# Logistic Regression

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

September 24th, 2007

**1**

---

# Generative v. Discriminative classifiers – Intuition

- **Want to Learn**: h:$X \mapsto Y$  ∈ $\{1,2,3,\ldots,k\}$
  - □ **X** – features
  - □ Y – target classes
- **Bayes optimal classifier** – P(Y|**X**)    *exactly*
- **Generative classifier**, e.g., Naïve Bayes:
  - □ Assume some **functional form for P(X|Y), P(Y)**
  - □ Estimate parameters of P(**X**|Y), P(Y) directly from training data
  - □ Use Bayes rule to calculate P(Y|X= x) $= \dfrac{P(Y, X=x)}{P(X=x)}$
  - □ This is a '***generative***' model
    - ■ **Indirect** computation of P(Y|X) through Bayes rule
    - ■ But, **can generate a sample of the data**, P(X) = $\sum_y$ P(y) P(X|y)
- **Discriminative classifiers**, e.g., Logistic Regression:
  - □ Assume some **functional form for P(Y|X)**
  - □ Estimate parameters of P(Y|X) directly from training data
  - □ This is the '***discriminative***' model
    - ■ Directly learn P(Y|X)
    - ■ But **cannot obtain a sample of the data**, because P(X) is not available

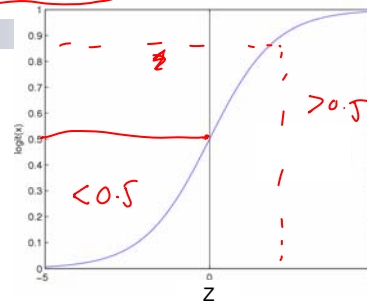**2**

1

# Logistic Regression

**Logistic function (or Sigmoid):** $\dfrac{1}{1 + exp(-z)}$

- Learn $P(Y|\mathbf{X})$ directly!
  - Assume a particular functional form
  - Sigmoid applied to a linear function of the data:

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

$z$

$> 0.5$

$< 0.5$

$z$

**Features can be discrete or continuous!**

©Carlos Guestrin 2005-2007
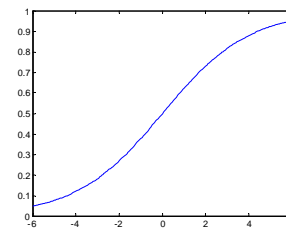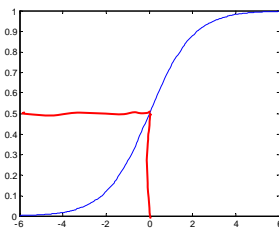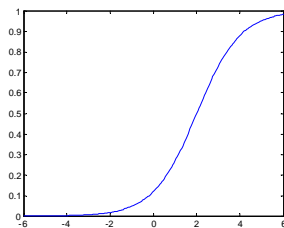
3

---

# Understanding the sigmoid

$$g\left(w_0 + \sum_i w_i x_i\right) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

Constant

$w_0 + w_i x_i$

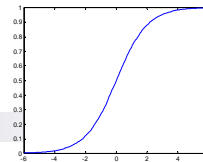| $w_0=-2, w_1=-1$ | $w_0=0, w_1=-1$ | $w_0=0, w_1=-0.5$ |
|---|---|---|

©Carlos Guestrin 2005-2007

4

# Logistic Regression – a Linear classifier

$$g\left(w_0 + \sum_{i=1}^{n} w_i x_i\right) = \frac{1}{1 + e^{w_0 + \sum_i w_i x_i}}$$

$-\infty$

$w_0 + \sum_i w_i x_i < 0$

$x_i$: $g > 0.5$

$+$ true

$w_0 + \sum_i w_i x_i = 0$

$g = 0.5 = \frac{1}{1+1}$

$x$: $w_0 + \sum_i w_i x_i > 0$

$g < 0.5$

false

n-dimensional space

5

---

# Very convenient!

$\ln 1 = 0$

$$P(Y=1 | X = < X_1, ... X_n >, w) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

implies

$= 1 - P(Y=1 | X, w)$

$$P(Y=0 | X = < X_1, ... X_n >, w) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

implies

output $y=1$, if $1 >$

$$\frac{P(Y=0 | X)}{P(Y=1 | X)} = exp\left(w_0 + \sum_i w_i X_i\right)$$

linear classification rule!

implies

$$\ln \frac{P(Y=0 | X)}{P(Y=1 | X)} = w_0 + \sum_i w_i X_i \quad < 0 \quad \text{return } Y=1$$

6

3

# What if we have continuous $X_i$?

Eg., character recognition: $X_i$ is i[th] pixel



*(handwritten annotations: $x_i$ pixel intensity; if class = b $x_i$ high; a $x_i$ low; $\sigma x_i$; mean pixel i class k; variance pixel i class k)*

Gaussian Naïve Bayes (GNB):

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

*(handwritten: labels)*

Sometimes assume variance

- is independent of Y (i.e., $\sigma_i$),
- or independent of $X_i$ (i.e., $\sigma_k$)
- or both (i.e., $\sigma$)

**7**

©Carlos Guestrin 2005-2007

---

# