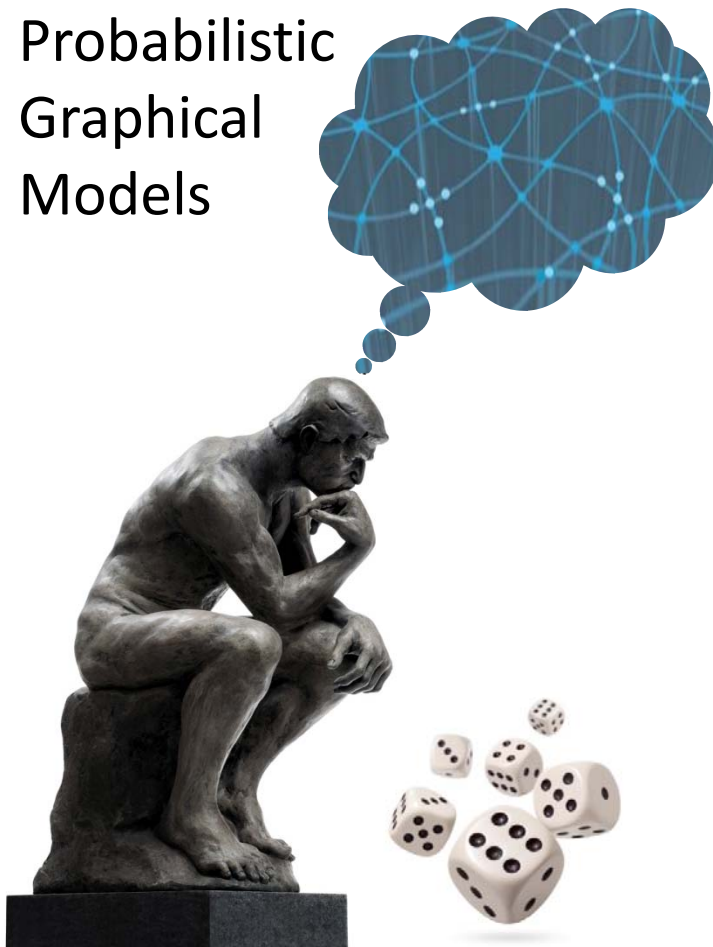


Probabilistic  
Graphical  
Models



Learning

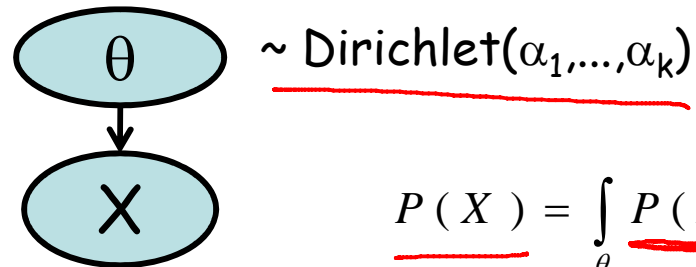
---

Parameter Estimation

---

# Bayesian Prediction

# Bayesian Prediction



$$\underline{P(X)} = \int_{\theta} \underline{P(X | \theta)} \underline{P(\theta)} d\theta$$

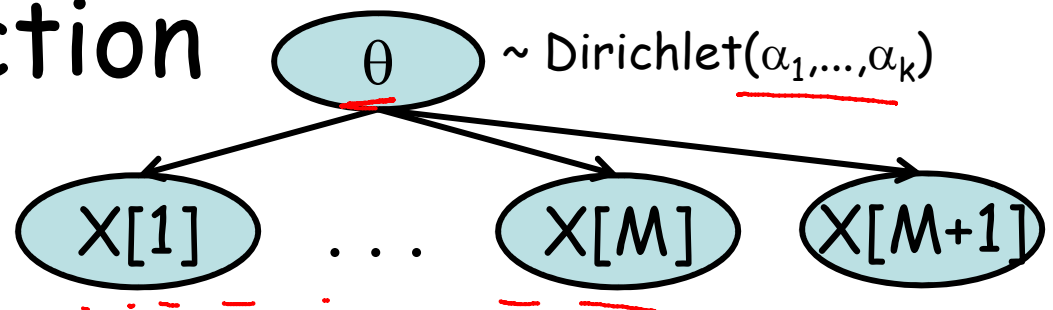
← marginalizing over  $\theta$

$$\begin{aligned} \underline{P(X = \underline{x^i} | \theta)} &= \frac{1}{Z} \int_{\theta} \theta_i \cdot \underbrace{\prod_j \theta^{\alpha_j - 1}}_{\text{prior}} d\theta \\ &= \frac{\alpha_i}{\sum_j \alpha_j = \alpha} \end{aligned}$$

fraction of instances we've seen where  $x^i$

- Dirichlet hyperparameters correspond to the number of samples we have seen

# Bayesian Prediction



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

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

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

~ Dirichlet( $\alpha_1 + M_1, \dots, \alpha_k + M_k$ )  
Posterior over  $\theta$  given  $D$

$$P(X[M+1] = \underline{x^i} | \theta, x[1], \dots, x[M]) = \frac{\alpha_i + M_i}{\underline{\alpha + M}}$$

$\alpha = \sum \alpha_i$   
 $M = \sum M_i$

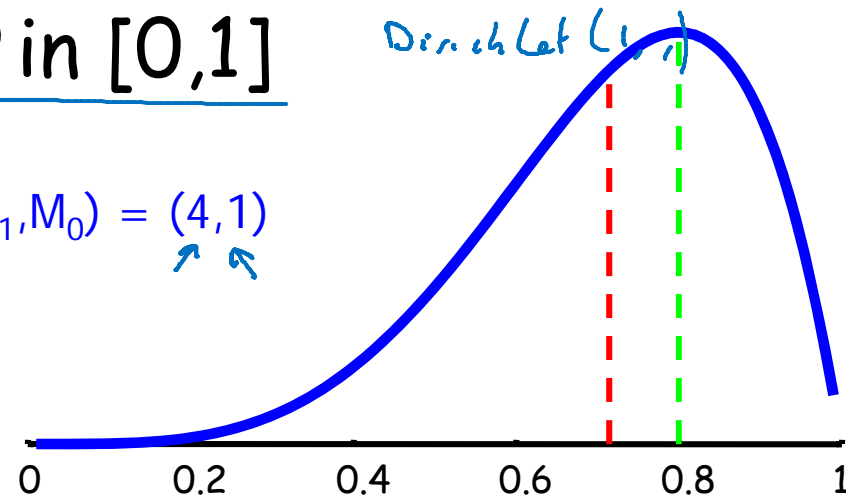
- Equivalent sample size  $\alpha = \alpha_1 + \dots + \alpha_k$ 
  - Larger  $\alpha \Rightarrow$  more confidence in our prior

# Example: Binomial Data

- Prior: uniform for  $\theta$  in  $[0,1]$

$$P(\theta) = \frac{1}{Z} \prod_k \theta^{\alpha_k - 1}$$

$$(M_1, M_0) = (4, 1)$$

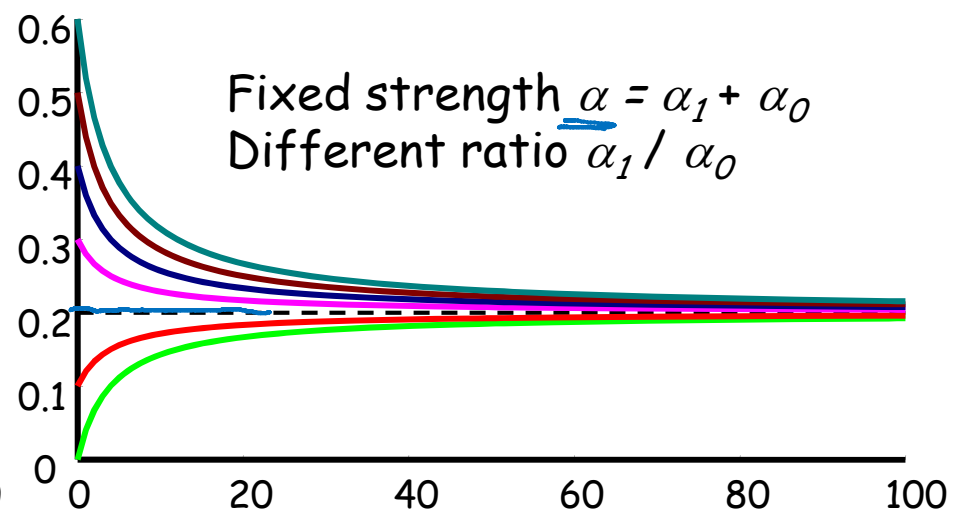
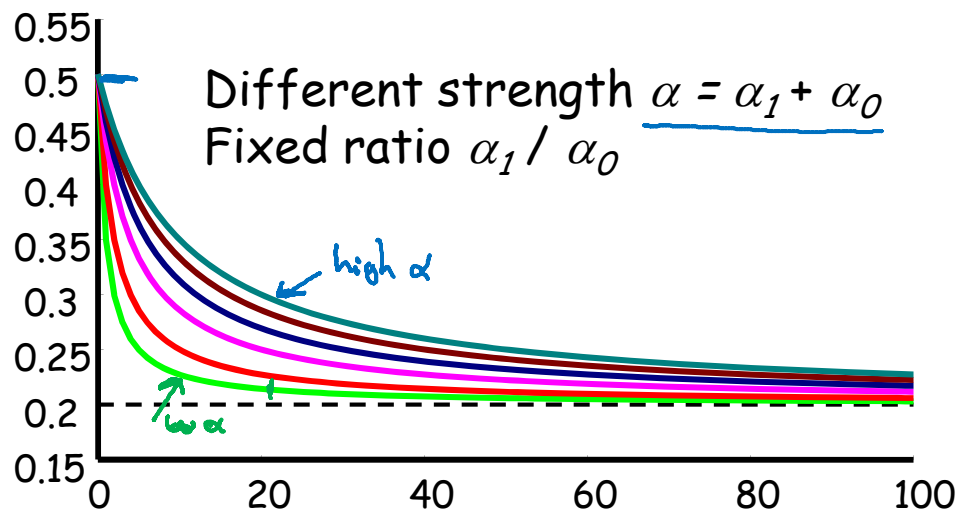


