

Lecture 17 4/7/21

Response Space

$$Y \subseteq \mathbb{R}$$

$$Y = \{c_1, c_2, \dots, c_K\}$$

$$\text{if } K=2, Y = \{c_1, c_2\}$$

$$Y \subseteq \mathbb{R}_{\geq 0}$$

$$Y \subseteq \{0, 1, 2, \dots\}$$

$$Y \subseteq (0, 1)$$

$$Y = \{c_1, c_2, \dots, c_K\}$$

$$Y \subseteq \{c_1, c_2, \dots, c_K\} \text{ ordinal}$$

$$K=2, Y = \{c_1, c_2\}$$

Types of Modeling

regression

classification

binary classification

Survival

Count

proportion

probability estimation

probability estimation

probability estimation

g return

$$\hat{y} \in Y$$

$$\hat{y} \in Y$$

$$\hat{y} \in Y$$

$$\hat{y} \in Y$$

$$\hat{y} \in Y$$

$$\hat{y} \in Y$$

$$\hat{p} := P(Y=c_1 | \vec{x})$$

$$P(Y=c_2 | \vec{x})$$

$$P(Y=c_K | \vec{x})$$

$$\hat{p} := P(Y=1 | \vec{x})$$

Example Alg.

OLS

KNN

SVM

Weibull regression

Poisson Regression

Beta regression

Multi-layer regression

proportional odds model

Logistic regression

If $y = \{0, 1\}$ for all i ,

$$y = t(\frac{z}{2})$$

$$= f(\vec{x}) + \delta \text{ where } \delta \in \{0, -1, +1\}$$

$$= h(\vec{x}) + \varepsilon \text{ where } \varepsilon \in \{0, -1, +1\}$$

$$= g(\vec{x}) + e \text{ where } e \in \{0, -1, +1\}$$

How do we build a probability estimation model?

Naturally,

$$g_0 = \bar{y}$$

		$g(\vec{x})$	
	0	1	
0	0	-1	False positive
1	+1	0	False negative

$$\Rightarrow Y \sim \text{Bern}(t(\frac{z}{2}))$$

We now view Y as a realization from a random variable (bernoulli). We will assume there exists a function $f_{pr}(\vec{x}): \mathbb{R}^{n+1} \rightarrow (0,1)$ and this function is the best guess of the probability $P(Y=1 | \vec{x}_{vec})$ you can create with \vec{x}_{vec} .

$$Y \sim \text{Bern}(f_{pr}(\vec{x}) + \underbrace{t(\vec{x}) - f_{pr}(\vec{x})}_{\delta_{pr}})$$

$\Rightarrow Y \sim \text{Bern}(f_{pr}(\vec{x}))$. f_{pr} is the model we want to find.

Let's assume that all the data (all the n observations) in D are independently realized.

$$\begin{aligned} P(D) &= P(Y_1=y_1, Y_2=y_2, \dots, Y_n=y_n | \vec{x}_1, \dots, \vec{x}_n) \\ &= \prod_{i=1}^n P(Y_i=y_i | \vec{x}_i) \\ &= \prod_{i=1}^n f_{pr}(\vec{x}_i)^{y_i} (1 - f_{pr}(\vec{x}_i))^{1-y_i} \end{aligned}$$

$$\begin{aligned} V &\sim \text{Bern}(\theta) \\ \text{241} &\quad \theta^V (1-\theta)^{1-V} \\ \text{Review} &\end{aligned}$$

Now we want to "fit" f_{pr} using our data (learning from data paradigm). How? Is it even possible? NO. We cannot fit arbitrary functions in any dimension. We need a set of candidate functions that we can fit. Call that \mathcal{H}_{pr} . Each element in this set maps $\mathbb{R}^{n+1} \rightarrow (0,1)$. How about:

$$\mathcal{H}_{pr} = \{ \vec{w} \cdot \vec{x} : \vec{w} \in \mathbb{R}^{n+1} \} ?$$

We can't use this since it returns values outside $(0,1)$, the space of legal probabilities. But... we really like $\vec{w}_{vec} \cdot \vec{x}_{vec}$ because (1) easy to interpret and we have lots of intuition about it from all of our previous modeling we've done and (2) monotonic in each of the x_j 's. How do we have our cake and eat it too? We need a function that takes $\vec{w}_{vec} \cdot \vec{x}_{vec}$ and maps it into the space $(0,1)$, i.e. $\phi: \mathbb{R} \rightarrow (0,1)$ which is called a "link function" I think because it links the two spaces (the reals and the probs). We restrict the link function to be strictly increasing. Thus,

$$\mathcal{H}_{pr} = \{ \phi(\vec{w} \cdot \vec{x}) : \vec{w} \in \mathbb{R}^{n+1} \}$$

These types of models are called "generalized linear models" (glm) because they retain $wvec \cdot xvec$ (the linear model) but then manipulate it in some way. Which link function should we use? There are three common ones. In order of use:

① Logistic / logit: $\phi(u) := \frac{e^u}{1+e^u} = \frac{1}{1+e^{-u}}$. Note: $1 - \phi(u) = \frac{1}{1+e^u}$

② Probit: $\phi(u) := F_{\Xi}(u)$ i.e. the CDF of the std. normal.

③ Complementary Log-Log (cloglog)

$$\begin{aligned} \phi(u) &= 1 - e^{-e^u} \Rightarrow 1 - \phi(u) = e^{-e^u} \Rightarrow \ln(1 - \phi(u)) = -e^u \\ \Rightarrow -\ln(1 - \phi(u)) &= e^u \Rightarrow u = \underbrace{\ln(-\ln(1 - \phi(u)))}_{\text{complement}} \end{aligned}$$

Let's employ the logistic link function:

$$\mathcal{H} = \left\{ \frac{1}{1+e^{-\vec{w} \cdot \vec{x}}} : \vec{w} \in \mathbb{R}^{r+1} \right\}$$

What is \mathcal{H} ? How to get $g \in \mathcal{H}$?

Why not find the $wvec$ that provides us the highest probability?

$$A: \vec{b} := \underset{\vec{w} \in \mathbb{R}^{r+1}}{\operatorname{argmax}} \left\{ \underbrace{\prod_{i=1}^n \left(\frac{1}{1+e^{-\vec{w} \cdot \vec{x}_i}} \right)^{y_i} \left(\frac{1}{1+e^{\vec{w} \cdot \vec{x}_i}} \right)^{1-y_i}}_{P(D)} \right\}$$

In OLS, we took the derivative and set it equal to zero to solve for $bvec$ and we found an analytical solution. However, there is no analytical solution here. You need to use a computer.

$$\vec{\nabla} P(D) \stackrel{\text{set}}{=} \vec{0}_{r+1} \text{ and approximate}$$

Usually this is done with "gradient descent". Computing $bvec$ is called "running a logistic regression". Once this is done ... we can predict using

$$\hat{p} = g_{rr}(\vec{x}) = \phi(\vec{b} \cdot \vec{x}) = \frac{1}{1+e^{-\vec{b} \cdot \vec{x}}} \text{ hopefully close to } f_{pr}(\vec{x})$$

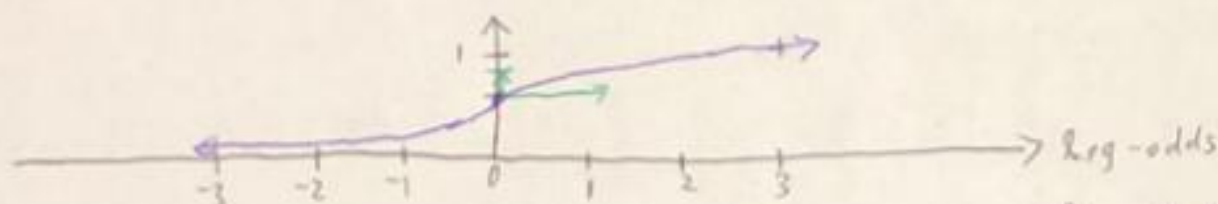
$$\hat{p}(Y=1 | \vec{x})$$

What is the interpretation of the slope coefficients (the cutates in the h-vet)?

$$\hat{p} = \frac{1}{1 + e^{-\vec{b} \cdot \vec{x}}} \Rightarrow \frac{1}{\hat{p}} = 1 + e^{-\vec{b} \cdot \vec{x}} \Rightarrow \frac{1}{\hat{p}} - 1 = e^{-\vec{b} \cdot \vec{x}} \Rightarrow \frac{1 - \hat{p}}{\hat{p}} = e^{-\vec{b} \cdot \vec{x}} \Rightarrow \ln\left(\frac{1 - \hat{p}}{\hat{p}}\right) = -\vec{b} \cdot \vec{x}$$

$$\Rightarrow \underbrace{\ln\left(\frac{\hat{p}}{1 - \hat{p}}\right)}_{\text{log-odds}} \quad \text{odds} = \frac{\hat{p}}{1 - \hat{p}}$$

$\Rightarrow b_j$ is the ~~change~~ change in the log-odds of $Y=1$ if x_j increases by 1.



$$0 = \ln\left(\frac{\hat{p}}{1 - \hat{p}}\right) \Rightarrow \hat{p}(0) = \frac{1}{1 + e^0} = \frac{1}{1 + 1} = 0.5$$

The logistic link function is highly nonlinear.

$0 \rightarrow 1$ logodds \Leftrightarrow 50% \rightarrow 73% in prob.

$3 \rightarrow 4$ logodds \Leftrightarrow 95% \rightarrow 98% in prob.

Log-odds	prob
0	0.5
-1	0.27 \approx 1/4
+1	0.73 \approx 3/4
-2	0.12 \approx 1/8
+2	0.88 \approx 7/8
\vdots	
$-\infty$	0
$+\infty$	1

Probability estimation models predict probabilities but we observe labels (i.e. 0 or 1). The true probabilities f_{pr} are unobserved! We need a metric called a "scoring rule" S that can compare a \hat{p} value to a y value.

A "proper scoring rule" $S(\hat{p}, y)$ is one where:

$$\forall_i \quad f_{pr}(x_i) = \arg\max \{ S(\hat{p}_i, y_i) \}$$

We will study two proper scoring rules:

① Brier score (1950). Let $S_i := -(y_i - \hat{p}_i)^2 \leq 0$

$$\bar{S} := \frac{1}{n} \sum_{i=1}^n S_i \leq 0$$

② Log scoring rule. Let $S_i := y_i \ln(\hat{p}_i) + (1 - y_i) \ln(1 - \hat{p}_i) \leq 0$

$$\bar{S} = \frac{1}{n} \sum S_i \leq 0$$

These scores are used as an " R^2 " of the model (but they're not between 0 and 1) in a conceptual sense. The closer to zero, the better the probability estimation model.