

FREQUENTIST STATISTICS

- ESTIMATORS, MONTE-CARLO BOOTSTRAPPING FOR PARAMETERS, SIMILAR TO SAMPLING FROM BAYESIAN POSTERIOR WHEN PRIORS NOT VERY STRONG
- VARIANCE OF MLE, FISHER INFO MATRIX
CURVATURE OF LIKELIHOOD SURFACE
SCORE: $S(\hat{\theta}) = \nabla \log P(D|\theta) \big|_{\hat{\theta}}$
OBSERVED INFORMATION MATRIX: $J(\hat{\theta}(D)) = -\nabla^2 S(\hat{\theta}) = -\nabla^2 \log P(D|\hat{\theta})$
- HESSIAN OF NLL
- GRADIENT OF NEGATIVE SCORE FCN
- FISHER INFO MATRIX $I_N(\hat{\theta}|D) = E_{\theta} [J(\hat{\theta}|D)]$, $I_N(\hat{\theta}) = N \cdot I_1(\hat{\theta})$
- SAMPLING DISTRIBUTION OF MLE IS ASYMPTOTICALLY NORMAL
- STD ERROR $\hat{S}_2 = I_N(\hat{\theta})^{-1/2}$

FREQUENTIST DECISION THEORY

NO PRIOR, POSTERIOR, OR EXPECTED LOSS \rightarrow NO AUTOMATIC WAY FOR OPTIMAL ESTIMATION

RISK/EXPECTED LOSS $R(\hat{\theta}, \delta) = E_{P(D|\theta)} [L(\theta^*, \delta(\tilde{D}))] = \int L(\theta^*, \delta(\tilde{D})) P(\tilde{D}|\theta^*) d\tilde{D}$

EXPECTATION WRT SAMPLING DISTRIBUTION

- FREQ. RISK AVERAGES OVER D AND CONDITIONS ON θ^*
- BAYESIAN EXPECTED LOSS AVG'S OVER θ AND CONDITIONS ON D
- FREQUENTIST MAKES NO SENSE $\rightarrow \theta^*$ IS UNKNOWN

BAYES RISK

CONVERTS $R(\theta^*, \delta)$ INTO $R(\theta)$ \rightarrow PRIOR ON θ^* $R_B = E_{P(\theta^*)} [R(\theta^*, \delta)] = \int R(\theta^*, \delta) P(\theta^*) d\theta^* \rightarrow \delta_B = \text{ARGMIN } R_B(\delta)$

ALSO INTEGRATED/PREPOSTERIOR RISK

THEOREMS:

- A BAYES ESTIMATOR CAN BE OBTAINED BY MINIMIZING POSTERIOR EXPECTED LOSS FOR EACH X
- EVERY ADMISSIBLE DECISION RULE IS A BAYES DECISION RULE WRT SOME, EV. IMPROPER, PRIOR.
- BEST WAY FOR MINIMIZING FREQUENTIST RISK IS TO BE BAYESIAN
- PICKING OPTIMAL ACTION CASE-BY-CASE (BAYESIAN) IS OPTIMAL ON AVERAGE (FREQUENTIST)

MINIMAX RISK

$R_{\max}(\delta) = \max_{\theta^*} R(\theta^*, \delta) \rightarrow \delta_{\text{MM}} = \text{ARGMIN } R_{\max}(\delta)$ - VERY PESSIMISTIC

- NOT USED XCEPT IN GAME-THEORETIC SCENARIOS B/C NATURE IS NOT ADVERSARIAL

ADMISSIBLE ESTIMATORS:

IF $R(\theta, \delta_1) \leq R(\theta, \delta_2) \forall \theta \rightarrow \delta_1$ DOMINATES δ_2 , AN ESTIMATOR IS ADMISSIBLE IF NOT STRICTLY DOMINATED BY ANY OTHER

- SAMPLE MEDIAN IS USUALLY BETTER THAN SAMPLE MEAN

STEIN'S PARADOX

WHEN 3 OR MORE PARAMS ARE ESTIMATED SIMULTANEOUSLY, COMBINED ESTIMATORS ARE ON AVG MORE ACCURATE THAN SEPARATE ON AVG

- CONSTANTS CAN BE \neq ADMISSIBLE ESTIMATORS, NOT ENOUGH THEN

CONSISTENT ESTIMATORS $\hat{\theta}(D) \rightarrow \theta^*, |D| \rightarrow \infty$

UNBIASED ESTIMATORS:

$\text{BIAS}(\hat{\theta}(\cdot)) = E_{P(D|\theta^*)} [\hat{\theta}(D) - \theta^*] = 0$ SAMPLING DISTRIBUTION IS CENTERED ON θ^*

MVUE $\text{VAR}[\hat{\theta}] \geq \frac{1}{4I(\theta)}$ CR BOUND

BIAS-VARIANCE TRADEOFF:

$\text{MSE} = \text{VARIANCE} + \text{BIAS}^2$; MIGHT MAKE SENSE TO USE BIASED EST. IF REDUCES VARIANCE

- IF 0-1 LOSS (CLASSIFICATION) BIAS AND VARIANCE COMBINE MULTIPLICATIVELY

IE, RIDGE REGRESSION, $\lambda \rightarrow 0$ RESULT IN BIAS

EMPIRICAL RISK MINIMIZATION

RISK RELIES ON TRUE DATA DISTRIBUTION

FREQUENTIST RISK CAN BE COMPUTED WHEN TASK IS ESTIMATING OBSERVABLE QUANTITIES AND NOT HIDDEN VARS/PARAMS.

$$L(\theta, S(D)) \longrightarrow L(y, \delta(x))$$

$$R(P_X, \delta) = E_{(x,y) \sim P_X} [L(y, \delta(x))] = \sum_x \sum_y L(y, \delta(x)) P_X(x,y)$$

P_X IS NATURAL, EMPIRICAL DISTRIBUTION

$$\text{EMPIRICAL RISK} = R_{\text{EMP}}(D, \delta) = R(P_{\text{EMP}}, \delta) = \frac{1}{N} \sum_{i=1}^N L(y_i, \delta(x_i))$$

- 0-1 LOSS \rightarrow MISCLASSIFICATION RATE
- L_2 LOSS \rightarrow MSE

$$\delta_{\text{EMP}}(D) = \text{ARGMIN}_{\delta} R_{\text{EMP}}(D, \delta)$$

- IN UNSUPERVISED LEARNING

$$L(y, \delta(x)) \rightarrow L(x, \delta(x)) \quad \text{RECONSTRUCTION ERROR}$$

$$\delta(x) = \text{DECODE}(\text{ENCODE}(x))$$

- IF PRIOR IS 'NATURAL' IS EQUAL TO EMPIRICAL \rightarrow OVERFITTING

$$\text{COMPLEXITY PENALTY: } R(D, \delta) = R_{\text{EMP}}(D, \delta) + \lambda C(\delta)$$

$\lambda \rightarrow$ PEN OF COMPLEXITY

- PICK C AS DOF OR VC DIMENSION

- PICK $\lambda \rightarrow \hat{\lambda} = \text{ARGMIN}_{\lambda} \hat{R}(\hat{S}_{\lambda}) \rightarrow$ VIA CROSS-VALIDATION (AVG RISK OVER FOLDS) NO TRY DIFFERENT λ VALUES
 \rightarrow VIA THEORETICAL UPPER BOUND

• VIA CV NOT FEASIBLE WITH > 2 PARAMS, EMPIRICAL BAYES W/ GRADIENT BASED OPTIMIZERS, NOT FEASIBLE FOR UNSUPERVISED LEARNING

• ONE-STANDARD ERROR RULES: PICK SIMPLEST MODEL W/ RISK NO MORE THAN 1 STD ERR ABOVE RISK OF BEST MODEL

• UPPER BOUND VIA STATISTICAL LEARNING THEORY (SLT)

USING CV IS INEFFICIENT. LET'S BOUND THE RISK FOR ANY DISTRIBUTION AND HYPOTHESIS SPACE WITH R_{EMP} , SAMPLE SIZE, AND $|H|$

$$- |H| \text{ IS FINITE } P(\max |R_{\text{EMP}}(D, h) - R(P_X, h)| > \epsilon) \leq 2|H|e^{-2N\epsilon^2}$$

- $|H|$ IS INFINITE (REAL VALUES) USE VC DIMENSION, NOT ALWAYS EASY TO COMPUTE, BOUNDS ARE LOOSE

• COLT COMPUTATIONAL LEARNING THEORY, TAKES COMPUTATIONAL COMPLEXITY OF LEARNER INTO ACCOUNT

• $|H|$ IS PAC LEARNABLE IF THERE'S A P-TIME ALGORITHM IDENTIFYING A FUNCTION PAC (PROBABLY APPROXIMATELY CORRECT)

SURROGATE LOSS FUNCTIONS

FOR OPTIMIZING 0-1 LOSS, OR OTHER METRICS (AUC, F1). USE MLE. CONSTRUCTS A DECISION FUNCTION AND USES LOSS ON THAT.

LOG-LOSS: $L_{\text{NLL}}(y, \eta) = -\log p(y|x, w) = \log(1 + e^{-y\eta})$, η IS LOG ODDS RATIO, DECISION FEN | MINIMIZING LOGLOSS = MAXIMIZING LIKELIHOOD

$$\text{HINGE LOSS: } L_{\text{HINGE}}(y, \eta) = \max(0, 1 - y\eta)$$

FREQUENTIST PATHOLOGIES

CONFIDENCE INTERVALS CONDITION ON UNKNOWN θ AND AVG ON FUTURE DATA D

P-VALUES ONLY USEFUL TO REJECT THE NULL HYPOTHESIS \rightarrow NEVER GATHER EVIDENCE IN FAVOR OF IT, SENSITIVE TO STOPPING RULE
LEAD TO NON EARLY TERMINATION, DIFFERENT P-VALUES FOR SAME SITUATION ON BINOMIAL AND NEG-BINOMIAL!

LIKELIHOOD PRINCIPLE IS VIOLATED: INFERENCE SHOULD BE BASED ON LIKELIHOOD OF OBSERVED DATA, AND ON EVENTS WHICH ACTUALLY HAPPENED, NOT WHICH MIGHT HAVE HAPPENED