# LINEAR REGRESSION

## GENERAL MODEL: $P(y|x,\theta) = N(y|w^Tx, \sigma^2)$

BASIS FUNCTION EXPANSION: $P(y|x,\theta) = N(y|w^T\phi(x), \sigma^2)$

STILL LINEAR WRT PARAMETERS

HIGHER DEGREES OF $\phi(x) \longrightarrow$ MORE COMPLEX FUNCTIONS

## MLE ESTIMATION

$\hat{\theta} = \text{ARGMAX } \log P(D|\theta)$. SAMPLES ARE ASSUMED IID. $\ell(\theta) = \sum_{i}^{N} \log P(y_i|x_i,\theta)$, BUT WE MINIMIZE NLL BECAUSE EQUIVALENT AND EASIER

— PLUG-IN GAUSSIAN FORM

$$\ell(\theta) = -\frac{1}{2\sigma^2} \sum_{i}^{N} (y_i - w^Tx_i)^2 - \frac{N}{2}\log(2\pi\sigma^2)$$

$\underbrace{\qquad\qquad\qquad}_{RSS\ (SSE)}$ , $\frac{SSE}{N} = MSE$ , $RSS(w) = ||\epsilon||_2^2 = \sum_{1}^{N}\epsilon_i^2$

— LEAST SQUARES BECAUSE MLE FOR $w$ MINIMIZES RSS      — NLL SURFACE IS QUADRATIC BOWL, CONVEX, UNIQUE MINIMUM, WE DERIVE

— DERIVATION $NLL(w) = \frac{1}{2} w^T(X^TX)w - w^T(X^Ty)$ . $X^TX = \sum_{i}^{N} x_i x_i^T = $ SUM-OF-SQUARES MATRIX

$$X^Ty = \sum_{i}^{N} x_i y_i$$

GRADIENT: $G(w) = [X^TXw - X^Ty] = \sum_{i}^{N} x_i(w^Tx_i - y_i) = 0 \longrightarrow X^TXw = X^Ty$

$$\hat{w}_{OLS} = (X^TX)^{-1} X^Ty$$

— GEOMETRIC INTERPRETATION

$N$ EXAMPLES, $D$ FEATURES, $N > D$. COLS OF $X$ ARE LINEAR SUBSPACE OF $|D|$ EMBEDDED IN $N$ DIMENSIONS.

WE SEEK $\hat{y} \in R^n$ WHICH LIES IN COLUMN LINEAR SUBSPACE AND IS CLOSER AS POSSIBLE TO $y = $ ARGMIN $||y - \hat{y}||_2$

IS ON BECAUSE $\hat{y} \in SPAN(x)$, $w$ EXISTS. $\longrightarrow$ RESIDUAL VECTOR TO BE ORTHOGONAL TO EVERY COLUMN IN $X$

$$\bar{x}_j^T(y - \hat{y}) = 0 \longrightarrow X^T(y - Xw) = 0 \longrightarrow \hat{w} = (X^TX)^{-1}X^Ty$$

$$\hat{y} = X\hat{w} = X(X^TX)^{-1}X^Ty \longrightarrow \text{ORTHOGONAL PROJECTION OF } y \text{ ON COLSPACE OF } X$$

$\underbrace{\qquad\qquad\qquad}_{\text{PROJECTION / HAT MATRIX}}$

— CONVEXITY

• A SET IS CONVEX IF WE DRAW LINE BETWEEN TWO POINTS AND LINE ALWAYS LIES INSIDE THE SET

• A FUNCTION IS CONVEX IF EPIGRAPH (SET OF PTS ABOVE FCN) IS CONVEX      • CONCAVE $\longrightarrow -f(\theta)$ IS CONVEX

• UNIQUE GLOBAL MINIMUM $\longrightarrow$ SECOND DERIVATIVE ALWAYS POSITIVE $\longrightarrow$ A TWICE-DIFF. $C_2$ CONTINUOUS, MULTIVAR FCN IS CONVEX IFF HESSIAN IS POSITIVE DEFINITE FOR ALL $\theta$

— ROBUST LINEAR REGRESSION

OLS IS VERY SENSITIVE TO OUTLIERS BECAUSE QUADRATIC LOSS $\longrightarrow$ HIGHER IMPACT. REPLACE GAUSSIAN RESPONSE W/ STH MORE HEAVY TAIL

I.E. LAPLACE DISTRIBUTION $P(y|x, w, b) = LAP(y|x^Tw, b) \propto EXP(-\frac{1}{b}|y - w^Tx|)$     $\ell(w) = \sum |R_i(w)|$ RESIDUALS

$\underbrace{\qquad}_{R}$     NONLINEAR!

