

Newton's Method and logistic regression, more least squares

least-squares polynomial regression

replace each x_i with feature vector $\Phi(x_i)$

ex. $\Phi(x_i) = [x_{i1}^2 \ x_{i1}x_{i2} \ x_{i2}^2 \ x_{i1} \ x_{i2} \ 1]^T$
↖ ↗
features
dimension

otherwise just like linear / logistic regression

very easy to overfit data

Weighted least squares regression

$$\sum_{i=1}^n w_i (\hat{y}_i - y_i)^2$$

$$w_i \rightarrow S_2 = \begin{bmatrix} w_1 & w_2 & \dots \\ 0 & & \\ & & w_n \end{bmatrix}$$

$$\min_w (X^T w - y)^T S_2 (X^T w - y)$$

w/ calculus: $X^T S_2 X w = X^T S_2 y$

Newton's Method

- smooth for $J(w)$
- faster

at v , approx $J(w)$ by quadratic, jump to critical point
↳ repeat about convergence

Taylor series about v :

$$\nabla J(w) = \nabla J(v) + \underbrace{(\nabla^2 J(v))}_{\text{Hessian of } J} (w-v) + O(\|w-v\|^2)$$

find critical point $\nabla J(w) = 0$:

$$w = v - (\nabla^2 J(v))^{-1} \nabla J(v)$$

Newton's Method:

$w \leftarrow \text{init}$

repeat until convergence:

$$e \leftarrow \text{sol } (\nabla^2 J(w))e = -\nabla J(w)$$

$$w \leftarrow w + e$$

- Doesn't know difference of min/max/saddle points
- depends on starting point

Logistic regression

recall: $s'(x) = s(x)(1-s(x))$, $s_i = s(X_i^T w)$ $s = \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix}$

$$\nabla J = -\sum_{i=1}^n (y_i - s_i) X_i = -X^T (Y - s)$$

$$\nabla_w^2 J = \sum_{i=1}^n s_i(1-s_i) X_i X_i^T = X^T S X \quad S := \begin{bmatrix} s_1(1-s_1) & & 0 \\ & \ddots & \\ 0 & & s_n(1-s_n) \end{bmatrix} \succ 0$$

$$\Rightarrow X^T S X \geq 0$$

$\Rightarrow J(w)$ is convex

Newton's method:

$$w \leftarrow 0$$

repeat until convergence

$$e \leftarrow \text{sol to } (X^T \Omega X) e = X^T (y - s)$$

$$w \leftarrow w + e$$

iteratively weighted least squares



LDA vs. Logistic Regression

advantage of LDA:

- for well separated classes, stable; log reg. unstable
- > 2 classes; log reg better for 2 class (softmax harder)
- LDA slightly more accurate

advantages of log reg:

- more emphasis on decision boundary (close points in LDA not given diff. prob.ity)
- always separates linearly separable pts
- more robust to non-gaussian dists. (longer skew)
- naturally fits in (0, 1)

ROC Curves (for test sets)

