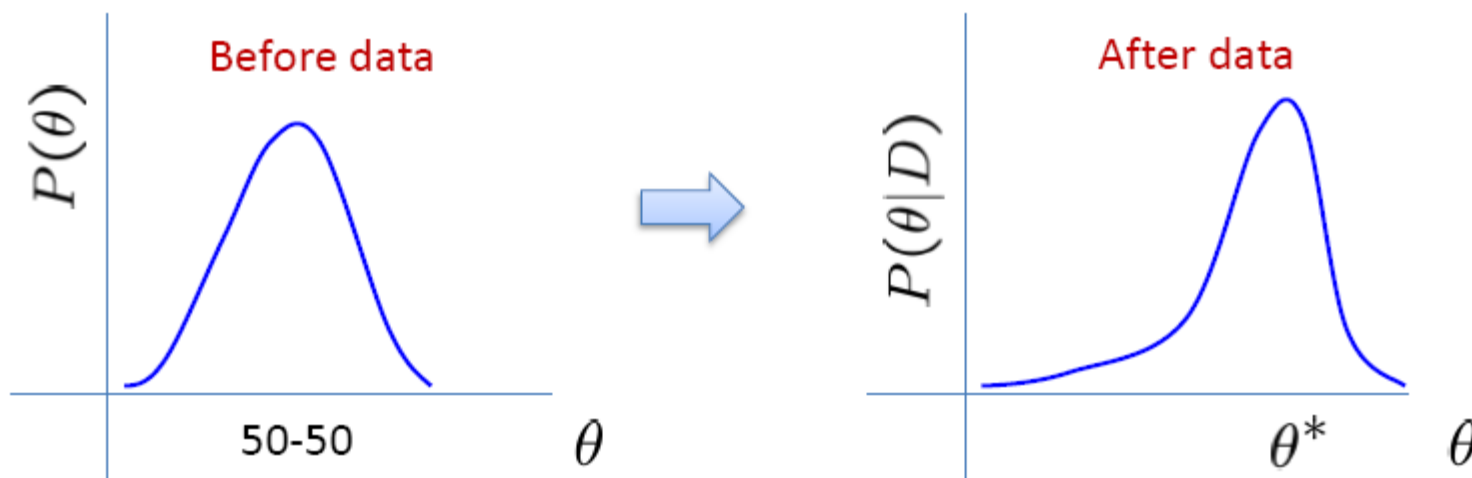
$$P(|\hat{\theta} - \theta^*| \geq \epsilon) \leq 2e^{-2n\epsilon^2}$$

Sample complexity

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

# What about prior knowledge?

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

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

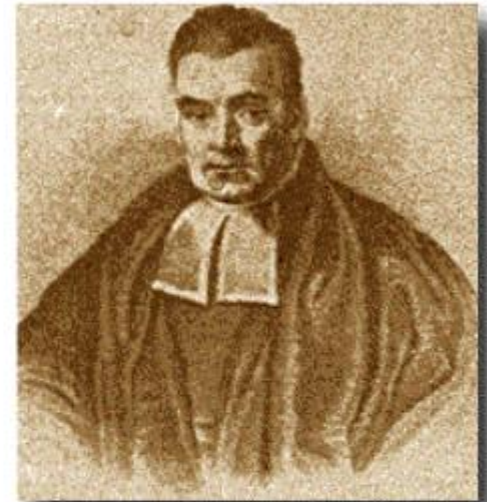- Rather than estimating a single $\theta$, we obtain a distribution over possible values of $\theta$

# Bayesian Learning

- Use Bayes rule:

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

- Or equivalently:

$$\underbrace{P(\theta \mid \mathcal{D})}_{\text{posterior}} \propto \underbrace{P(\mathcal{D} \mid \theta)}_{\text{likelihood}}\underbrace{P(\theta)}_{\text{prior}}$$
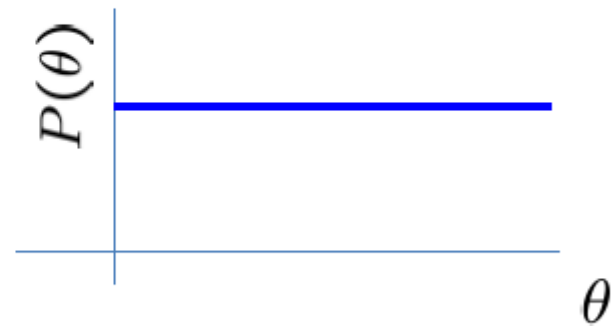
**Bayes, Thomas (1763)** An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London,* **53:370-418**

# Prior Distribution

- What about prior?
  - Represents expert knowledge (philosophical approach)
  - Simple posterior form (engineer's approach)

