

## Lec 5

Thursday, September 12, 2019 10:52

Recap: Logistic regression (Gnan)

$$\text{Posit logit}(P(Y=1|X)) = \beta^T X$$

$$\text{alt: } P(Y=1|X) = \sigma(\beta^T X)$$

Fitting logistic regression

Logistic regn model specifies a generative model for the data.

- First draw  $X_i$

- Then compute  $\sigma(\beta^T X_i)$

-  $Y_i \sim \text{Bernoulli}(\sigma(\beta^T X_i))$

(i.e. generate  $U_i \sim \text{Unif}(0,1)$ )

$$\text{set } Y_i = \mathbb{I}[U_i \leq \sigma(\beta^T X_i)]$$

Assuming that the data is independent,

for every  $\beta$ , there is a particular likelihood for observing our data

$$\text{Lik}(\beta) = P(X_1, Y_1, \dots, X_n, Y_n; \beta)$$

$$= \prod_{i=1}^n P(X_i, Y_i; \beta)$$

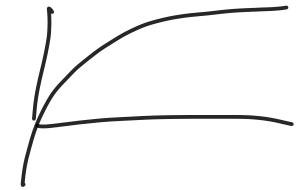
$$= \prod_{i=1}^n P(Y_i | X_i; \beta) P(X_i)$$

$$= \prod_{i=1}^n P(X_i) \cdot \left( \begin{cases} \sigma(\beta^T X_i) & \text{if } Y_i=1 \\ 1 - \sigma(\beta^T X_i) & \text{if } Y_i=0 \end{cases} \right)$$

Max likelihood principle:

choose parameters that max  
 $L(\beta)$  if observing the data

The veriness of ...  
that we observe

log =  - always increasing

$$\begin{aligned} \Rightarrow \arg \max_{\beta} \text{Lik}(\beta) \\ &= \arg \max_{\beta} \log(\text{Lik}(\beta)) \\ &= \arg \min_{\beta} (-\log(\text{Lik}(\beta))) \leftarrow \text{neg log lik} \end{aligned}$$

$$\begin{aligned} -\log \text{Lik}(\beta) &= \sum_{i=1}^n \left( -\log P(x_i) - \log P(y_i | x_i; \beta) \right) \\ &= -\sum_i \log P(x_i) \\ &\quad + \sum_i \begin{cases} -\log \sigma(\beta^T x_i) & y_i = 1 \\ -\log(1 - \sigma(\beta^T x_i)) & y_i = 0 \end{cases} \end{aligned}$$

$$\begin{aligned} &= -\sum_i \log(P(x_i)) \\ &\quad + \underbrace{\sum_i \left( y_i (-\log \sigma(\beta^T x_i)) + (1 - y_i) (-\log(1 - \sigma(\beta^T x_i))) \right)}_{\mathcal{L}(\beta) \leftarrow \text{neg log lik fun}} \end{aligned}$$

$$\arg \max_{\beta} \text{Lik}(\beta) = \arg \min_{\beta} \mathcal{L}(\beta)$$

$$-\log(\sigma(\beta^T x_i)) = -\log\left(\frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}\right) = -\beta^T x_i + \log(1 + e^{\beta^T x_i})$$

$$-\log(1 - \sigma(\beta^T x_i)) = -\log\left(\frac{1}{1 + e^{\beta^T x_i}}\right) = +\log(1 + e^{\beta^T x_i})$$

$$\begin{aligned} \mathcal{L}(\beta) &= \sum_i \left( \cancel{y_i} \log(1 + e^{\beta^T x_i}) - y_i \beta^T x_i \right. \\ &\quad \left. + (1 - \cancel{y_i}) \log(1 + e^{\beta^T x_i}) \right) \end{aligned}$$

$$= \sum_i (\log(1 + e^{\beta^T x_i}) - y_i \beta^T x_i)$$

Fitting logistic regn: solving  $\min_{\beta} \mathcal{L}(\beta)$

call  $\hat{\beta}$  the optimal  $\beta$

Can solve w/  $\nabla \mathcal{L}(\beta) = 0$

but computationally more difficult than OLS

Revisit after fall break

More than two categories: Multinomial Logistic Regression

$$Y \in \{1, \dots, m\} \quad m > 2$$

$$\text{Posit } P(Y=j | X=x) \propto e^{\beta_j^T x}$$

$$\beta_1, \dots, \beta_m \in \mathbb{R}^m$$

$$1 = \sum_j P(Y=j | X=x) \Rightarrow P(Y=j | X=x) = \frac{e^{\beta_j^T x}}{\sum_{j'=1}^m e^{\beta_{j'}^T x}}$$