• OPTIMIZE NLL (SPLIT VARIABLE TRICK)     • OPTIMIZE HUBER LOSS

$$L_H(R, \delta) \begin{cases} R^2/2 & |R| \le \delta & \ell_2 \\ \delta|R| - \delta^2/2 & |R| > \delta & \ell_1 \end{cases}$$ EVERYWHERE DIFF, $C_1$, FASTER 2 OPTIMIZE BC SMOOTH OPT. METHODS

# VARIANTS OF LINEAR REGRESSION

|  | LIKELIHOOD | PRIOR |
|---|---|---|
| LEAST SQUARES | GAUSSIAN | UNIFORM |
| RIDGE | GAUSSIAN | GAUSSIAN |
| LASSO | GAUSSIAN | LAPLACE |
| ROBUST | LAPLACE | UNIFORM |
| ROBUST | STUDENT | UNIFORM |

## RIDGE REGRESSION

RESILIENT TO OVERFITTING. GAUSSIAN PRIOR. $P(w) = \prod_j N(w_j | 0, \tau^2)$ ENCOURAGES PARAMS TO BE SMALL

- ARGMAX $\sum_j^N \log N(y_i | w_0 + w^T x_i, \sigma^2) + \sum^D \log N(w_j | 0, \tau^2)$ $\xrightarrow{\text{MINIMIZE}}$ $\overbrace{J(w) = \underbrace{\frac{1}{N} \sum (y_i - (w_0 + w^T x_i))^2}_{\text{MLE/NLL}} + \underbrace{\lambda ||w||_2^2}_{\text{PENALTY}}}$

$\hat{w}_{RIDGE} = (\lambda I_0 + X^T X)^{-1} X^T y$

$\lambda = \sigma^2 / \tau^2$, $||w||_2^2 = w^T w$ (SQUARED TWO NORM)

- GAUSSIAN PRIOR $\longrightarrow$ $\ell_2$ REGULARIZATION / WEIGHT DECAY. PENALIZES SUM OF MAGNITUDE OF $w$s

- $\lambda$ PICKED WITH CV

### — COMPUTATIONAL EFFICIENCY

$(\lambda I_0 + X^T X)$ IS BETTER CONDITIONED / MORE LIKELY TO BE INVERTED BUT 4 NUMERICAL STABILITY IT'S BETTER NOT TO INVERT MATRICES ALTOGETHER

AUGMENT X WITH DATA FROM PRIOR $\tilde{X} = \begin{pmatrix} X/\sigma \\ \sqrt{\lambda} \end{pmatrix}$ $\tilde{y} = \begin{pmatrix} y/\sigma \\ 0_{D \times 1} \end{pmatrix}$ $\hat{X} = QR$, Q IS ORTHONORMAL $(Q^T Q = QQ^T = I)$ AND R IS UPPER TRIANGULAR (EASY TO INVERT)

- QR DECOMP.

   $w_{RIDGE} = R^{-1} R^{-T} R^T Q^T \tilde{y} = R^{-1} Q^T \tilde{y}$

   OK EVEN IF X, y (IMPLEMENTATION ALL USE QR DECOMPOSITION) $O(ND^2)$

- $D >> N$ DO SVD DECOMPOSITION FIRST. $w_{RIDGE} = V(Z^T Z + \lambda I_N)^{-1} Z^T y$ REPLACE X, (D-DIM) WITH Z, (N-DIM), THEN RETRANSFORM TO |D| WITH V, $O(DN^2)$

### — SHRINKAGE

RELATION BETWEEN RIDGE PREDICTIONS AND SINGULAR VALUES OF X (VIA SVD)   $DOF(\lambda) = \sum^D \frac{\sigma^2}{\sigma_j^2 + \lambda}$   $\lambda = 0 \to D$

