# Logistic Regression

IS DISCRIMINATIVE CLASSIFIER, ONLY FITS THE MODEL, DIRECTLY FIT $P(Y|X)$, NO JOINT $P(X,Y)$ MODEL

MODEL: $P(y|x,w) = BER(y|SIGM(w^T x))$

MLE ESTIMATE: $NLL(w) = -\frac{1}{N}\sum \log\left[\mu_i^{I(y_i=1)} \times (1-\mu_i)^{I(y_i=0)}\right] = -\sum_{i}^{N}\left[y_i \log\mu_i + (1-y_i)\log(1-\mu_i)\right]$

$\beta(u) = \hat{y} \text{ st } \hat{y} - y$
CROSS ENTROPY !
ERROR FCN

• $NLL = \sum^{N} \log\left(1 + EXP(-\tilde{y}_i w^T x_i)\right)$ — NO CAN WRITE IN CLOSED FORM → OPTIMIZATION ALGORITHM, WE NEED GRADIENT AND HESSIAN

$g(w) = \frac{d}{dw} f(w) = \sum_i (\mu_i - y_i)x_i = X^T(\mu - y)$

$H(w) = \frac{d}{dw} g(w)^T = \sum_i (\nabla_w \mu_i)x_i^T = \sum_i \mu_i(1-\mu_i)x_i x_i^T = X^T S X$ ; $S = DIAG(\mu_i(1-\mu_i))$, $H$ IS POSITIVE DEFINITE OFC SO NLL IS CONVEX, YAY!

## GRADIENT DESCENT (STEEPEST DESCENT)

$\theta_{n+1} = \theta_n - \eta_n g_n$ , $\eta_n$ LEARNING RATE ; HOW DO WE SET IT?

— CONSTANT → TOO LOW, TAKES LONG TIME
→ TOO LARGE, FAILS TO CONVERGE, OSCILLATES

— LINE SEARCH

GUARANTEED TO CONVERGE (GLOBAL CONVERGENCE) NO MATTER WHERE IT STARTS

$f(\theta + \eta d) \approx f(\theta) + \eta g^T d$ , TAYLOR; $d$ IS DESCENT DIRECTION. WE WANT IT SMALL SO THAT $f(\theta + \eta d) < f(\theta)$ BUT NOT TOO SMALL. $\eta$ TO MINIMIZE: $\phi(\eta) = f(\theta_n + \eta d_n)$. $\phi'(\eta) = 0$ || $\phi'(\eta) = d^T g$ BY CHAIN RULE

• ZIG-ZAGS

— $g = 0$ : WE ARRIVED
— $g \perp d$ : STEP STOPS WHERE GRADIENT IS $\perp$ TO SEARCH DIRECTION

TO MINIMIZE ZIGZAGGING

• MOMENTUM $\theta_{n+1} = \theta_n - \eta_n g_n + \mu_n(\theta_n - \theta_{n-1})$ , $0 \leq \mu_n \leq 1$
(HEAVY BALL METHOD)

• CON)UGATE GRADIENTS QUADRATIC OBJECTIVES $f(\theta) = \theta^T A \theta$

## REGULARIZATION

ACTUAL MLE IS WHEN $\|w\| \to \infty$, INFINITELY STEEP SIGMOID $I(w^T x > w_0)$, BRITTLE SOLUTION, POOR GENERALIZATION

$L_2$ REGULARIZE: $f'(w) = NLL(w) + \lambda w^T w$      AND PLUG IN GRADIENT OPTIMIZER

$g'(w) = g(w) + \lambda w$

$H'(w) = H(w) + \lambda I$

# NEWTON'S METHOD

SECOND ORDER METHOD, TAKES CURVATURE OF SPACE INTO ACCOUNT (HESSIAN), FASTER $\qquad \theta_{n+1} = \theta_n - \eta_n H_n^{-1} g_n$ . $\boxed{\text{STEP} \quad d = -H^{-1} g_n}$

MINIMIZES $2^{ND}$ ORDER

APPROXIMATION OF F

## STEPS

— UNTIL CONVERGENCE, STEP = K

— $g_n = \nabla f(\theta_n)$ , $H_n = \nabla^2 f(\theta_n)$

— $H_n d_n = -g_n$ , SOLVE FOR $d_n$

— LINE SEARCH TO MINIMIZE $\eta_n$ ON $d_n$

— $\theta_{n+1} = \theta_n + \eta_n d_n$

• IF $H_n$ IS NOT POSITIVE DEFINITE, FCN NOT CONVEX, GO BACK TO STEEPEST DESCENT

### LEVENBERG - MARQUARDT ADAPTS BETWEEN NEWTON AND S.D. STEPS

— $\theta_{n+1} = \theta_n - \left( H_n + \lambda \, \text{DIAG}(H_n) \right)^{-1} g_n$

$\begin{cases} \lambda >> \longrightarrow \text{VANILLA GRAD DESCENT} \\ \lambda << \longrightarrow \text{NEWTON'S} \end{cases}$

• ERROR INCREASE $\rightarrow$ REJECT UPDATE STEP AND INCREASE $\lambda$

ERROR DECREASE $\rightarrow$ ACCEPT UPDATE AND DECREASE $\lambda$

# IRLS - ITERATIVELY REWEIGHTED LEAST SQUARES

NEWTON'S TO FIND MLE FOR BINARY LOGISTIC REGRESSION $\qquad W_{n+1} = W_n - H^{-1} g_n = (X^T S X)^{-1} X^T S_n z_n$

HESSIAN IS EXACT SO $\eta_n = 1$, COOL.

$$z_n = X w_n + S^{-1}(y - \mu_n) \quad \text{WORKING RESPONSE}$$

IS WEIGHED LEAST SQUARE MINIMIZING $\sum S_{ni} (z_{ni} - W^T x)^2$ S IS DIAGONAL $\longrightarrow$ $z_{ni} = W_n^T x_i + \dfrac{y_i - \mu_{ni}}{\mu_{ni}(1 - \mu_{ni})}$

## STEPS

• $W_* = \theta_0$

• $W_0 = \log\left( \bar{y} / (1 - \bar{y}) \right)$

=

$\eta_i = W_0 + W^T x_i \qquad\qquad z_i = \eta_i + \dfrac{y_i - \mu_i}{S_i}$

$\mu_i = \text{SIGM}(\eta_i) \qquad\qquad S = \text{DIAG}(S_{i:N})$

$S_i = \mu_i(1 - \mu_i) \qquad\quad W = (X^T S X)^{-1} X^T S z$

UNTIL CONVERGENCE

# BFGS

IS QUASI - NEWTON METHODS, $H$ CAN BE EXPENSIVE TO COMPUTE EXPLICITLY. APPROXIMATE $H$ USING GRADIENT AT EACH STEP

$B_n \approx H_n$

$$B_{n+1} = B_n + \frac{y_n y_n^T}{y_n^T s_n} - \frac{(B_n s_n)(B_n s_n)^T}{s_n^T B_n s_n} \quad \text{BFGS FORMULA}$$

$s_n = \theta_n - \theta_{n-1} \qquad y_n = g_n - g_{n-1}$

• $B_0 = I$

ENSURES MATRIX REMAINS POSITIVE DEF.

|DIAGONAL + LOW RANK APPROXIMATION|

• ALTERNATIVE! BFGS APPROXIMATES $C_n \approx H_n^{-1}$
  INVERSE HESSIAN RIGHT AWAY

• $O(D^2)$ SPACE

# L - BFGS

LIMITED MEMORY VERSION

IS ACTUALLY DIAGONAL + LOW RANK, USES ONLY M MOST RECENT $(s_n, y_n)$ • $O(mD)$ • $m \sim 20$

OBS: BROYDEN FAMILY FORMULAS: $H_{n+1} = (1 - \phi) H_{n+1}^{DFP} + \phi H_{n+1}^{BFGS}$ H DFP IS SIMILAR BUT LESS ROBUST THAN BFGS

