

what if the response Y is categorical?

Still interested in discriminative models, i.e. models for $P\{Y|X\}$

Y	Probability
1	$P\{Y=1\} = \pi$
0	$P\{Y=0\} = 1-\pi$

$Y \sim \text{Bernoulli}(\pi)$, $E\{Y\} = \pi$, $\text{Var}\{Y\} = \pi(1-\pi)$

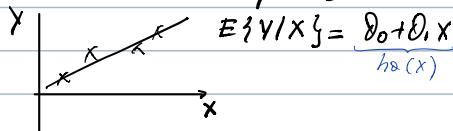
Unlike Normal distribution,
mean and variance depend
on the same parameter

Question: how does $P\{Y=1\}$ depend on predictor X ?

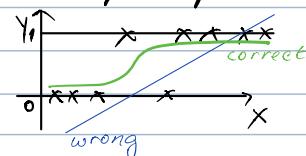
Specify hypothesis $h_\theta(x) : X \rightarrow Y = \{0, 1\}$

mean response function
or activation function

Continuous response:



Binary response:



$$\pi = E\{Y|X\} = g(\theta_0 + \theta_1 x)$$

$$g(E\{Y|X\}) = \theta_0 + \theta_1 x$$

link function

Good choice for $g(\cdot)$: a CDF of a probability distribution

$g(t) = \text{CDF of } N(0,1) \Rightarrow \text{probit regression}$ (no analytical expression)

$g(t) = \frac{e^t}{1+e^t} = \text{CDF of logistic distribution} \Rightarrow \text{logistic regression}$

$$E\{Y|X\} = g(\theta_0 + \theta_1 x)$$

$\underbrace{\text{linear function}}_{\text{non-linear function}}$

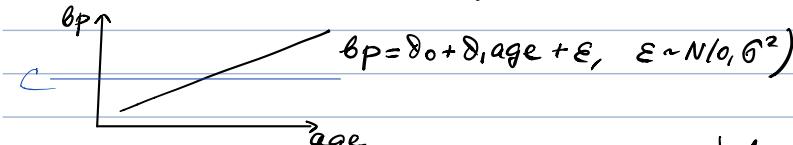
not a linear model

a generalized linear model

no analytical solution for param. estimate
(numeric solution only)

Motivation for probit regression: Latent variable

Example: a linear regression for blood pressure as function of age:



Only observe hypertension: $Y = \begin{cases} 1, & BP > c \\ 0, & BP \leq c \end{cases}$

$$\begin{aligned} P\{Y=1\} &= P\{BP > c\} = P\{\alpha + \beta \text{age} + \epsilon > c\} = P\{\epsilon < (\alpha - c) + \beta \text{age}\} \\ &= P\left\{\frac{\epsilon}{\sigma} < \frac{\alpha - c}{\sigma} + \frac{\beta}{\sigma} \text{age}\right\} = \Phi\left(\frac{\theta_0 + \theta_1 \text{age}}{\sigma}\right) \end{aligned}$$

$\underbrace{\sim N(0,1)}_{\text{~} \sim N(0,1)}$ $\underbrace{\theta_0}_{\text{~} \theta_0}$ $\underbrace{\theta_1}_{\text{~} \theta_1}$ Φ CDF of $N(0,1)$

Logistic regression:

$$\pi = E\{Y|X\} = \frac{e^{\theta_0 + \theta_1 x}}{1 + e^{\theta_0 + \theta_1 x}} = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}} \quad \begin{array}{l} \text{- monotonic increasing if } \theta_1 > 0 \\ \text{- monotonic decreasing if } \theta_1 < 0 \end{array}$$

$$\log \frac{E\{Y|X\}}{1 - E\{Y|X\}} = \log \frac{\pi}{1-\pi} = \theta_0 + \theta_1 x$$

$$\log \frac{E\{Y|X=x+1\}}{1 - E\{Y|X=x+1\}} / \frac{E\{Y|X=x\}}{1 - E\{Y|X=x\}} = \log \frac{\text{odds}(X=x+1)}{\text{odds}(X=x)} = \text{log(ratio)}$$

Interpretation of θ_1 :

$$\log \frac{E\{Y|X=x+1\}}{1 - E\{Y|X=x+1\}} - \log \frac{E\{Y|X=x\}}{1 - E\{Y|X=x\}} = \theta_1$$

$\underbrace{\text{odds}(X=x+1)}_{\text{odds}(X=x)}$

(2)

Parameter estimation: no error term; least squares
not appropriate \Rightarrow maximum likelihood

$y_i \stackrel{\text{ind}}{\sim} \text{Bernoulli}(\pi_i)$, where $\pi_i = \frac{e^{\delta_0 + \delta_1 x_i}}{1 + e^{\delta_0 + \delta_1 x_i}}$