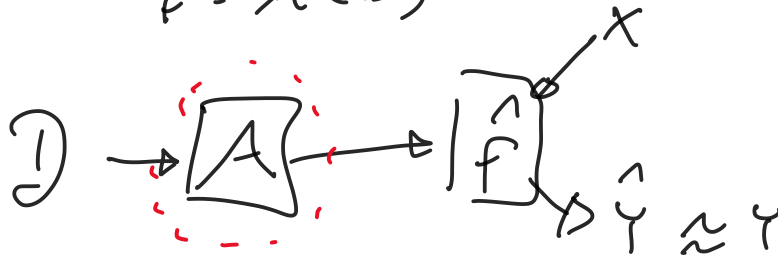
$$\text{Fix } \beta_m = \vec{0}.$$

Fit by maximum likelihood

The bias-variance tradeoff

A supervised learning algo can be understood as a fn  $A$  from  $\mathcal{D} = \{x_1, y_1, \dots, x_n, y_n\}$  to a prediction rule  $\hat{f}$

$$\hat{f} = A(D)$$



What makes a good algo?

A good algo should do well over all over all dataset  $D$

$$\begin{aligned} R(A) &= \mathbb{E}_D R(A(D)) \\ &= \mathbb{E}_D [R(\hat{f})] \\ &= \mathbb{E}_D [ \mathbb{E}_{x,y} [ \ell(y, \hat{f}(x)) ] ] \end{aligned}$$

random b/c  
new random  
test pt

random b/c  
new random  
training datasets

for regression

$$= \mathbb{E}_D [ \underbrace{\mathbb{E}_{x,y} [ (y - \hat{f}(x))^2 ]}_{\otimes} ]$$

$$\otimes = \mathbb{E}_{x,y} [ (y - f(x))^2 ]$$

$$= \mathbb{E}_x [ \mathbb{E}_y [ (y - f(x))^2 | x ] ]$$

$$= \mathbb{E}_x [ \mathbb{E}_y [ y^2 - 2yf(x) + f(x)^2 | x ] ]$$

$$= \mathbb{E}_x [ \mathbb{E}_y [ y^2 | x ] - 2f(x) \underbrace{\mathbb{E}[y|x]}_{f^*(x)} + f(x)^2 ]$$

$$= \mathbb{E}_x \left[ \underbrace{\mathbb{E}_y (Y^2 | X) - f^{*2}(x)}_{\text{Var}(Y|X)} + \underbrace{f^2(x) - 2f(x)f^*(x) + f^*(x)^2}_{(f^*(x) - f(x))^2} \right]$$

$$= \underbrace{\mathbb{E}_x \left[ (f^*(x) - f(x))^2 \right]}_{\substack{\text{reducible error} \\ \text{can be made} \\ \text{small w/ } f \approx f^*}} + \underbrace{\mathbb{E}_x \text{Var}(Y|X)}_{\substack{\text{irreducible error} \\ \text{no matter what } f \text{ is}}}$$

$$\text{Hence } R(f) = \mathbb{E}_D \left[ \begin{array}{c} \uparrow \\ \end{array} \right]$$

$$= \mathbb{E}_x \text{Var}(Y|X)$$

$$+ \mathbb{E}_x \left[ \underbrace{\mathbb{E}_D \left[ (f^*(x) - \hat{f}(x))^2 \right]}_{\text{Err}(x)} \right]$$

$$\text{Err}(x) = \mathbb{E}_D \left[ (\hat{f}(x) - f^*(x))^2 \right]$$

$$= \mathbb{E}_D \left[ \left( (\hat{f}(x) - \mathbb{E}_D \hat{f}(x)) + (\mathbb{E}_D \hat{f}(x) - f^*(x)) \right)^2 \right]$$

$$= \mathbb{E}_D \left[ (\hat{f}(x) - \mathbb{E}_D \hat{f}(x))^2 \right] + (\mathbb{E}_D \hat{f}(x) - f^*(x))^2$$

$$+ 2(\mathbb{E}_D [\hat{f}(x) - \mathbb{E}_D \hat{f}(x)]) (\mathbb{E}_D \hat{f}(x) - f^*(x))$$

$$= \text{Var}_D \hat{f}(x) + (\mathbb{E}_D \hat{f}(x) - f^*(x))^2$$

$$\underbrace{\hspace{10em}}_{\text{bias}^2}$$