# MULTI-CLASS LOGISTIC REGRESSION

OR MAX-ENTROPY CLASSIFIER $\qquad P(y=c \mid x, w) = \dfrac{\exp(w_c^T x)}{\sum_c \exp(w^T x)}$

- CONDITIONAL LOGIT MODEL NORMALIZES OVER DIFFERENT CLASSES FOR EACH DATACASE

- $\mu_{ic} = P(y_i = c \mid x, w_i) = S(\eta_i)_c$ , $\eta_i = w^T x_i$ is $C \times 1$ VECTOR , $y_{ic} = I(y_i = c)$ ONE OF C ENCODING , BIT = 1 IFF $y_i = c$

- $\ell(w) = \log \prod_i^N \prod_c^C \mu_{ic}^{y_{ic}} = \sum_i^N \sum_c^C y_{ic} \lg \mu_{ic} = \sum_i^N \left[ \left( \sum_c^C y_{ic} w_c^T x_i \right) - \log \left( \sum_c^C \exp(w_c^T x_i) \right) \right]$ $\qquad$ • NLL $= -\ell(w)$

GRADIENT: $\quad G(w) = \nabla f(w) = \sum_i^N (\mu_i - y_i) \otimes x_i$

$y_i = \left( I(y=1) \ldots I(y=C-1) \right)$ $\quad \mu_i(w) = \left[ P(y_i = 1 \mid x_i, w_i) \ldots P(y_i = C-1 \mid x_i, w_i) \right]$

$A \otimes B = \begin{bmatrix} a_{11} B \ldots \ldots a_{1N} B \\ \vdots \qquad \qquad \vdots \\ a_{m1} B \ldots \ldots a_{mN} B \end{bmatrix}$ KRONECKER PRODUCT

$G(w) = \sum_i \begin{cases} (\mu_{i1} - y_{i1}) x_{i1} \\ (\mu_{i1} - y_{i1}) x_{i2} \\ \vdots \\ \vdots \end{cases}$ $\quad \nabla_{w_c} f(w) = \sum_i (\mu_{ic} - y_{ic}) x_i$

SAME FORM AS PRIMARY LOGISTIC REGRESSION ; ERROR TIMES $x_i \longrightarrow$ BECAUSE IS GLM

HESSIAN: $\quad H(w) = \nabla f(w) = \sum_i^N \left( DIAG(\mu_i) - \mu_i \mu_i^T \right) \otimes \left( x_i x_i^T \right)$ . ALSO BLOCK, $H_{cc'}(w) = \sum_i \mu_{ic} \left( \delta_{c,c'} - \mu_{ic'} \right) x_i x_i^T$

IS POSITIVE DEFINITE $\longrightarrow$ UNIQUE MLE

MINIMIZING $\quad f'(w) = -\log P(D \mid w) - \lg P(w)$ , $\quad P(w) = \prod_c N(w_c \mid 0, V_0)$

- $f'(w) = f(w) + \frac{1}{2} \sum_c w_c V_0^{-1} w_c$

HESSIAN IS $O\left( (CD) \times (CD) \right)$ SO HERE BFGS IS THUMBS UP!

- $g'(w) = g(w) + V_0^{-1} \left( \sum_c w_c \right)$

- $H'(w) = H(w) + I_C \otimes V_0^{-1}$

# Bayesian Logistic Regression

We want full posterior over params $P(w|D)$, not convenient because L.R. has no conjugate prior. Approximations i.e. MCMC, variational inference, expectation propagation. But later.

