

Classification

Logistic Regression

(1) HW 1 due 4/21 (last update 4/13)

Thus far, regression:

predict a continuous value given some inputs

Given $x \in \mathbb{R}^d$

$y \in \mathbb{R}$

continuous

predict $y = f(x)$

1, 0.1, π - - -

Reading Your Brain, Simple Example

Pairwise classification accuracy: 85%

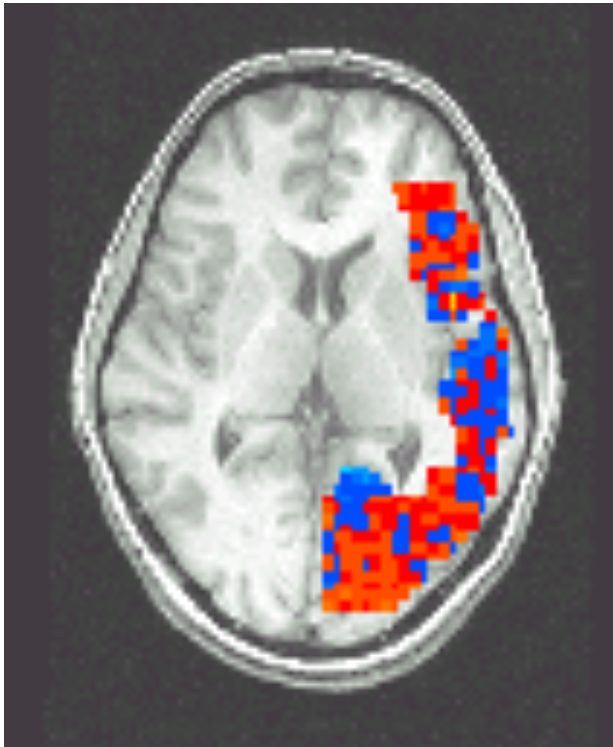
[Mitchell et al.]

$f(x_1) = \text{"Person"}$
 $f(x_2) = \text{"Animal"}$
encoding:
"Person" = 1
"Animal" = 0
 $\begin{cases} f(x_1) = 1 \\ f(x_2) = 0 \end{cases}$

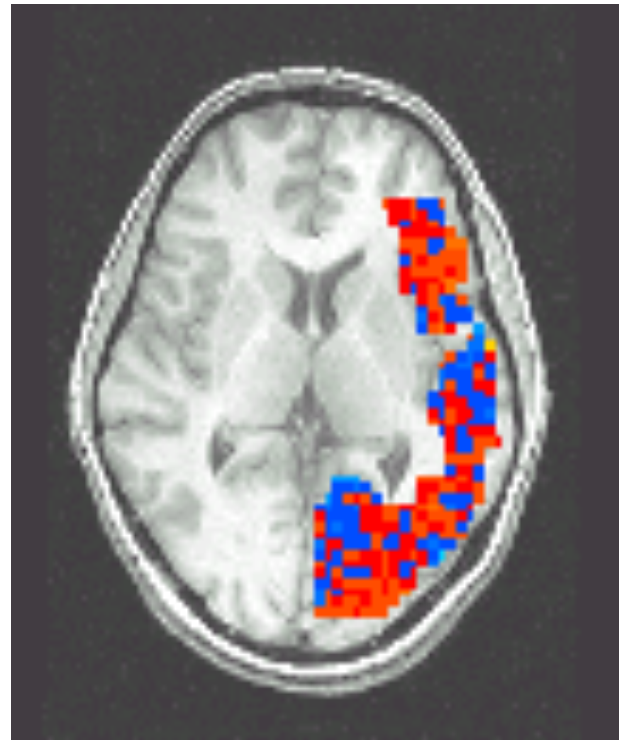
Person



Animal



x_1



x_2

Classification

- Learn $f: X \rightarrow Y$

- X - features

- Y - target classes $\{0, 1\}$, $\{1, \dots, K\}$ discrete
binary multi-class

- Loss Function

$$l(f(x), y) = \mathbb{1}\{f(x) \neq y\} \quad \begin{cases} 1 & \text{wrong} \\ 0 & \text{correct} \end{cases}$$

- Expected loss of f :

performance measure of f

$$\begin{aligned} \text{joint } p_{XY} \quad \mathbb{E}_{XY} [l(f(x), y)] &= \mathbb{E}_{XY} [\mathbb{1}\{f(x) \neq y\}] = \mathbb{E}_X [\mathbb{E}_{Y|X} [\mathbb{1}\{f(x) \neq y\} | X=x]] \\ &\stackrel{(\text{def})}{=} \sum_{i=1}^K p(Y=i | X=x) \cdot \mathbb{1}\{f(x) \neq i\} \\ &= \sum_{i \neq f(x)} p(Y=i | X=x) \quad \text{prob of wrong pred} \\ &= \left(\sum_{i=1}^K p(Y=i | X=x) \right) - p(Y=f(x) | X=x) \quad \text{prob of correct pred} \end{aligned}$$

Classification

- Learn $f: X \rightarrow Y$
 - X - features
 - Y - target classes

- Loss Function $\ell(f(x), y) = \mathbf{1}\{f(x) \neq y\}$

- Expected loss of f :

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\begin{aligned}\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] &= \sum_i P(Y = i|X = x) \mathbf{1}\{f(x) \neq i\} = \sum_{i \neq f(x)} P(Y = i|X = x) \\ &= 1 - P(Y = f(x)|X = x)\end{aligned}$$

- Suppose you knew $P(Y|X)$ exactly, how should you classify?

Classification

- Learn $f: X \rightarrow Y$
 - X - features
 - Y - target classes

- **Loss Function** $\ell(f(x), y) = \mathbf{1}\{f(x) \neq y\}$

- **Expected loss of f :**

$$\mathbb{E}_{XY}[\mathbf{1}\{f(X) \neq Y\}] = \mathbb{E}_X[\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x]]$$

$$\mathbb{E}_{Y|X}[\mathbf{1}\{f(x) \neq Y\}|X = x] = \sum_i P(Y = i|X = x)\mathbf{1}\{f(x) \neq i\} = \sum_{i \neq f(x)} P(Y = i|X = x) \\ = 1 - P(Y = f(x)|X = x)$$

- **Suppose you knew $P(Y|X)$ exactly, how should you classify?**

- **Bayes-Optimal classifier:**

based on Bayes
Theorem

$$f(x) = \arg \max_y \mathbb{P}(Y = y|X = x)$$

minimize ℓ / loss

$$= \frac{P(x, y)}{P(x)}$$

Bayes Optimal Binary Classifier

$$Y \in \{0, 1\}$$

- Suppose you knew $P(Y|X)$ exactly, how should you classify?
 - Bayes-Optimal classifier:

$$f(x) = \arg \max_y \mathbb{P}(Y = y | X = x)$$

- Suppose we don't know $P(Y|X)$, but have n iid examples

$$\{(x_i, y_i)\}_{i=1}^n$$

- What is a natural estimator for $P(Y | X)$?

Bayes Optimal Binary Classifier

- Suppose we don't know $P(Y|X)$, but have n iid examples

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n \quad \mathcal{D} \sim (X, Y)^n \quad Y \in \{0, 1\}$$

- What is a natural estimator for $P(Y | X)$?

Fix some $\tilde{x} \in X$, (x', y) $x' \neq \tilde{x}$

Suppose $x_i = \tilde{x}$ for $m \leq n$ samples

What is a natural estimator for $\theta_* := \mathbb{P}(Y = 1 | X = \tilde{x})$?

If k of the m labels are equal to $Y = 1$ then

$$\hat{\beta}(Y=1 | X=\tilde{x}) = \frac{k}{m}$$

$$\begin{aligned} P(Y=1) &= \frac{P(X,Y)}{P(X)} \\ k/n &\approx P(X,Y) \\ m/n &\approx P(\tilde{x}) \end{aligned}$$

$$\star \mathbb{E}_{\mathcal{D}} [\hat{\beta}(Y=1 | X=\tilde{x})] = P(Y=1 | X=\tilde{x}) \quad \text{unbiased}$$

Bayes Optimal Binary Classifier

- Suppose we don't know $P(Y|X)$, but have n iid examples

$$\{(x_i, y_i)\}_{i=1}^n$$

$$Y \in \{0, 1\}$$

- What is a natural estimator for $\arg\max_y P(Y = y | X)$?

If $X = \{0, 1\}^d$, or is generally discrete

$$\hat{f}(x) = \arg\max_{y \in \{0, 1\}} \frac{\sum_{i=1}^n \mathbf{1}[x_i = x, y_i = y]}{\sum_{i=1}^n \mathbf{1}[x_i = x]}$$

Issues?

- (1) May not see all (x, y) pairs $2^d \cdot 2$ possibilities
 - (2) may not even see some x , 2^d
 - (3) $\hat{P}(y/x) \approx P(y/x)$ require a lot of data for each (x, y)
 - \Rightarrow require HUGE n
 - (4) continuous X
- ML: use a model for $P(Y|X) \Rightarrow$ predict on unseen data

Process

Collect a **dataset** $\{(x_i, y_i)\}_{i=1}^n$, $y \in \{0, 1\}$

Decide on a **model** $f: X \rightarrow \mathcal{P}(Y=1|X)$, $f \in \mathcal{F}$

Find the function which fits the data best

Choose a loss function $\ell \in \{f \mapsto \ell(f, y)\}$

Pick the function which minimizes loss on data

$$\frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i) \Rightarrow f$$

Use function to make prediction on new

examples x_{new}

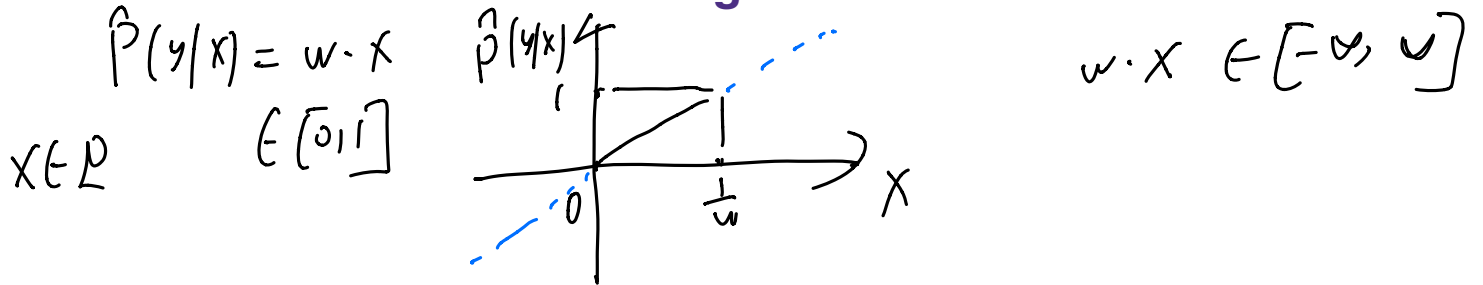
$$\hat{p}(Y=1|X=x) = f(x_{\text{new}})$$

Decide on a model, Binary Classification

To make predictions for unseen inputs (x s),

need a **general** model for $\mathbb{P}(Y = 1|X = x)$

- What about standard linear regression model?



$$\sigma(\cdot) : [-\infty, \infty] \rightarrow [0, 1]$$

- Need to map real values to $[0, 1]$

- We call such maps “link functions”

$$f(x) = \sigma(w^T x)$$

Logistic Regression

$$x \in \mathbb{R}^d, w \in \mathbb{R}^d$$

$$z = -\infty \Rightarrow \sigma(z) = 0$$

$$z = \infty \Rightarrow \sigma(z) = 1$$

Actually classification, not regression :)

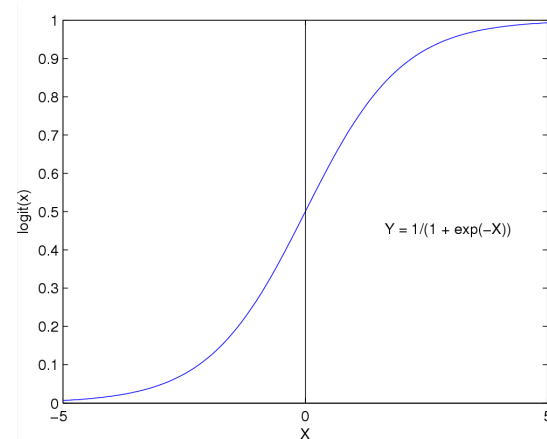
Learn $\mathbb{P}(Y = 1|X = x)$ using $\sigma(w^T x)$, for link function $\sigma =$

Logistic function(or Sigmoid): *"S-shaped"*
 $\sigma(z) = \frac{1}{1 + \exp(-z)}$

$$\mathbb{P}[Y = 1|X = x, w] = \sigma(w^T x) = \frac{1}{1 + \exp(-w^T x)}$$

$$\mathbb{P}[Y = 0|X = x, w] = 1 - \sigma(w^T x) = \frac{\exp(-w^T x)}{1 + \exp(-w^T x)}$$

$$= \frac{1}{1 + \exp(w^T x)}$$



Features can be discrete or continuous!

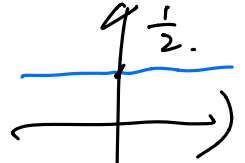
$$x \mapsto \frac{\exp(-w_j^T x)}{\sum_{i=1}^K \exp(-w_i^T x)}$$

$$\hat{p}(y=j|x) = \frac{\exp(-w_j^T x)}{\sum_{i=1}^K \exp(-w_i^T x)}$$

Understanding the sigmoid

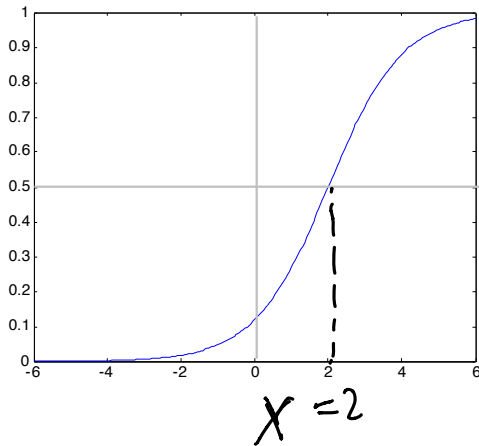
$$\hat{p}(y=1|x) \begin{cases} \geq 0.5 \Rightarrow 1 \\ \leq 0.5 \Rightarrow 0 \end{cases}$$

$$\sigma(w_0 + \sum_k w_k x_k) = \frac{1}{1 + e^{-(w_0 + \sum_k w_k x_k)}}$$

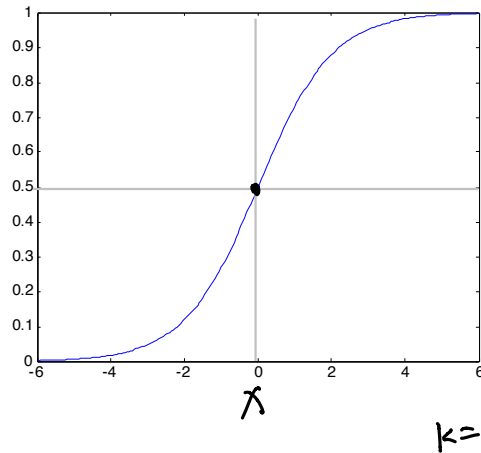


$$\begin{aligned} w_0 &\rightarrow 0 \\ w_1 &\rightarrow 0 \end{aligned} \Rightarrow \frac{1}{2}$$

$$w_0 = -2, w_1 = -1$$



$$\hookrightarrow w_0 = 0, w_1 = -1$$



$$w_0 = 0, w_1 = -0.5$$

