

BAYESIAN STATISTIC

POSTERIOR $P(\theta|D)$

MAP $\text{ARGMAX}_{\theta} P(\theta|D)$ IS POSTERIOR MODE, POPULAR BECAUSE REDUCES TO OPTIMIZATION PROBLEM. HAS DRAWBACKS

- IS POINT ESTIMATE, NO UNCERTAINTY
- NO REPARAMETERIZATION INVARIANT, MODE IS NOT CONSERVED, $\hat{\theta} = \text{ARGMAX}_{\theta} P(D|\theta)P(\theta) \propto |I(\theta)|^{-1/2}$ L FISHER INFO OF $P(x|\theta)$
- MEASURE D , NO TAKES SPACE INTO ACCOUNT \rightarrow USE LOSS FUNCTIONS

0-1 LOSS = $1(\theta \neq \hat{\theta}) \rightarrow$ MAP POSTERIOR MODE

SQUARED LOSS = $(\theta - \hat{\theta})^2 \rightarrow$ POSTERIOR MEAN

$|\theta - \hat{\theta}| \rightarrow$ POSTERIOR MEDIAN

CREDIBLE INTERVALS \rightarrow LIKE CONFIDENCE INTERVALS BUT BAYESIAN, CONTINUOUS REGION W/ $(1-\alpha)$ OF POSTERIOR PROBABILITY MASS

HPD/HDI: SET OF MOST PROBABLE POINTS 100(1- α) OF PROBABILITY MASS

BAYESIAN MODEL SELECTION

CROSSVALIDATION REQUIRES FITTING EACH MODEL M TIMES, WHICH IS INEFFICIENT, LET'S COMPUTE POSTERIOR OVER MODELS

$$P(M|D) = \frac{P(D|M)P(M)}{\sum_M P(D,M)} \rightarrow \hat{M} = \text{ARGMAX}_M P(M|D) \xrightarrow[\text{PICK M MAXIMIZING}]{\text{UNIF. PRIOR}} P(D|M) = \int P(D|\theta)P(\theta|M)d\theta$$

MARGINAL LIKELIHOOD, INTEGRATED LIKELIHOOD, EVIDENCE

WE AREN'T MAXIMIZING PARAMS OR MLE, BUT INTEGRATING THEM AWAY/MARGINALIZING THEM, NO OVERFITTING D/E MORE PARAMS \neq BETTER MARGINAL LIKELIHOOD 'OCCAM'S RAZOR', SIMPLEST MODEL EXPLAINING DATA

NO RESTRICTED TO GRID OF VALUES, NUMERICAL OPTIMIZATION, EMPIRICAL BAYES/ TYPE II MAX. LIKELIHOOD $\hat{\lambda} = \text{ARGMAX}_{\lambda} P(D|\lambda)$ \rightarrow POINT ESTIMATE OF POSTERIOR HYPERPARAMS

WE NEED TO COMPUTE MARGINAL LIKELIHOOD $P(D|M) \rightarrow$ EASY WHEN WE HAVE CONJUGATE PRIOR

$$P(\theta) = \frac{Q(\theta)}{Z_0} \leftarrow \text{NORM. CONSTANT} \quad \text{LIKELIHOOD} \quad P(D|\theta) = \frac{Q(D|\theta)}{Z_L} \quad \text{POSTERIOR} \quad P(\theta|D) = \frac{Q(\theta|D)}{Z_N} \quad Q(D|\theta) = Q(D|\theta)Q(\theta) \quad \left| \quad \frac{Q(\theta|D)}{Z_N} = \frac{Q(D|\theta)Q(\theta)}{Z_L Z_0 P(D)} \right. \quad P(D) = \frac{Z_N}{Z_0 Z_L}$$

MARGINAL LIKELIHOOD CAN BE TOUGH \rightarrow BAYESIAN INFORMATION CRITERION

$$BIC = \log P(D|\hat{\theta}) - \frac{\text{DOF}(\hat{\theta})}{2} \log N \propto \log P(D), \quad \hat{\theta} \text{ IS MLE} \quad \text{IS PENALIZED LOG LIKELIHOOD}$$

$$BIC - \text{COST} = -2 \log P(D|\hat{\theta}) + \text{DOF}(\hat{\theta}) \log N \approx -2 \log P(D) \rightarrow \text{MINIMIZATION}$$

'MDL' PRINCIPLE, KOLMOGOROV COMPLEXITY, AKAIKE INFO CRITERION, ALTERNATIVE TO BIC $\rightarrow AIC = \log P(D|\hat{\theta}_{MLE}) - \text{DOF}(M)$

WHEN DOING MARGINAL LIKELIHOODS PRIOR IS IMPORTANT AS WE WEIGHT LIKELIHOODS W/ IT, IF WE DON'T KNOW PRIOR, WE PUT A PRIOR ON IT, HYPERPRIOR,

$$P(D|M) = \iint P(D|w)P(w|\alpha, M)P(\alpha, M)dw d\alpha \quad \text{USUALLY HYPERPRIOR IS UNINFORMATIVE}$$

$$\text{USUALLY } \alpha \text{ IS OPTIMIZED AND NOT INTEGRATED OUT} \quad P(D|M) = \int P(D|w)P(w|\hat{\alpha}, M)dw, \quad \hat{\alpha} = \text{ARGMAX}_{\alpha} P(D|\alpha, M) = \text{ARGMAX}_{\alpha} \int P(D|w)P(w|\alpha)dw$$

EMPIRICAL BAYES!

