# 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

- Consider learning $f: X \to Y$
  - X is a vector of real-valued features $<X_1, X_2, \ldots, X_n>$
  - Y is boolean
  - Assume all $X_i$'s are conditionally independent given Y
  - Model $P(X_i \mid Y = y_k)$ as Gaussian $\sim \mathcal{N}(\mu_{ik}, \sigma_i)$
  - Model P(Y) as Bernoulli ($\pi$)

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

# Parametric form of P(Y|X)

○ $P(Y=1|X) = \dfrac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$

# Parametric form of P(Y|X)

$$P(x \mid y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{\left(-\frac{1}{2}\left(\frac{x-\mu_{ik}}{\sigma_{ik}}\right)^2\right)}$$

○ $\displaystyle\sum_i \ln \frac{P(X_i \mid Y = 0)}{P(X_i \mid Y = 1)}$

# Parametric form of P(Y|X)

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

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

# Very convenient!

- $P(Y = 1 \mid X = \, <X_1, \ldots, X_n> \,) = \dfrac{1}{1 + exp(w_0 + \sum_{i=1}^{n} w_i x_i)}$

- implies
  - $P(Y = 0 \mid X = \, <X_1, \ldots, X_n> \,) =$

- implies

  - $\dfrac{P(Y = 0 \mid X)}{P(Y = 1 \mid X)} =$

- or equivalently

  - $\ln \dfrac{P(Y = 0 \mid X)}{P(Y = 1 \mid X)} =$

# Very convenient!

- $P(Y = 1 \mid X = <X_1, \ldots, X_n>) = \dfrac{1}{1 + exp(w_0 + \sum_{i=1}^{n} w_i x_i)}$

- implies

- $P(Y = 0 \mid X = <X_1, \ldots, X_n>) = \dfrac{exp(w_0 + \sum_{i=1}^{n} w_i x_i)}{1 + exp(w_0 + \sum_{i=1}^{n} w_i x_i)}$

- implies

- $\dfrac{P(Y = 0 \mid X)}{P(Y = 1 \mid X)} = exp\left(w_0 + \sum_{i=1}^{n} w_i x_i\right)$
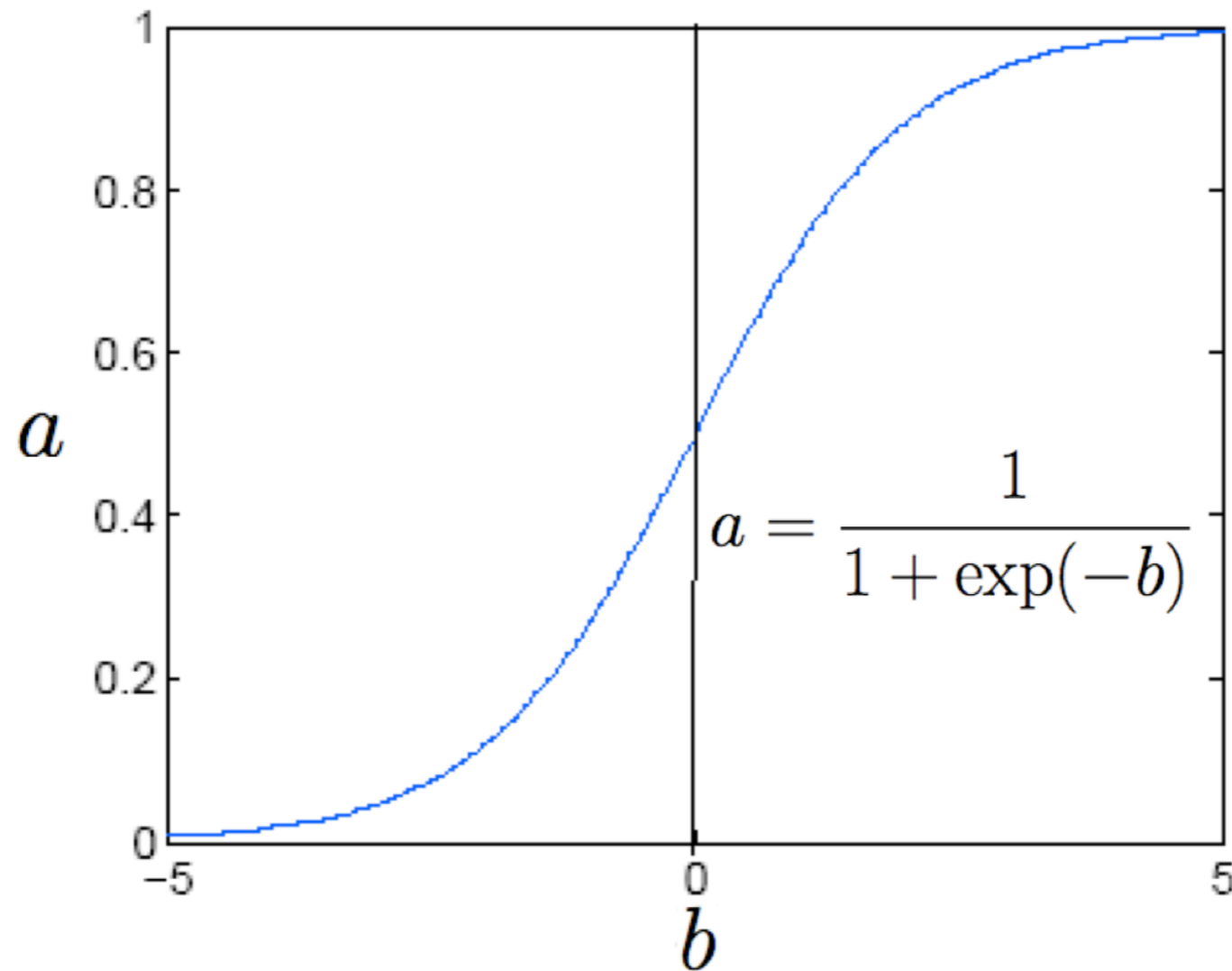
**linear classification rule!**

- or equivalently

- $\ln \dfrac{P(Y = 0 \mid X)}{P(Y = 1 \mid X)} = w_0 + \sum_{i=1}^{n} w_i x_i$ ⟵ dot product of weights and the features

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

# Logistic function



$$a = \frac{1}{1 + \exp(-b)}$$

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

# Logistic regression more generally

◦ Logistic regression when Y not boolean, but still discrete valued

◦ Now $Y \in \{y_1, \ldots y_R\}$ and so we need to learn $R$-$1$ sets of weights

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

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

# Training logistic regression: MCLE

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

  - $$W_{MLE} = \arg\max_{W} P( <X^1, Y^1> \ldots <X^L, Y^L> \mid W)$$

  - $$= \arg\max_{W} \prod_{l} P( <X^l, Y^l> \mid W)$$

- Maximum <u>conditional</u> likelihood estimate (MCLE)

# Training logistic regression: MCLE

- We have L training examples $\{<X^1, Y^1>,\ldots, <X^L, Y^L>\}$
- Maximum likelihood estimate (MLE) for parameters W
    - $W_{MLE} = \arg\max_{W} P(<X^1, Y^1> \ldots <X^L, Y^L> \,|\, W)$
    - $= \arg\max_{W} \prod_{l} P(<X^l, Y^l> \,|\, W)$
- Maximum <u>conditional</u> likelihood estimate (MCLE)

# Training logistic regression: MCLE

◦ We need to choose W = <$w_0$,…,$w_n$> to <u>maximize the conditional likelihood</u> of training data

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

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

◦ Training data D = {<$X^1, Y^1$>,…, <$X^L, Y^L$>}

◦ Data likelihood is $\displaystyle\prod_l P(< X^l, Y^l > \,|\, W)$

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

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

# Expressing conditional log likelihood

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

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

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

- $l(W) = \sum_l Y^l \ln P(Y^l = 1 \,|\, X^l, W) + (1 - Y^l)\ln P(Y^l = 0 \,|\, X^l, W)$

  $= \sum_l Y^l \ln \dfrac{P(Y^l = 1 \,|\, X^l, W)}{P(Y^l = 0 \,|\, X^l, W)} + \ln P(Y^l = 0 \,|\, X^l, W)$

  $= \sum_l Y^l (w_0 + \sum_{i}^{n} w_i X_i^l) - \ln(1 + exp(w_0 + \sum_{i}^{n} w_i X_i^l))$

# Maximizing conditional log likelihood

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

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

- **Good news**: $l(W)$ is a concave function of $W$

- **Bad news**: no closed-form solution to maximize $l(W)$

### What do we do?