## Laplace Approximation

Let's approximate the posterior to a Gaussian $P(\theta|D) = \frac{1}{Z} e^{-E(\theta)}$  $E(\theta)$ energy function  $E(\theta) = -\log P(\theta, D)$, $Z$ norm constant

$$E(\theta) \approx E(\theta^*) + (\theta - \theta^*)^T g + \frac{1}{2}(\theta - \theta^*)^T H (\theta - \theta^*)$$

Taylor expansion around the mode

$$g = \nabla E(\theta)\big|_{\theta^*} \quad H = \frac{\partial^2 E(\theta)}{\partial \theta \, \partial \theta^T}\big|_{\theta^*}$$

mode $\to$ gradient is $0$

$$P(\theta|D) = N(\theta|\theta^*, H^{-1})$$

$$Z \approx \int \hat{P}(\theta|D) d\theta = e^{-E(\theta)}(2\pi)^{D/2}|H|^{-1/2} = P(D) \quad \leftarrow \text{Laplace approximation of ML}$$

also Gaussian approximation, saddle point approximation

$$\log P(D) \approx \log P(D|\theta^*) + \log P(\theta^*) - \frac{1}{2}\log|H|$$

Occam factor, measures model complexity. If uniform prior $P(\theta) \propto 1$, no $H$ term (no curvature) and $\theta^* \to \hat{\theta}_{MLE}$

$$\log P(D) \approx \log P(D|\hat{\theta}) - \frac{D}{2}\log N \quad \leftarrow \text{BIC score}$$

## Gaussian Approximation of LogReg

Gaussian prior $P(w|D) \approx N(w|\hat{w}|H^{-1})$  $\hat{w} = \text{argmin}_w E(w)$, $E(w) = -(\log P(D|w) + \log P(w))$, $H = \nabla^2 E(w)|_{\hat{w}}$

$P(w) = N(w|0, V_0)$  Issues when data is linearly separable $\to$ MLE not well defined, goes to infinity, sigmoid very steep.

Regularize w/ spherical prior $N(w|0, 100I)$. Better than MAP alone

## Posterior Predictive

$$P(y|x, D) = \int P(y|x, w) P(w|D) dw$$

intractable in this case. plug-in approx $P(y|1|x, D) \approx P(y=1|x, E[w])$  $E[w]$ posterior mean

Bayes point, underestimates uncertainty

— Monte Carlo Approximation  $P(y=1|x, D) \approx \frac{1}{S}\sum \text{sigm}((w^s)^T x)$, $w^s \sim P(w|D)$ samples from the posterior

if MC posterior $\to$ samples 4 prediction.   if Gaussian posterior $\to$ samples from it

— Probit Approximation

$P(w|D) \approx N(w|m_N, V_N)$  Gaussian approx posterior | $P(y=1, x, D) = \int \text{sigm}(a) N(a|\mu, \sigma_a^2) da$  posterior predictive

$= \frac{1}{2} x^T V_N x$  [Murphy P. 259]

Approx sigm with Probit = CDF of normal distribution

$\text{sigm}(u) \approx \Phi(\lambda a)$  analytically convolvable with Gaussians

$$\int \text{sigm}(a) N(a|\mu, \sigma^2) da \approx \text{sigm}(\kappa(\sigma^2)\mu), \quad \kappa(\sigma^2) = (1 + \pi\sigma^2/8)^{-1/2}$$

• $P(y=1|D, x) \approx \text{sigm}(\kappa(\sigma_a^2)\mu_a)$  $\leftarrow$ moderated output : less extreme than plug-in, but same decision boundary

# BAYESIAN LOGISTIC REGRESSION — OUTLIER DETECTION

- **USUALLY:** WITH RESIDUALS $R: y_i - \hat{y}_i$ SHOULD FOLLOW $N(0, \sigma^2)$, QQ-PLOT THEORETICAL vs EMPIRICAL QUANTILES

- **BINARY DATA:** STATISTICS NOT ASYMPTOTICALLY NORMAL. BAYESIAN $\rightarrow$ POINTS FOR WHEN $P(y|\hat{y})$ IS SMALL

  OUTLIERS: POINTS WITH LOW PROBABILITY UNDER X-VALIDATED POSTERIOR PREDICTIVE

$$P(y_i | x_i, x_{-i}, y_{-i}) = \int P(y_i | x_i, w) \prod_{i' \neq i} P(y_{i'} | x_{i'}, w) P(w) \, dw$$