- MLE for  $P(X[6]=1)=4/5$
- Bayesian prediction is  $5/7$

$$\frac{\alpha_1 + m_1}{\alpha + m} = \frac{1 + 4}{2 + 5}$$

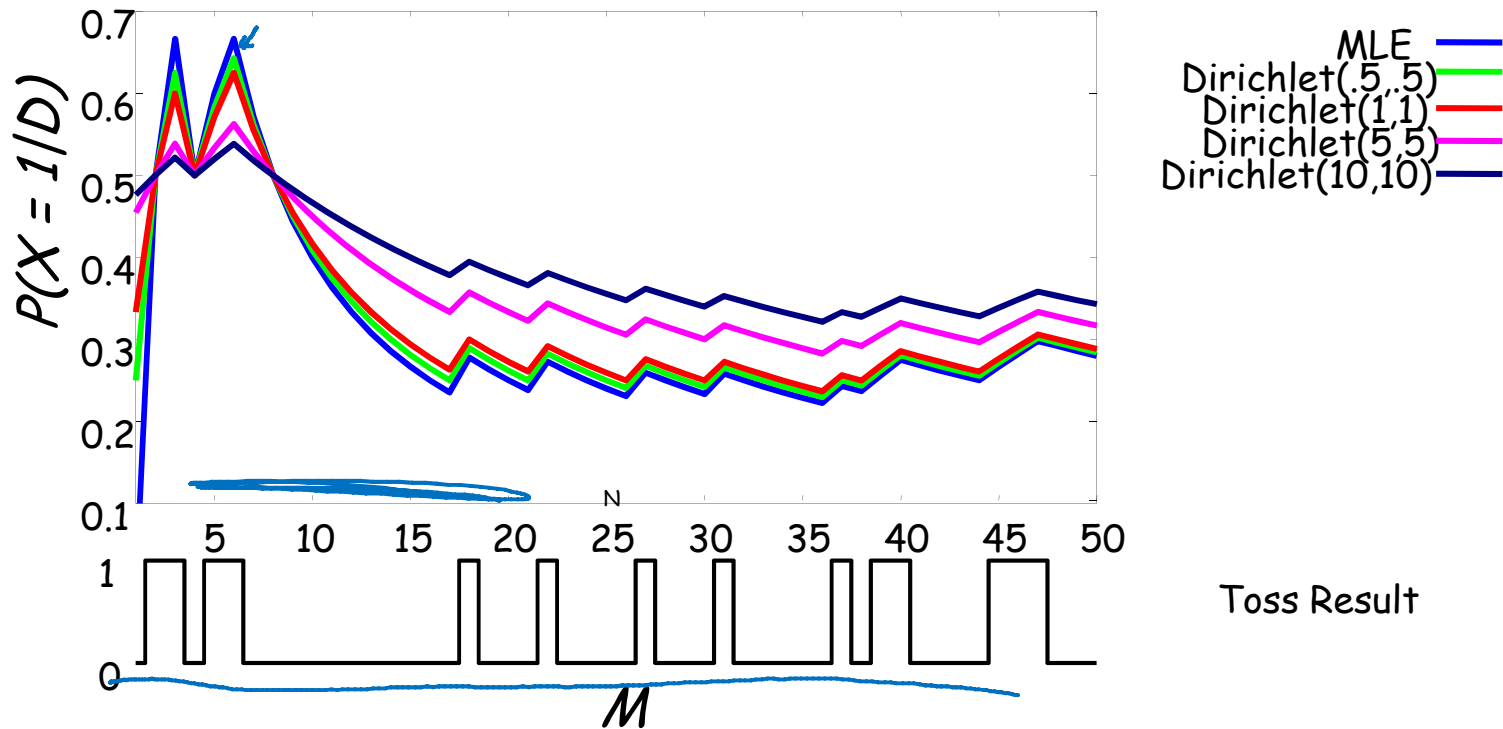
# Effect of Priors

- Prediction of  $P(X=1)$  after seeing data with  $M_1 = \frac{1}{4}M_0$  as a function of sample size  $M$



# Effect of Priors

- In real data, Bayesian estimates are less sensitive to noise in the data



# Summary

- Bayesian prediction combines sufficient statistics from imaginary Dirichlet samples and real data samples
- Asymptotically the same as MLE
- But Dirichlet hyperparameters determine both the prior beliefs and their strength