Likelihood $L = \prod_{i=1}^N \underbrace{\pi_i^{y_i} (1-\pi_i)^{1-y_i}}_{\substack{\text{Bernoulli distribution} \\ \text{indep.}}}$

$$\log L = \sum_{i=1}^N y_i \log(\pi_i) + (1-y_i) \log(1-\pi_i) = \sum_{i=1}^N \left[y_i \log \frac{1}{1+e^{-(\delta_0+\delta_1 x_i)}} + (1-y_i) \log \left(1 - \frac{1}{1+e^{(\delta_0+\delta_1 x_i)}}\right) \right]$$

plug in π_i

convex

To maximize $\log L$, set $\nabla_{\delta} \log L = 0$. It helps to use chain rule:

$$\text{Define } \pi(z) = \frac{1}{1+e^{-z}}, \pi'(z) = \frac{d}{dz} \frac{1}{1+e^{-z}} = \frac{1}{(1+e^{-z})^2} e^{-z} = \frac{1}{1+e^{-z}} \cdot \left(1 - \frac{1}{1+e^{-z}}\right) = \pi(z)(1-\pi(z))$$

$$\text{In this notation, } \log L = \sum_{i=1}^N y_i \log(\pi(\delta_0 + \delta_1 x_i)) + (1-y_i) \log(1-\pi(\delta_0 + \delta_1 x_i))$$

$$\begin{aligned} \frac{d \log L}{d \delta_1} &= \frac{d \log L}{d \pi} \cdot \frac{d \pi}{d z} \cdot \frac{d z}{d \delta_1} = \sum_{i=1}^N \left[y_i \frac{1}{\pi(\delta_0 + \delta_1 x_i)} - (1-y_i) \frac{1}{1-\pi(\delta_0 + \delta_1 x_i)} \right] \cdot \cancel{\pi(\delta_0 + \delta_1 x_i)(1-\pi(\delta_0 + \delta_1 x_i))} \cdot x_i \\ &= \sum_{i=1}^N (y_i - \pi(\delta_0 + \delta_1 x_i)) x_i \end{aligned}$$

In multivariate case, Batch gradient descent:

Initialize δ_{px1}

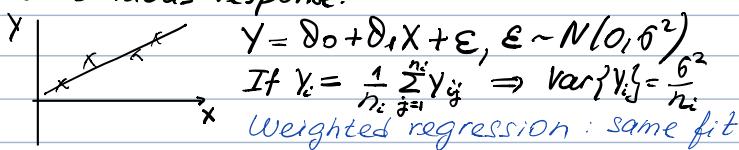
Repeat until convergence {

$$\text{For } j = 1, \dots, p \quad \delta_j \leftarrow \delta_j + \alpha \sum_{i=1}^N (y_i - \pi(x_i^T \delta)) x_{ij}$$

}

Practical comment: grouped data

Continuous response:



Log-likelihood:

$$\log \prod_{i=1}^N \binom{n_i}{y_i} \pi_i^{y_i} (1-\pi_i)^{n_i-y_i}$$

$$= \sum_{i=1}^N \log \binom{n_i}{y_i} + \sum_{i=1}^N [y_i \log \pi_i + (n_i - y_i) \log(1 - \pi_i)] = \text{const} + \sum_{i=1}^N \sum_{j=1}^{n_i} [y_{ij} \log(\pi_i) + (1 - y_{ij}) \log(1 - \pi_i)]$$

does not depend on π_i , can ignore

$\sum_{j=1}^{n_i} y_{ij}$, where $y_{ij} = \begin{cases} 1 & \text{if } x_{ij} \\ 0 & \text{otherwise} \end{cases}$

proportion of $Y=1$ out of n_i obs with same X

Binary response: Data: (x_i, y_i, n_i)
 $y_i \sim \text{Binomial}(n_i, \pi_i)$
 $E\{y_i\} = n_i \pi_i$
 $\text{Var}\{y_i\} = n_i \pi_i (1 - \pi_i)$

\Rightarrow individual and grouped data give same δ if we know n_i

Logistic regression for classification



- this is the reason for the name "discriminative model"

Summary of classification:

		Decision	0	1	Total
Truth	0	TN	FP	N	
	1	FN	TP	P	

/ counts of observations
confusion matrix

$$\text{Accuracy} = \frac{TP+TN}{P+N}$$

Prevalence: $\frac{P}{P+N}$

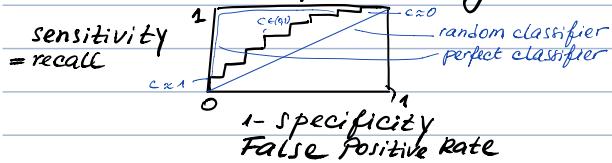
False positive rate = $\frac{FP}{N}$

Sensitivity = recall = true positive rate = $\frac{TP}{P}$

Specificity = selectivity = true negative rate = $\frac{TN}{N}$

Precision = positive predictive value = $\frac{TP}{TP+FP}$
False positive rate = $\frac{FP}{TP+FP}$