SMALLER SINGULAR VALUES ARE DIRECTIONS WITH HIGHER POSTERIOR VARIANCE $\to$ MOST SHRUNKED   $\lambda = \infty \to 0$

(THOSE ARE THE ONES WE $\overset{\text{RIDGE}}{\text{AXE}}$) SINGULAR VALUES ARE EIGENVECTORS OF $X^T X$.   OBS PCA

- CHOLESKY DECOMPOSITION $\Lambda = \sqrt{\Lambda} \cdot \sqrt{\Lambda}^T$
- QR DECOMPOSITION $X = QR$, Q ORTHONORMAL, R UPPER TRIANGULAR
- SVD DECOMPOSITION $X = U \Sigma V^X$, $V^T V = I$, $UU^T = U^T U = I$, $\Sigma$ DIAG, $Z = UD$

# LMS ALGORITHM (ONLINE LINEAR REGRESSION)

ALSO DELTA RULE / WIDROW - HOFF RULE

$$y_n = x_i(\theta_n^T x_i - y_i) \longrightarrow \text{GRADIENT ACTS AS ERROR SIGNAL}$$

$$\theta_{n+1} = \theta_n - \eta_n(\hat{y}_n - y_n)x_n$$

NO PROJECTION STEP BECAUSE UNCONSTRAINED.

USUALLY $0.1 < \eta < 0.4$

# LINEAR SEPARABILITY

$W^T \perp$ TO DECISION BOUNDARY   MARGIN: DISTANCE BETWEEN (OPTIMAL) HYPERPLANE AND ANY DATAPOINT

# BAYESIAN LIN REG

FULL POSTERIOR OVER $W$ AND $\sigma^2$

- $\sigma^2$ KNOWN

  - POSTERIOR: $P(w \mid X, y, \sigma^2) = N(w \mid w_N, V_N)$ $\quad$ [ LIKELIHOOD: $P(y \mid X, w, \mu, \sigma^2) = N(y \mid \mu + Xw, \sigma^2 I)$

    $w_N = V_N V_0^{-1} w_0 + \frac{1}{\sigma^2} V_N X^T y$ $\qquad$ CONJ. PRIOR: $P(w) = N(w \mid w_0, V_0)$ ]

    $V_N = \sigma^2 (\sigma^2 V_0^{-1} + X^T X)^{-1}$ $\quad$ IF $w_0 = 0$ AND $V_0 = \tau^2 I$ $\longrightarrow$ RIDGE ESTIMATE

  - POSTERIOR PREDICTIVE

    $$P(y \mid x, D, \sigma^2) = \int N(y \mid x^T w, \sigma) N(w \mid w_N, V_N) \, dw = N(y \mid w_N^T x, \sigma_N^2(x)) \quad, \ \sigma_N^2 = \sigma^2 + X^T V_N X$$

    $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ DEPENDS ON HOW CLOSE $x$ IS TO TRAINING DATA

- $\sigma^2$ UNKNOWN

  ~~BORING TOTAGS~~ MURPHY PP. 234

# EMPIRICAL BAYES ( EVIDENCE PROCEDURE )

PICK THE HYPERPARAMETERS OF THE PRIOR $\eta = (\alpha, \beta)$ TO MAXIMIZE MARGINAL LIKELIHOOD $\lambda = 1/\sigma^2$ PRECISION OF NOISE, $\alpha$ PRECISION OF PRIOR $P(w) = N(w \mid 0, \alpha^{-1}, I)$. ALTERNATIVE TO CROSS. VALIDATION. • BETTER BECAUSE E B ENABLES COMPUTING OF DIFFERENT $\alpha_i$ FOR EVERY FEATURE $\longrightarrow$ FEATURE SELECTION VIA ARD ( AUTOMATIC RELEVANCY DETERMINATION ) . IMPOSSIBLE WITH CV

- USEFUL TO COMPARE DIFFERENT KINDS OF MODELS

$$P(D \mid m) = \iint P(D \mid w, m) P(w \mid m, \eta) P(\eta \mid m) \, dw \, d\eta \approx \max_\eta \int P(D \mid w, m) P(w \mid m, \eta) P(\eta \mid m) \, dw$$

# ARD: