

CS 4824/ECE 4424: Logistic Regression

Acknowledgement:

Many of these slides are derived from Tom Mitchell, Pascal Poupart, Pieter Abbeel, Eric Eaton, Carlos Guestrin, William Cohen, and Andrew Moore.

Logistic Regression

Idea:

- Naïve Bayes allows estimating $P(Y|X)$ by learning $P(Y)$ and $P(X|Y)$
- Why not learn $P(Y|X)$ directly?

Problem setting

$$x_1 = \text{Temp.} \xrightarrow{\mathbb{R}} Y = \text{Play T} \in \mathbb{N}$$

- Consider learning $f: X \rightarrow Y$

- X is a vector of real-valued features $\langle X_1, X_2, \dots, X_n \rangle$

$$x_i \in \mathbb{R}$$

- Y is boolean $Y \in \{0, 1\}$

- Assume all X_i 's are conditionally independent given Y

C.I.

- Model $P(X_i | Y = y_k)$ as Gaussian $\sim \mathcal{N}(\mu_{ik}, \sigma_i)$

not σ_{ik}

- Model $P(Y)$ as Bernoulli (π)

feature x_i class label y_k

- Given that, what's the parametric form of $P(Y|X)$? $\propto P(Y) P(X|Y)$

Parametric form of $P(Y|X)$

$$\begin{aligned}
 P(Y=1|X) &= \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)} \\
 &= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \\
 &= \frac{1}{1 + \exp\left(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}\right)} \quad \begin{array}{l} \exp(\ln x) = x \\ \prod_i P(x_i|Y=0) \end{array} \\
 \hat{P}(Y=1) = \pi &= \frac{1}{1 + \exp\left(\ln \frac{P(Y=0)}{P(Y=1)} + \sum_i \ln \frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)} \quad \text{c.I.} \\
 &= \frac{1}{1 + \exp\left(\ln \frac{1-\pi}{\pi} + \sum_i \ln \frac{P(x_i|Y=0)}{P(x_i|Y=1)}\right)}
 \end{aligned}$$

Parametric form of $P(Y|X)$

$$\begin{aligned}
 \circ \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)} &= \sum_i \ln \frac{\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(\frac{-(x_i - \mu_{i0})^2}{2\sigma_i^2}\right)}{\frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(\frac{-(x_i - \mu_{i1})^2}{2\sigma_i^2}\right)} \quad \left[P(x|y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu_{ik}}{\sigma_{ik}}\right)^2} \right] \\
 &= \sum_i \ln \exp\left(\frac{(x_i - \mu_{i1})^2 - (x_i - \mu_{i0})^2}{2\sigma_i^2}\right) \quad \left[\frac{e^a}{e^b} = e^{a-b} \right] \\
 &= \sum_i \ln \left(\frac{(x_i - \mu_{i1})^2 - (x_i - \mu_{i0})^2}{2\sigma_i^2} \right) \\
 &= \sum_i \ln \left(\frac{(x_i^2 - 2x_i\mu_{i1} + \mu_{i1}^2) - (x_i^2 - 2x_i\mu_{i0} + \mu_{i0}^2)}{2\sigma_i^2} \right) \\
 &= \sum_i \ln \left(\frac{2x_i(\mu_{i0} - \mu_{i1}) + \mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right) \quad \left[\sum_i \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right]
 \end{aligned}$$

Parametric form of $P(Y|X)$

Therefore, $P(Y = 1 | X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

where, $w_0 = \ln \frac{1 - \pi}{\pi} + \sum_i \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}$; and $w_i = \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2}$ for $i = 1 \dots n$

Very convenient!

- $P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

- implies

- $P(Y = 0 | X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_{i=1}^n x_i w_i)}{1 + \exp(w_0 + \sum_{i=1}^n x_i w_i)}$

- implies

- $\frac{P(Y = 0 | X)}{P(Y = 1 | X)} = \exp(w_0 + \sum_{i=1}^n x_i w_i) \geq 1$

linear decision rule.

- or equivalently

- $\ln \frac{P(Y = 0 | X)}{P(Y = 1 | X)} = w_0 + \sum_{i=1}^n x_i w_i \geq 0$

Very convenient!

- $P(Y = 1 | X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

◦ implies

- $P(Y = 0 | X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_{i=1}^n w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

◦ implies

- $\frac{P(Y = 0 | X)}{P(Y = 1 | X)} = \exp(w_0 + \sum_{i=1}^n w_i x_i)$

linear classification rule!