- Uninformative priors:
  - Uniform distribution

- Conjugate priors:
  - Closed-form representation of posterior
  - $P(\theta)$ and $P(\theta|D)$ have the same form

# Conjugate Prior (I)

- P($\theta$) and P($\theta$|D) have the same form

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

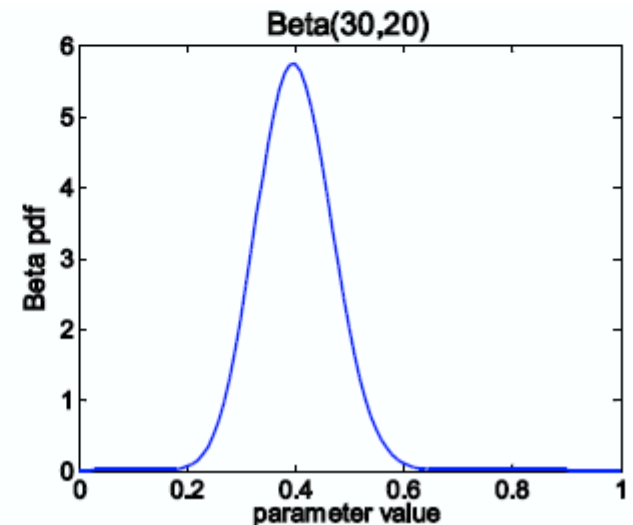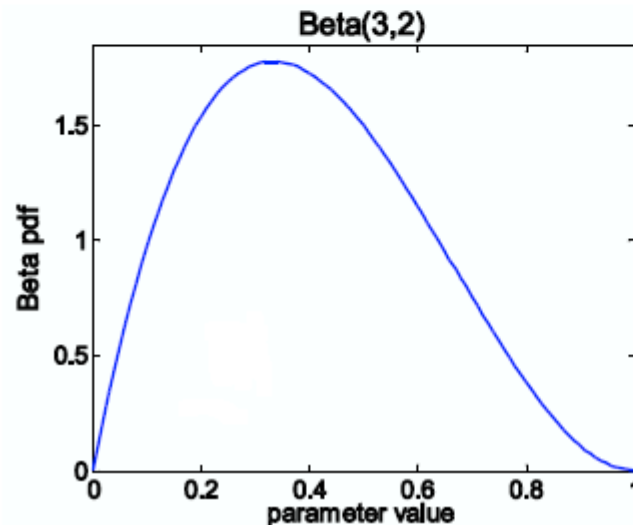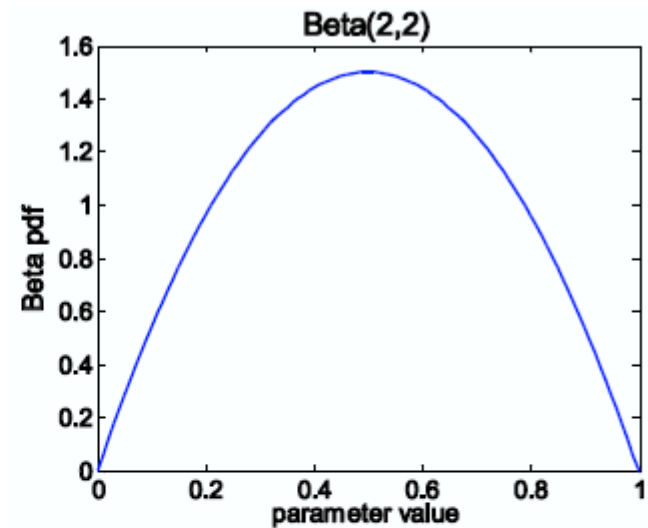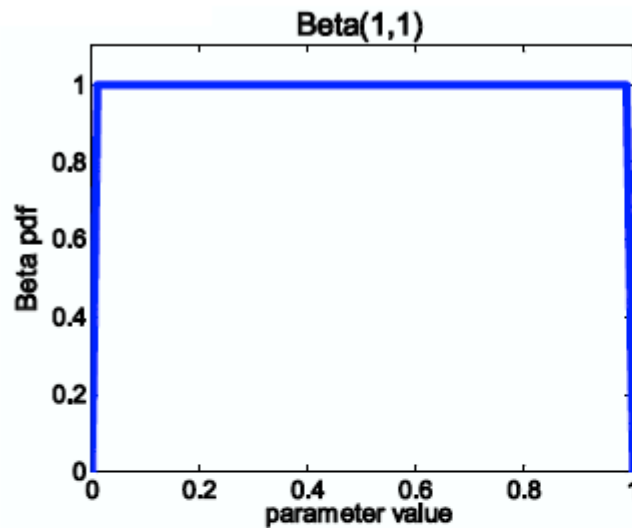$$P(\theta|D) \sim Beta(\beta_H + \alpha_H, \beta_T + \alpha_T)$$