Receiver operating characteristic (ROC) curve



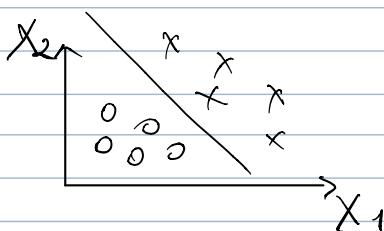
- Area under the ROC curve summarizes the predictive ability over all cutoffs

- More optimistic (i.e., higher) on the training set than on the validation set

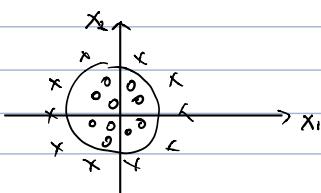
Decision boundary

The classification cutoff separates points by a hyperplane:

$$\hat{Y} = 1 \text{ if } \hat{\pi} = \frac{1}{1 + e^{-(\delta_0 + \delta_1 X)}} \geq \frac{1}{2} \iff 1 + e^{-(-\delta_0 - \delta_1 X)} \leq 2 \\ e^{-(-\delta_0 - \delta_1 X)} \leq 1 \\ -(\delta_0 + \delta_1 X) \leq 0 \\ \text{decision boundary is a hyperplane} \rightarrow \delta_0 + \delta_1 X \geq 0$$



$$Y = 1 \text{ if } \hat{\pi} = \frac{1}{1 + e^{-(3 - X_1 + X_2)}} \geq \frac{1}{2} \\ \Rightarrow 3 - X_1 + X_2 \geq 0 \Rightarrow X_2 \leq 3 - X_1$$



Higher-order terms can express complex decision boundaries

$$\frac{1}{1 + e^{(-1 + 0 \cdot X_1 + 0 \cdot X_2 + 1 \cdot X_1^2 + 1 \cdot X_2^2)}} \geq 0.5 \iff X_1^2 + X_2^2 \geq 1$$

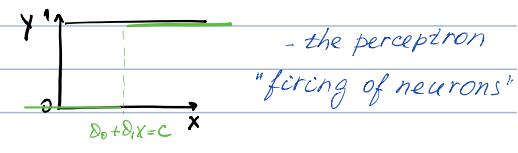
Comments:

- Linear separability. MLE does not have a unique solution
 $\text{Var}(\delta) \rightarrow \infty$

- Alternative ways to derive Linear Boundary:

$$f(X) = \begin{cases} 1, & \text{if } X_i^\top \delta \geq c \\ 0, & \text{if } X_i^\top \delta < c \end{cases} \quad \left\{ \begin{array}{l} g(x): \\ \text{step function} \end{array} \right.$$

No probabilistic interpretation; no MLE



- Variable selection : no residuals \rightarrow no R^2 , no MSE

$$\begin{aligned} AIC &= -2\log L + 2p \\ BIC &= -2\log L + \log N \cdot p \end{aligned} \quad \left. \begin{array}{l} \text{find subsets of predictors} \\ \text{that minimize AIC or BIC} \end{array} \right.$$

train/dev/test sets or cross-validation - as in linear regression

- Regularization

$$\text{maximize } \log L - 2\|\beta\|^2 \rightarrow \text{other penalties/algos as in regression}$$

Multi-class classification

		y				
		1	2	...	J	# possible classes
X_1	X_2	n_{11}	n_{12}	\dots	n_{1J}	grouped data
		n_{21}	n_{22}	\dots	n_{2J}	multinomial
X_M		n_{M1}	n_{M2}	\dots	n_{MJ}	

Multinomial Logistic regression:

$$\log \frac{\pi_j(x)}{\pi_1(x)} = \alpha_j + \beta_j x, \quad j = 2, \dots, J \quad \left. \begin{array}{l} \text{useful if interested} \\ \text{in param. interpretation} \end{array} \right.$$

arbitrary category as baseline,
to satisfy $\sum_j \pi_j = 1$

separate params
for each category

Solving for π_j :

$$\pi_1(x) = \frac{1}{1 + \sum_{j=1}^J e^{\alpha_j + \beta_j x}}, \quad \pi_j = \frac{e^{\alpha_j + \beta_j x}}{1 + \sum_{j=1}^J e^{\alpha_j + \beta_j x}}$$

$$\text{MLE: } \log \text{Likelihood}_i = \log \left[\prod_{j=1}^J \pi_j(x)^{y_{ij}} \right] = \sum_{j=2}^J y_{ij} \log \pi_j(x) + \left(1 - \sum_{j=2}^J y_{ij} \right) \log \pi_1(x)$$

for one data point to simplify

Alternative parametrization: softmax regression

$$\pi_j(x) = \frac{e^{\alpha_j + \beta_j x}}{\sum_{j=1}^J e^{\alpha_j + \beta_j x}} \quad \left. \begin{array}{l} \text{more convenient} \\ \text{if goal is prediction} \end{array} \right.$$

- only $J-1$ free parameter sets

- set $\alpha_1 = \beta_1 = 0$ to have the previous case