BAYES FACTORS

IF PRIOR IS UNIFORM \rightarrow MODEL SELECTION IS PICKING MODEL WITH HIGHEST MARG. LIKELIHOOD

BAYES FACTOR IS RATIO OF ML $BF = \frac{P(D|M_1)}{P(D|M_0)} = \frac{P(M_1|D)}{P(M_0|D)} / \frac{P(M_1)}{P(M_0)}$ • FOR COMPARING TWO MODELS, NULL/ALTERNATIVE

— IT'S A LIKELIHOOD RATIO SO WE CAN COMPARE MODELS OF DIFFERENT COMPLEXITY

— BAYESIAN ALTERNATIVE TO P-VALUE

— JEFFREY - LINDLEY PARADOX! USE PROPER PRIORS (INTEGRATING TO 1) WHEN DOING MODEL SELECTION

PRIORS

— NO STRONG BELIEFS ABOUT $\theta \rightarrow$ UNINFORMATIVE PRIORS, CAN BE NON-INTUITIVE TO FIGURE OUT; VALIDATE CONCLUSIONS VIA SENSITIVITY ANALYSIS

— JEFFREY'S PRIOR: TECHNIQUE FOR CREATING NON-INFORMATIVE PRIORS. ANY REPARAMETRIZATION OF $P(\phi)$ SHALL ALSO BE UNINFORMATIVE.

GENERALLY: $P(\theta) = P(\phi) \left| \frac{d\phi}{d\theta} \right|$, IF $P(\phi) \propto I(\phi)^{1/2}$ FISHER INFO. $I(\phi) = -E \left[\left(\frac{d \log P(X|\phi)}{d\phi} \right)^2 \right]$ CURVATURE OF EXPECTED LOG-LIKELIHOOD, MEASURE OF STABILITY OF MLE

$I(\theta) = -E \left[\left(\frac{d \log P(X|\theta)}{d\theta} \right)^2 \right] = I(\phi) \left| \frac{d\phi}{d\theta} \right|^2 \rightarrow I(\theta)^{1/2} = I(\phi)^{1/2} \left| \frac{d\phi}{d\theta} \right|$ TRANSFORMED PRIOR IS THE SAME

• JEFFREYS FOR MULTIPLE PARAMS: $P(\theta) \propto \sqrt{\det I(\theta)}$

• CAN ALSO BE MADE TRANSLATION INVARIANT OR SCALE INVARIANT, CAN BE IMPROPER AS LONG AS POSTERIOR IS PROPER

— ROBUST PRIOR! WE AREN'T CONFIDENT, NO MUCH INFLUENCE, HEAVY TAILS SO THINGS AREN'T CLOSE TO MEAN

— MIXTURES OF CONJUGATE PRIORS IS ALSO CONJUGATE, GOOD COMPROMISE BETWEEN COMPUTATIONAL TRFE. AND FLEXIBILITY

$$P(\theta) = \sum_M \underbrace{P(z=M)}_{\text{WEIGHT}} \underbrace{P(\theta|z=M)}_{\text{CONJUGATE}}$$

• ALSO POSTERIOR CAN BE WRITTEN AS MIXTURE OF CONJUGATES

HIERARCHICAL BAYES

WHAT IF WE DON'T KNOW PRIOR? PRIOR ON PRIOR! HIERARCHICAL MODEL $\eta \rightarrow \theta \rightarrow D$, COMMON IN GRAPHICAL MODELS

EMPIRICAL BAYES

LINE FORWARD

• MLE \rightarrow MAP \rightarrow MLE-II (EB) \rightarrow MAP II \rightarrow FULL BAYES

BAYESIANNESSE \rightarrow

BAYESIAN DECISION THEORY

ACTION SPACE A , LOSS $L(y, a)$, HOW COMFORTABLE a IS WITH STATE y **GOAL:** $\delta: X \rightarrow A$ DEVISE POLICY, OPTIMAL ACTION FOR EACH POSSIBLE INPUT

OPTIMAL = LOSS-MINIMIZING = $\delta(x) = \text{ARGMIN}_a E[L(y, a)]$ || UTILITY = $U = -L = \delta(x) = \text{ARGMAX}_a E[U(y, a)]$ MAXIMIZING

MAX EXPECTED UTILITY, RATIONAL BEHAVIOR

POSTERIOR EXPECTED LOSS: $P(a|x) = E_{P(y|x)}[L(y, a)] = \sum L(y, a)P(y|x)$ MINIMIZE \rightarrow BAYES ESTIMATOR $\delta(x): \text{ARGMIN}_a P(a|x)$

\rightarrow 0-1 LOSS $L(y, a) = 1(y \neq a) = \begin{cases} 0 & a=y \\ 1 & a \neq y \end{cases}$ $P(a|x) = P(a \neq y|x) = 1 - P(y|x)$ MAXIMIZE $y^* = \text{ARGMAX}_y P(y|x)$ POSTERIOR MODE, MAP

\bullet REJECT OPTION WHEN $P(y|x)$ IS VERY UNCERTAIN; RISK-AVERSE DOMAINS $U = C + 1$ IS REJECTION $\lambda_R = \text{COST REJECTION}$ $\lambda_S = \text{SUBST. ERROR}$

$L(y=, a=) = \begin{cases} 0 & I= \\ \lambda_R & I=C+1 \\ \lambda_S & \text{ELSE} \end{cases}$ REJECT IF MOST PROBABLE CLASS HAS $P < 1 - \frac{\lambda_R}{\lambda_S}$

L_2 LOSS, SQUARED ERROR

$L(y, a) = (y - a)^2$ $P(a|x) = E[(y - a)^2|x] = E[y^2|x] - 2aE[y|x] + a^2$

$\frac{\partial}{\partial a} P(a|x) \rightarrow \hat{y} = E[y|x] = \int y \cdot P(y|x) dy$ POSTERIOR MEAN, MIN. MEAN SQ. ERROR, MMSE

L_1 LOSS

$L(y, a) = |y - a| \rightarrow$ OPT. ESTIMATE IS POSTERIOR MEDIAN $P(y < a|x) = P(y > a|x) = 0.5$

\rightarrow IN SUPERVISED LEARNING LOSS OF ACTION δ WHEN STATE IS $\theta = L(\theta, \delta) = \sum_x \sum_y L(y, \delta(x)) P(x, y|\theta)$ GENERALIZATION ERROR

$\rightarrow P(\delta|\theta) = \int P(\theta|\theta) L(\theta, \delta) d\theta$

\rightarrow IN BINARY DECISION PROCESSES WITH F_{POS}/F_{NEG} COSTS

$P(\hat{y}=0|x) = L_{FN} \cdot P(y=1|x)$ PICK 1 IFF $\frac{P(y=1|x)}{P(y=0|x)} > \frac{L_{FP}}{L_{FN}}$

$P(\hat{y}=1|x) = L_{FP} \cdot P(y=0|x)$

CONFUSION MATRIX

	$y=1$	$y=0$
$\hat{y}=1$	TP	FP
$\hat{y}=0$	FN	TN

$N_+ = TP + FP$
 $N_- = FN + TN$
 $N_+ = TP + FN$
 $N_- = FP + TN$

$TP/N_+ = \text{TPR, SENSITIVITY, RECALL}$
 $FN/N_+ = \text{FNR, MISS RATE, TYPE II}$
 $FP/N_- = \text{FPR, TYPE I}$
 $TN/N_- = \text{TNR, SPECIFICITY}$

\bullet ON DIFFERENT THRESHOLDS τ TO PICK OPTIMAL

- ROC CURVES

PLOT TPR vs FPR FOR DIFFERENT τ , MEASURE USING AUC OR EQUAL ERROR RATE EER FPR = TPR @ τ

- PRECISION / RECALL CURVES

FOR RARE EVENTS ROC IS NOT MUCH INFORMATIVE

$$\text{PRECISION} = TP / \hat{N}_+$$

$$\text{RECALL} = TP / N_+$$

$$P = \frac{\sum y_i \hat{y}_i}{\sum \hat{y}_i} \quad R = \frac{\sum y_i \hat{y}_i}{\sum y_i}$$

F-SCORE: $\frac{2PR}{R+P}$, HARMONIC MEAN

FOR SINGLE THRESHOLD

- FALSE DISCOVERY RATES, MANY BINARY DECISIONS $P(y_i = 1 | D) > \tau$, $D = \{x_i\}_1^N$, IS MULTIPLE HYPOTHESIS TESTING

HOW TO SET τ ? \rightarrow MINIMIZE FALSE POSITIVES

$$FD(\tau, D) = \sum_i \underbrace{(1 - p_i)}_{\text{ERR}} \cdot \underbrace{I(p_i > \tau)}_{\text{DISCOVERY}}$$

$$FOR(\tau, D) = \frac{FD(\tau, D)}{N(\tau, D)}$$

CHOOSE τ