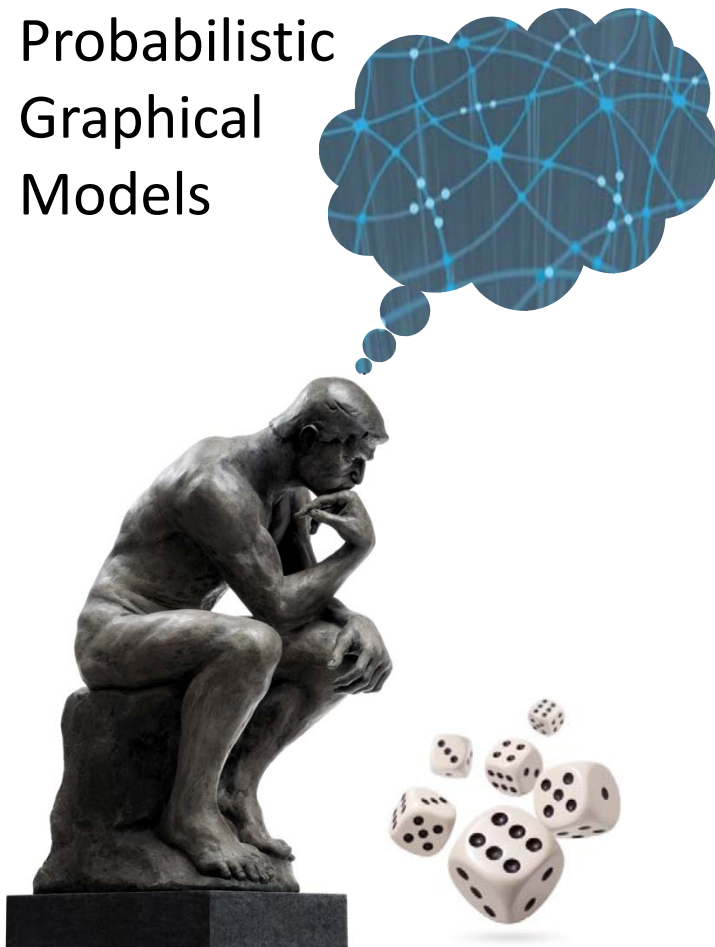


Probabilistic  
Graphical  
Models



Learning

---

Parameter Estimation

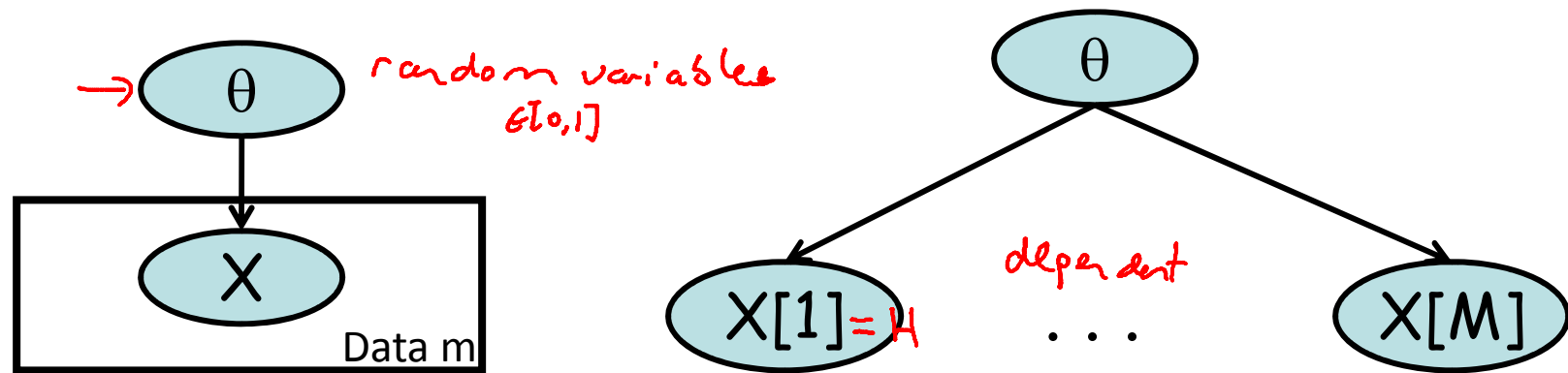
---

# Bayesian Estimation

# Limitations of MLE

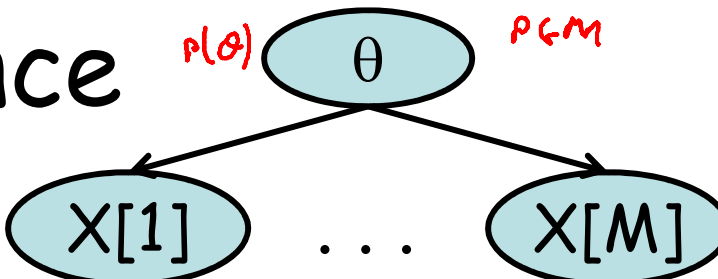
- Two teams play 10 times, and the first wins 7 of the 10 matches
  - ⇒ Probability of first team winning = 0.7
- A coin is tossed 10 times, and comes out 'heads' 7 of the 10 tosses
  - ⇒ Probability of heads = 0.7
- A coin is tossed 10000 times, and comes out 'heads' 7000 of the 10000 tosses
  - ⇒ Probability of heads = 0.7

# Parameter Estimation as a PGM



- Given a fixed  $\theta$ , tosses are independent
- If  $\theta$  is unknown, tosses are not marginally independent
  - each toss tells us something about  $\theta$

# Bayesian Inference

- Joint probabilistic model 

$$\underline{P(x[1], \dots, x[M], \theta)} = \underline{P(x[1], \dots, x[M] | \theta)} \underline{P(\theta)}$$

$$= P(\theta) \prod_{i=1}^M P(x[i] | \theta)$$

$$= \underline{P(\theta)} \underbrace{\theta^{M_H} (1 - \theta)^{M_T}}_{\text{likelihood function}} \leftarrow \text{likelihood function}$$

$$\underline{P(\theta | x[1], \dots, x[M])} = \frac{\underbrace{P(x[1], \dots, x[M] | \theta)}_{\text{likelihood}} \underbrace{P(\theta)}_{\text{prior}}}{\underbrace{P(x[1], \dots, x[M])}_{\text{constant relative to } \theta} \underbrace{\theta}_{\text{data } D}}$$

