| Response Space | Type of Modeling | g return | Example alg |
|---|---|---|---|
| $y \subseteq \mathbb{R}$ | regression | $\hat{y} \in y$ | OLS |
| $y \subseteq \{c_1, c_2, ..., c_p\}$ | classification | $\hat{c} \in y$ | KNN |
| $K=2,\ y = \{\begin{smallmatrix}c_1\\c_2\end{smallmatrix}, \begin{smallmatrix}\\0,1\end{smallmatrix}\}$ | binary classification | $\hat{y} \in y$ | SVM |
| $y \subseteq \mathbb{R}_{\geq 0}$ | survival | $\hat{y} \in y$ | Weibull regression |
| $y \subseteq \{0,1,2,...\}$ | count | $\hat{y} \in y$ | Poisson regression |
| $y \in (0,1)$ | proportion | $\hat{y} \in y$ | Beta regression |
| $y = \{c_1, c_2, ..., c_k\}$ | probability estimation | $\hat{P} := \{\begin{smallmatrix}P(Y=c_1|\vec{x})\\P(Y=c_2|\vec{x})\\P(Y=c_k|\vec{x})\end{smallmatrix}\}$ | Multinomial regression |
| $K=2,\ y = \{0,1\}$ | probability estimation | $\hat{p} := P(Y=1|\vec{x})$ | Logistic regression |
| $y \in \{c_1,...,c_p\}$ Ordinal | probability estimation | | Proportional odds model |

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

$y = t(\vec{z})$
$= f(\vec{x}) + \delta$ where $\delta \in \{0, -1, +1\}$
$= h(\vec{x}) + \xi$ where $\xi \in \{0, -1, +1\}$
$= g(\vec{x}) + e$ where $e \in \{0, -1, +1\}$

false positive, false negative

How do we build a probability estimation model? $g_0 = \bar{y}$   Naively,

$\iff Y \sim Bern(t|\vec{z})$

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

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

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

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

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

$Y \sim Bern(\theta)$, $\rho \propto \theta^y (1-\theta)^{1-y}$

Now we want to "fit" f_pr using our data (learning from data paradigm). How? Is this 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 $\mathbb{R}^{p+1} \to (0,1)$. How about:

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

We can't use this since it returns values outside (0,1), the space of legal probabilities. But... we really like wvec dot xvec 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 we do we have our cake and eat it too?

We need a function that takes wvec dot xvec and maps it into the space (0,1) i.e. $\phi: \mathbb{R} \to (0,1)$ which is called a "link function". I think because it links the two spaces (the reals and the prob's). We restrict the link function to be strictly increasing. Thus,

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

These types of models are called "generalized linear models" (glm) because they retain wvec dot 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: $\phi_{(G)} := \frac{e^u}{1+e^u} = \frac{1}{1+e^{-u}}$.  Note: $1 - \phi = \frac{1}{1+e^u}$

(2) Probit: $\phi(u) := F_z(u)$ i.e. the CDF of the std. normal.

(3) Complementary Log-Log (cloglog)
$$\phi(u) = 1 - e^{-e^u} \Rightarrow 1 - \phi(u) = e^{-e^u} \Rightarrow \ln(1 - \phi(u)) = -e^u$$
$$\Rightarrow -\ln(1-\phi(u)) = e^u \Rightarrow u = \ln(-\ln(1 - \phi(u)))$$  cloglog

Let's employ the logistic link function:
$$\mathcal{H} = \left\{ \frac{1}{1+e^{-\vec{w} \cdot \vec{x}}} : \vec{w} \in \mathbb{R}^{p+1} \right\}$$

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

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

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

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.

$$\vec{\nabla} P(D) \overset{set}{=} \vec{0}_{p+1} \quad \text{and approximate}$$

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

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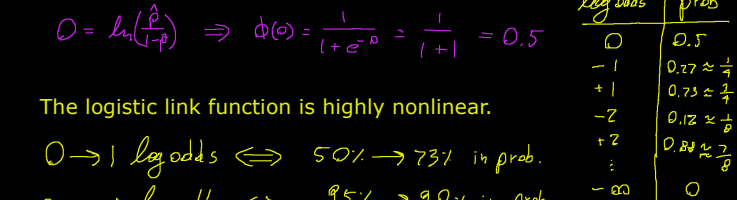$$\overset{\shortparallel}{\hat{p}(Y=1|\vec{x})}$$

What is the interpretation of the slope coefficients (the entries in the b-vec)?

$$\hat{p} = \frac{1}{1+e^{-\vec{b}\cdot\vec{x}}} \Rightarrow \frac{1}{\hat{p}} = 1 + e^{-\vec{b}\cdot\vec{x}} \Rightarrow \frac{1}{\hat{p}} - 1 = e^{-\vec{b}\cdot\vec{x}}$$

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

$Odds = \frac{\hat{p}}{1-\hat{p}}$   odds,   log-odds

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

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

| log-odds | prob |
|---|---|
| 0 | 0.5 |
| -1 | 0.27 ≈ 1/4 |
| +1 | 0.73 ≈ 3/4 |
| -2 | 0.12 ≈ 1/8 |
| +2 | |
| -∞ | 0 |
| +∞ | 1 |

The logistic link function is highly nonlinear.

$0 \to 1$ log-odds $\iff 50\% \to 73\%$ in prob.
$3 \to 4$ log-odds $\iff 95\% \to 98\%$ in prob.

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 p-hat value to a y value.

A "proper scoring rule" S(p-hat, y) is one where:
$$\forall i \quad f_{pr}(\vec{x}_i) = argmax \{ S(\hat{p}_i, y_i) \}$$

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.