**For Binomial, conjugate prior is Beta distribution.**
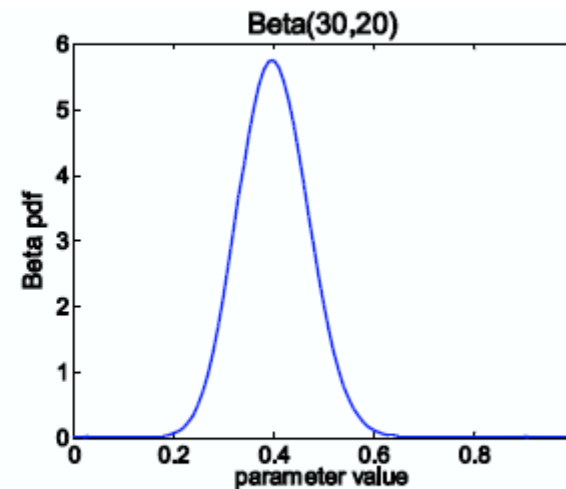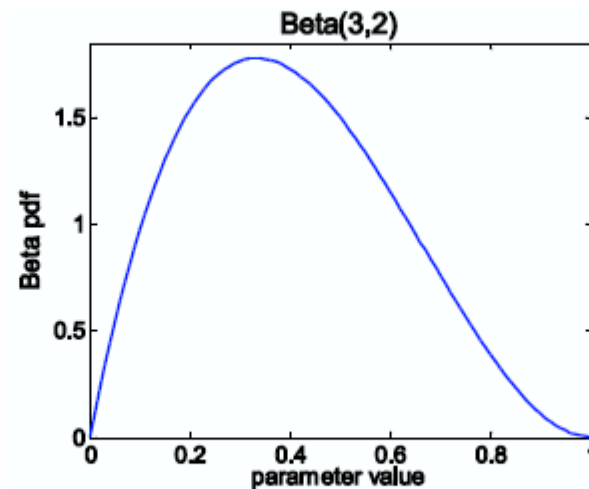
# Beta Distribution

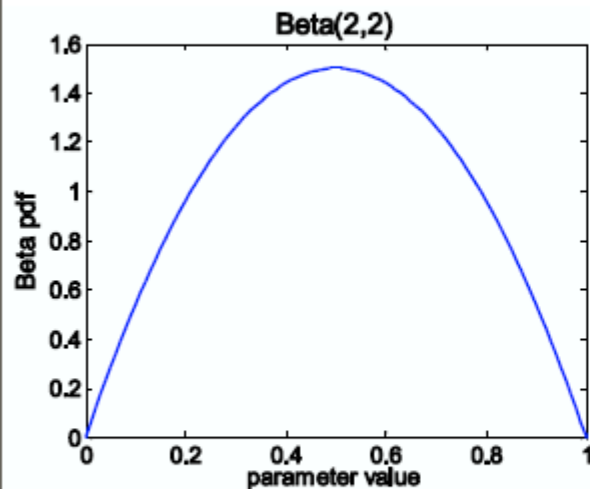$Beta(\beta_H, \beta_T)$ More concentrated as values of $\beta_H$, $\beta_T$ increase

# Beta conjugate prior

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



As $n = \alpha_H + \alpha_T$ increases

As we get more samples, effect of prior is "washed out"

# Conjugate Prior (II)

- $P(\theta)$ and $P(\theta|D)$ have the same form

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,

$$P(\theta) = \frac{\prod_{i=1}^{k} \theta_i^{\beta_i - 1}}{B(\beta_1, \ldots, \beta_k)} \sim \text{Dirichlet}(\beta_1, \ldots, \beta_k)$$

Then posterior is Dirichlet distribution

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

**For Multinomial, conjugate prior is Dirichlet distribution.**

# Maximum A Posterior Estimation

Choose $\theta$ that maximizes a posterior probability

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

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

MAP estimate of probability of head:

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

$$\hat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

Mode of Beta distribution

# MLE vs. MAP

- Maximum Likelihood estimation (MLE)

  Choose value that maximizes the probability of observed data

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

- Maximum *a posteriori* (MAP) estimation

  Choose value that is most probable given observed data and prior belief

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

When is MAP same as MLE?

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

What if we toss the coin too few times?

