# Bayesian Concept Learning

- POSTERIOR = LIKELIHOOD TIMES ~~PRIOR~~ PRIOR, NORMALIZED BY DATA/EVIDENCE

- WHEN DATA ENOUGH $\longrightarrow$ POSTERIOR PEAKS ON SINGLE DATA/CONCEPT: **MAP ESTIMATE**

- WITH MORE DATA $\longrightarrow$ MAP CONVERGES TO MLE, DATA OVERWHELMS THE PRIOR

- SMALL OR AMBIGUOUS DATASET $\longrightarrow$ PLUG-IN APPROXIMATION

**GENERATIVE MODEL**: FULLY PROBABILISTIC MODEL OF ALL VARIABLES

**DISCRIMINATIVE MODEL**: MODEL ONLY FOR TARGET VARIABLES CONDITIONED ON OBSERVATION

# Beta Binomial Model

"INFERRING PROBABILITY COIN SHOWS HEADS, GIVEN SERIES OF OBSERVED TOSSES"

- LIKELIHOOD $P(D|\theta) = \theta^{N_1}(1-\theta)^{N_0}$, $N_1 = \sum(x=1)$ $N_0 = \sum(x=0)$, SUFFICIENT STATISTICS

- PRIOR: CONJUGATE PRIOR WHEN HAS SAME FORM AS LIKELIHOOD

"BINOMIAL P FOLLOWS BETA DISTRIBUTION"

$$BETA(\theta|a,b) \propto \theta^{a-1}(1-\theta)^{b-1}$$ PRIOR PARAMETERS = HYPERPARAMETERS: ENCODE PRIOR BELIEFS

- POSTERIOR: $P(\theta|D) \propto BIN(N_1|\theta, N_0+N_1) \cdot BETA(\theta|a,b) \cdot BETA(\theta|N_1+a, N_0+b)$

  BATCH UPDATE = SEQUENTIAL UPDATE

  - MEAN: $\frac{a+N_1}{a+b+N}$   - VARIANCE: $\frac{\hat{\theta}(1-\hat{\theta})}{N}$

- POSTERIOR PREDICTIVE: $P(x=1|D) = \int_0^1 P(x=1|\theta)P(\theta|D)d\theta = \int_0^1 \theta\, BETA(\theta|a,b)d\theta = E[\theta|D] = \frac{a}{a+b}$

  PREDICT FUTURE OBSERVATIONS

  - EQUIV TO PLUG-IN POSTERIOR MEAN PARAMETERS

  - USING MLE IS POOR WHEN SAMPLE COUNT IS SMALL: ZERO COUNT / SPARSE DATA PROBLEM

    $\hookrightarrow$ ADD ONE SMOOTHING: ADD 1 TO COUNTS

# Dirichlet Multinomial

- LIKELIHOOD $P(D|\theta) = \prod^u \theta_u^{N_u}$

- PRIOR: $DIR(\theta|\alpha) = \frac{1}{B(\alpha)} \prod^u \theta_u^{\alpha_u-1} \cdot I(x \in S_u)$ $\rightarrow$

- POSTERIOR: $P(\theta|D) = DIR(\theta|\alpha_1+N_1, \dots, \alpha_u+N_u)$ $\left(\prod^u \theta^{\alpha_u+N_u-1}\right)$

- POSTERIOR PREDICTIVE: $P(x=j|D) = \frac{\alpha_j+N_j}{\alpha_0+N}$

- APPLICATION: LANGUAGE MODELING | BAG OF WORDS

- DISCRETE MULTIVARIATE VARIABLE

# General Generative Classifier

$$P(y=c|x,\theta) = \frac{P(y=c|\theta)P(x|y=c,\theta)}{\sum_{c'} P(y=c'|\theta)P(x|y=c',\theta)}$$

GENERATE DATA USING CLASS-CONDITIONAL DENSITY $P(x|y=c)$ AND CLASS PRIOR $P(y=c)$