# Dirichlet Distribution

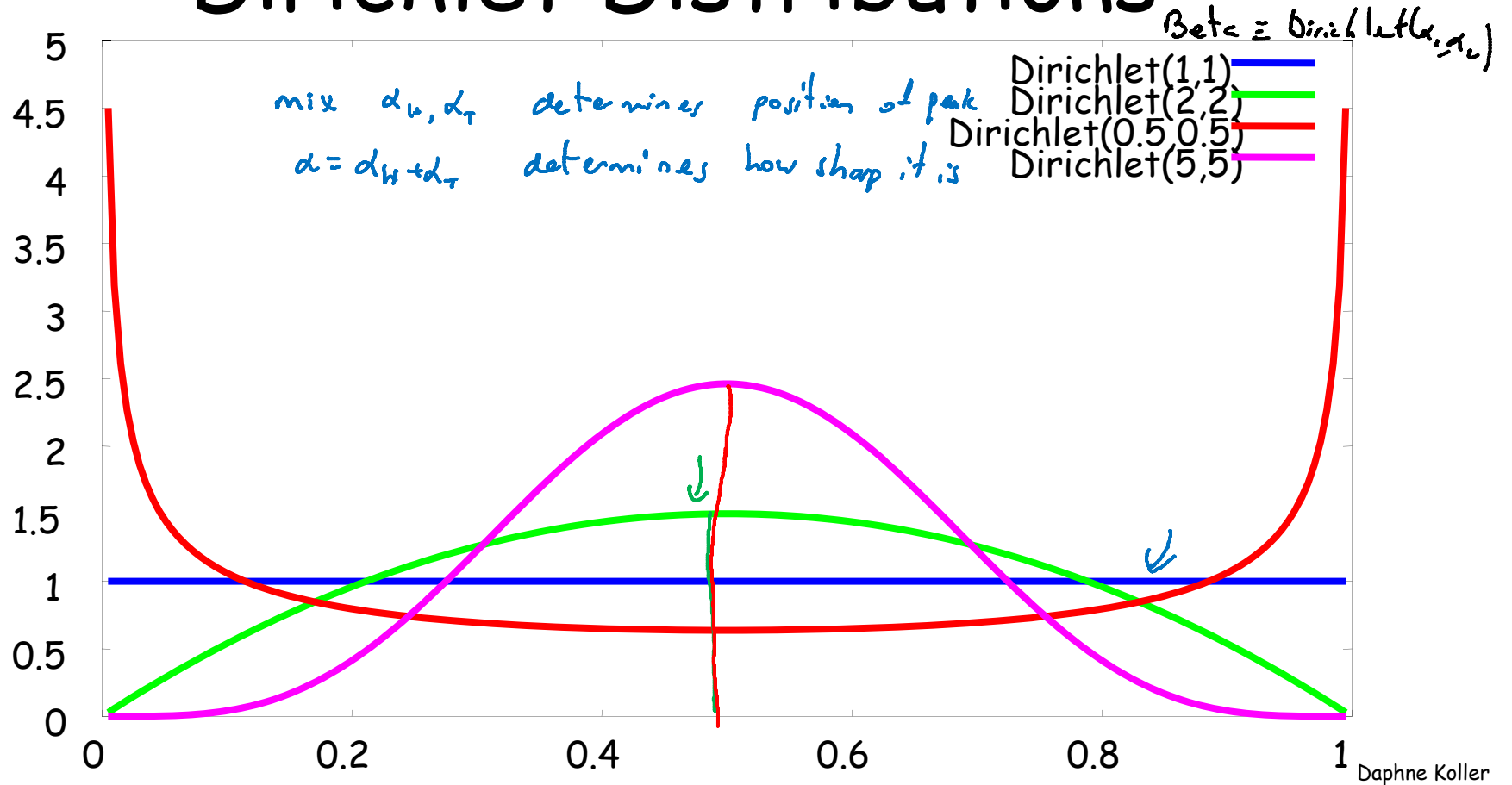
- $\theta$  is a multinomial distribution over  $k$  values
- Dirichlet distribution  $\theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$

– where  $P(\theta)$  =  $\frac{1}{Z} \prod_{i=1}^k \theta_i^{\alpha_i - 1}$  and  $Z = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)}$   $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$

*Handwritten notes: A red arrow points from the  $\alpha_i$  in the Dirichlet distribution to the  $\alpha_i$  in the probability formula. A blue box encloses the product term in the probability formula. A blue box encloses the gamma function ratio in the normalization constant. A red bracket under the  $\alpha_1, \dots, \alpha_k$  in the Dirichlet distribution is labeled "hyperparameters". A red  $\alpha_i$  is written above the word "hyperparameters".*

- Intuitively, hyperparameters correspond to the number of samples we have seen

# Dirichlet Distributions



# Dirichlet Priors & Posteriors

$$\overbrace{P(\theta | D)}^{\text{posterior}} \propto \overbrace{P(D | \theta)}^{\text{likelihood}} \overbrace{P(\theta)}^{\text{prior}}$$

$$m_i = \# \text{ instances with } x_i \quad P(D | \theta) = \prod_{i=1}^k \underbrace{\theta_i^{M_i}}_{\text{multinomial } \theta} \quad \theta_i^{m_i + \alpha_i - 1} \quad P(\theta) \propto \prod_{i=1}^k \underbrace{\theta_i^{\alpha_i - 1}}$$

- If  $P(\theta)$  is Dirichlet and the likelihood is multinomial, then the posterior is also Dirichlet
  - Prior is  $\text{Dir}(\alpha_1, \dots, \alpha_k)$
  - Data counts are  $M_1, \dots, M_k$
  - Posterior is  $\text{Dir}(\alpha_1 + M_1, \dots, \alpha_k + M_k)$
- Dirichlet is a conjugate prior for the multinomial

prior, posterior have the same form

# Summary

- Bayesian learning treats parameters as random variables
  - Learning is then a special case of inference
- Dirichlet distribution is conjugate to multinomial
  - Posterior has same form as prior
  - Can be updated in closed form using sufficient statistics from data