# ONLINE LEARNING AND OPTIMIZATION

- **OFFLINE:** $f(\theta) = \frac{1}{N} \sum_i^N \ell(\theta, z_i)$, $z_i = \underset{\text{SUP}}{(x_i, y_i)}$ OR $z_i = \underset{\text{UNSUP}}{x_i}$, $\ell(z_i, \theta) = \text{LOSS}$, I.E. $-\log P(y_i | x_i, \theta)$ OR $L(y_i, h(x_i, \theta))$

- **ONLINE: REGRET MINIMIZATION** AVG LOSS RELATIVE TO THE BEST THAT COULD'VE GOTTEN IN HINDSIGHT WITH FIXED PARAMETERS

  REGRET: $\frac{1}{N} \sum_i^N \ell(\theta_t, z_t) - \min_{\theta \in \Theta} \frac{1}{N} \sum \ell(\theta_*, z_t)$   IS OBJECTIVE / LOSS

  ONLINE GRADIENT DESCENT $\theta_{n+1} = \text{PROJ}_\theta (\theta_n - \eta_n g_n)$ | $\text{PROJ}_v = \underset{w \in V}{\text{ARGMIN}} \|w - v\|_2$ IS PROJECTION OF $v$ ON $V$, $g_n = \nabla \ell(\theta_n, z_n)$, $\eta_n$ STEPSIZE

- **ONLINE: MINIMIZE EXPECTED LOSS IN THE FUTURE** $\langle$ STOCHASTIC OPTIMIZATION $\rangle$

  SOME VARS IN OBJECTIVE ARE RANDOM

$$f(\theta) = E[\ell(\theta, z)]$$

  **SGD:** $\bar{\theta}_n = \frac{1}{N} \sum_i^N \theta_t$ (RUNNING AVERAGE) $\xrightarrow{\text{ONLINE}}$ $\bar{\theta}_n = \bar{\theta}_{n-1} - \frac{1}{N}(\bar{\theta}_{n-1} - \bar{\theta}_n)$ POLYAK-RUPPERT AVERAGING

  HOW TO SET STEP SIZE: ROBBINS-MONRO CONDITIONS: $\sum \eta_n = \infty$, $\sum \eta_n^2 < \infty$ FOR CONVERGENCE

  $\eta_n = \frac{1}{N}$   OR   $\eta_n = (\tau_0 + n)^{-n}$

  $\underset{\geq 0}{\phantom{x}}$ $\underset{\substack{\text{FORCES OLD VALUES} \\ \text{CONTROL}}}{(0.5, 1]}$

  **HEURISTIC:** • TRY A RANGE OF $\eta$ VALUES
  • CHOOSE ONE W/ FASTEST DECREASE IN OBJECTIVE
  • APPLY TO REST OF DATA

  **DRAWBACKS:** MANUAL TUNING OF PARAMETERS
  SAME $\eta$ SIZE FOR ALL STEPS; SAME STEP FOR ALL PARAMS

  EARLY STOPPING: STOP AT PLATEAU, NOT NEC.LY CONVERGENCE

# ADAGRAD

ALINE TO DIAGONAL HESSIAN APPROXIMATION, PER PARAMETER STEP VALUE ADAPTING TO CURVATURE OF LOSS FCN

$$\theta_i(n+1) = \theta_i(n) - \eta \frac{g_i(n)}{\tau_0 + \sqrt{s_i(n)}} \qquad s_i(n) = s_i(n-1) + g_i(n)^2$$

EFFICIENT TO COMPUTE GRADIENT ON MINIBATCHES. $B=1 \rightarrow$ SGD
$B=N \rightarrow$ STEEPEST DESCENT
BECAUSE TAKES FEW STEPS TO DETERMINE DIRECTION. AM UNCERTAINTY HELPS AVOID LOCAL MINIMA

• ADADELTA, RMSPROP, ADAM

SGD STEPS
- INIT $\theta, \eta$
- REPEAT
  - PERMUTE DATA
  - $i = 1 : N$
    - $g = \nabla \ell(\theta, z)$
    - $\theta \leftarrow \text{PROJ}_\theta (\theta - \eta g)$
  - UPDATE $\eta$
- CONVERGENCE