- You say: Probability next toss is a head = 0

- Billionaire says: You're fired!        ...with prob 1 ☺

$$\widehat{\theta}_{MAP} = \frac{\alpha_H + \beta_H - 1}{\alpha_H + \beta_H + \alpha_T + \beta_T - 2}$$

- Beta prior equivalent to extra coin flips **(regularization)**
- As $n \to \infty$, prior is "forgotten"
- **But, for small sample size, prior is important!**
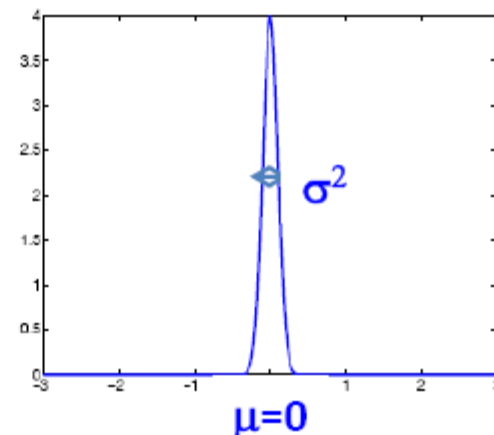
# Bayesian vs. Frequentists

# 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...**

$$P(x \mid \mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{\frac{-(x-\mu)^2}{2\sigma^2}} \quad = N(\mu, \sigma^2)$$
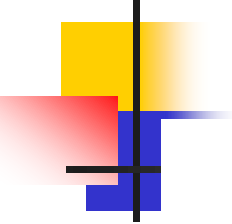
# Properties of Gaussians

- affine transformation (multiplying by scalar and adding a constant)
  - $X \sim N(\mu, \sigma^2)$
  - $Y = aX + b \rightarrow Y \sim N(a\mu + b, a^2\sigma^2)$

- Sum of Gaussians
  - $X \sim N(\mu_X, \sigma^2_X)$
  - $Y \sim N(\mu_Y, \sigma^2_Y)$
  - $Z = X+Y \rightarrow Z \sim N(\mu_X + \mu_Y, \sigma^2_X + \sigma^2_Y)$

# MLE for Gaussian mean and variance

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

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

Note: MLE for the variance of a Gaussian is **biased**

- Expected result of estimation is **not** true parameter!
- Unbiased variance estimator:

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

# Proof

$$\theta = [\theta_1, \theta_2]^T, \ \theta_1 = \mu, \ \theta_2 = \sigma^2$$

$$p(x \mid \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right]$$

$$X = \{x_1, x_2, \cdots, x_N\}$$

$$l(x) = p(X \mid \theta) = \prod_{k=1}^{N} p(x_k \mid \theta)$$

$$H(\theta) = \ln l(x) = \sum_{k=1}^{N} \ln P(x_k \mid \theta)$$

# Proof

$$\nabla_\theta H(\theta) = \sum_{k=1}^{N} \nabla_\theta \ln p(x_k \mid \theta) = 0$$

$$\ln p(x_k \mid \theta) = -\frac{1}{2} \ln 2\pi\theta_2 - \frac{1}{2\theta_2} (x_k - \theta_1)^2$$

$$\nabla_\theta \ln p(x_k \mid \theta) = \begin{bmatrix} \dfrac{1}{\theta_2}(x_k - \theta_1) \\[2ex] -\dfrac{1}{2\theta_2} + \dfrac{1}{2\theta_2^2}(x_k - \theta_1)^2 \end{bmatrix} \quad \begin{cases} \displaystyle\sum_{k=1}^{N} \dfrac{1}{\hat{\theta}_2}(x_k - \hat{\theta}_1) = 0 \\[3ex] -\displaystyle\sum_{k=1}^{N} \dfrac{1}{\hat{\theta}_2} + \sum_{k=1}^{N} \dfrac{(x_k - \hat{\theta}_1)^2}{\theta_2^2} = 0 \end{cases}$$

$$\hat{\mu} = \hat{\theta}_1 = \frac{1}{N}\sum_{k=1}^{N} x_k$$

$$\hat{\sigma}^2 = \hat{\theta}_2 = \frac{1}{N}\sum_{k=1}^{N} (x_k - \hat{\mu})^2$$

# MAP for Gaussian mean and variance

- Conjugate priors

  - Mean: Gaussian prior

  - Variance: Wishart Distribution

- Prior for mean:

$$P(\mu \mid \eta, \lambda) = \frac{1}{\lambda\sqrt{2\pi}} e^{\frac{-(\mu-\eta)^2}{2\lambda^2}} \qquad = N(\eta, \lambda^2)$$

# MAP for Gaussian Mean

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

$$\widehat{\mu}_{MAP} \;=\; \frac{\frac{1}{\sigma^2}\sum_{i=1}^{n} x_i + \frac{\eta}{\lambda^2}}{\frac{n}{\sigma^2} + \frac{1}{\lambda^2}}$$

(Assuming known variance $\sigma^2$)

# What you should know…

- Learning parametric distributions: form known, parameters unknown

  - Bernoulli ($\theta$, probability of flip)

  - Gaussian ($\mu$, mean and $\sigma^2$, variance)

- MLE

- MAP