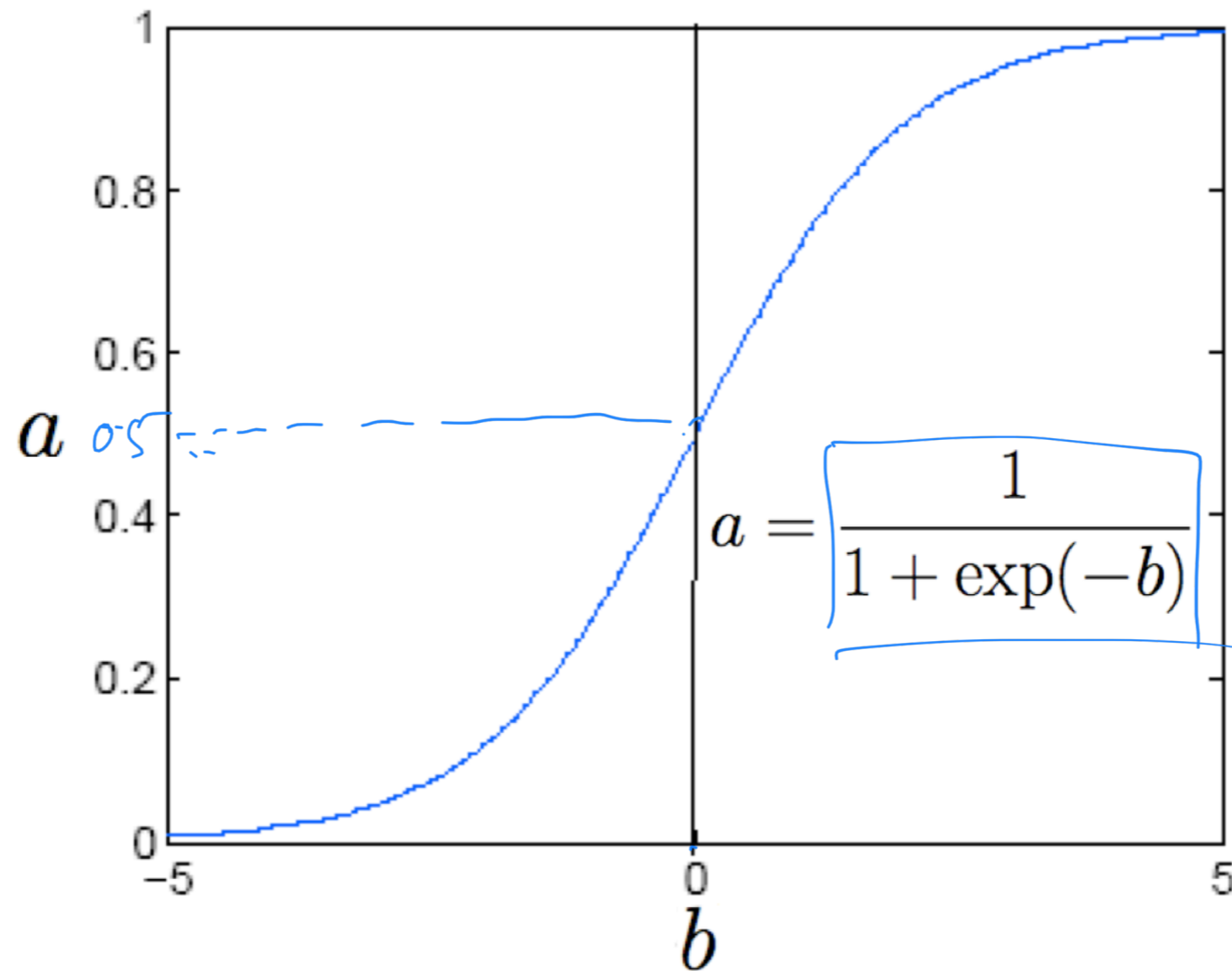
◦ or equivalently

- $\ln \frac{P(Y = 0 | X)}{P(Y = 1 | X)} = w_0 + \sum_{i=1}^n w_i x_i$

dot product of weights and the features

$$a \cdot b = \sum_{i=1}^n a_i b_i$$

Logistic function



$$P(Y = 1 | X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

Logistic regression more generally

- Logistic regression when Y not boolean, but still discrete valued
- Now $Y \in \{y_1, \dots, y_R\}$ and so we need to learn $R-1$ sets of weights

- for $k < R$:
$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^n w_{ki} x_i)}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} x_i)}$$

- for $k = R$:
$$P(Y = y_R | X) = \frac{1}{1 + \sum_{j=1}^{R-1} \exp(w_{j0} + \sum_{i=1}^n w_{ji} x_i)}$$

Training logistic regression: MCLE

- We have L training examples $\{\langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle\}$
- Maximum likelihood estimate (MLE) for parameters W

- $W_{MLE} = \arg \max_W P(\langle X^1, Y^1 \rangle \dots \langle X^L, Y^L \rangle | W)$

$w_0, w_i, \forall i$

- $= \arg \max_W \prod_l P(\langle X^l, Y^l \rangle | W)$

- Maximum conditional likelihood estimate (MCLE)

Data conditional likelihood $\prod_l P(Y^l | X^l, W)$

$$W_{MCLE} = \arg \max_W \prod_l P(Y^l | X^l, W)$$

Training logistic regression: MCLE

- We have L training examples $\{\langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle\}$
- Maximum likelihood estimate (MLE) for parameters W
 - $W_{MLE} = \arg \max_W P(\langle X^1, Y^1 \rangle \dots \langle X^L, Y^L \rangle | W)$
 - $= \arg \max_W \prod_l P(\langle X^l, Y^l \rangle | W)$
- Maximum conditional likelihood estimate (MCLE)

Training logistic regression: MCLE

- We need to choose $W = \langle w_0, \dots, w_n \rangle$ to maximize the conditional likelihood of training data

- where $P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

- and $P(Y = 1 | X, W) = \frac{\exp(w_0 + \sum_{i=1}^n w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$

- Training data $D = \{ \langle X^1, Y^1 \rangle, \dots, \langle X^L, Y^L \rangle \}$

- Data likelihood is $\prod_l P(\langle X^l, Y^l \rangle | W)$

- Data conditional likelihood is $\prod_l P(Y^l | X^l, W)$

- Therefore we need to estimate $W_{MCLE} = \arg \max_W \prod_l P(Y^l | X^l, W)$

Expressing conditional log likelihood

$$l(W) = \ln \prod_l P(Y^l | X^l, W) = \sum_l \ln P(Y^l | X^l, W)$$

$$\text{where } P(Y = 0 | X, W) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

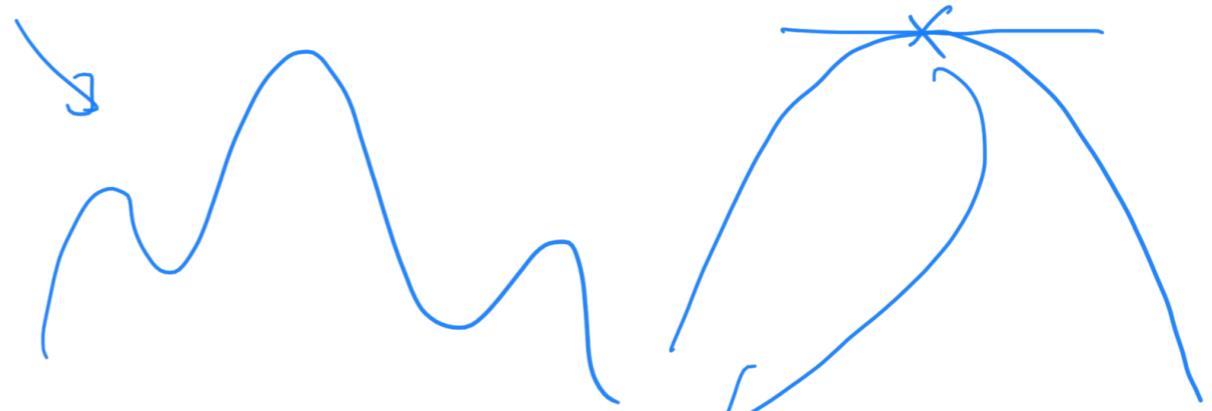
$$\text{and } P(Y = 1 | X, W) = \frac{\exp(w_0 + \sum_{i=1}^n w_i x_i)}{1 + \exp(w_0 + \sum_{i=1}^n w_i x_i)}$$

$$\begin{aligned} l(W) &= \sum_l \left[Y^l \ln P(Y^l = 1 | X^l, W) + (1 - Y^l) \ln P(Y^l = 0 | X^l, W) \right] \leftarrow \\ &= \sum_l Y^l \ln \frac{P(Y^l = 1 | X^l, W)}{P(Y^l = 0 | X^l, W)} + \ln P(Y^l = 0 | X^l, W) \\ &= \sum_l Y^l (w_0 + \sum_i w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i w_i X_i^l)) \end{aligned}$$

Maximizing conditional log likelihood

$$l(W) = \ln \prod_l P(Y^l | X^l, W)$$

$$= \sum_l Y^l (w_0 + \sum_i^n w_i X_i^l) - \ln(1 + \exp(w_0 + \sum_i^n w_i X_i^l))$$



- **Good news:** $l(W)$ is a concave function of W
- **Bad news:** no closed-form solution to maximize $l(W)$

What do we do?

Optimization