# PERCEPTRON

ONLINE BINARY LOGISTIC REGRESSION $\theta_n = \theta_{n-1} - \eta_n g_n = \theta_{n-1} - \eta_n(\mu_i - y_i)x_i$ | $\mu_i = P(y=1|x_i,\theta_i) = E[y_i|x_i,\theta_n]$

HAS SAME FORM AS **LMS** BECAUSE GENERALIZED LINEAR MODELS

$\hat{y}_i = \text{ARGMAX } P(y|x_i,\theta)$ , $\mu_i = \text{SIGM}(\theta^T,x)$ , REPLACE W $\hat{y}$     $g \approx (\hat{y}_i - y)x_i \to y \in \{-1,1\} \to \hat{y} = \text{SIGN}(\theta^T x_i)$

UPDATE   NO CHANGE IF CLASSIFICATION IS RIGHT, UPDATE ONLY IF WRONG    $\theta_n = \theta_{n-1} + \eta_n y_i x_i$

- **CONVERGES** IFF LINEARLY SEPARABLE DATA     • HISTORICALLY IMPORTANT BUT MORE MODERN ALGOS ARE BETTER

# BAYESIAN ONLINE LEARNING

RECURSIVE APPLICATION OF BAYES RULE    $P(\theta|D_{1:n}) \propto P(D_n|\theta) P(\theta|D_{1:n-1})$    • RETURNS A FULL POSTERIOR
ONLINE FRIENDLY FOR HYPERPARAMS

- CAN BE **QUICKER** THAN **SGD**     • DIFFERENT RATE FOR EACH PARAM.     • **2ND** ORDER MODELS ARE TRICKY ONLINE

SIMPLE APPROXIMATION OF CURVATURE OF SPACE

EXAMPLES:   KALMAN FILTER (ONLINE LINEAR REGRESSION ... CONVERGE TO OFFLINE (OPTIMAL) VALUE IN SINGLE PASS OVER THE DATA)

PARTICLE FILTER

DENSITY FILTER

# ADADELTA

- RESTRICTS ACCUMULATION WINDOW
- EXPONENTIAL MOVING AVERAGE  $E[g^2] = \rho E[g^2]_{t-1} + (1-\rho)g^2_t \longrightarrow RMS[y_t] = \sqrt{E[g^2]_t + \epsilon}$
- $\Delta x_t = \dfrac{-\eta}{RMS[g_t]} \cdot g_t$  'UNIT NORMALIZATION'  $\Delta x_t = \dfrac{- RMS[\Delta x]_{t-1}}{RMS[g]_t} \cdot g_t$

# RMSPROP

- KEEP MA OF SQUARED GRADIENT FOR EACH WEIGHT
- $MS = 0.4 \, MS(w,t-1) + 0.1\left(\dfrac{\partial E}{\partial w[t]}\right)^2$
- $\Delta x_t = -\dfrac{g_t}{\sqrt{MS(w,t)}}$

# ADAM

- EXPONENTIAL DECAY RATES $\beta_1, \beta_2$ . 1ST MOMENT $m$, 2ND MOMENT $v$.
- $m_t = \beta_1 m_{t-1} + (1-\beta_1)g_t$  EMA MEAN
- $v_t = \beta_2 v_{t-1} + (1-\beta_2)g^2_t$  EMA VAR $\to$ DIAG FISHER MATRIX APPROX | ADAGRAD IF $\beta_2 \to 1$   $\beta_1 \to 0$
- BIAS CORRECT: $\hat{m}_t = m_t/(1-\beta_1^t)$, $\hat{v}_t = v_t/(1-\beta_2^t)$
- $\Delta x_t = -\dfrac{\alpha \hat{m}_t}{(\sqrt{\hat{v}_t} + \epsilon)} \approx$ SNR-ISH
- ADAMAX: SCALES GRADIENTS PROPORTIONALLY TO INFINITY NORM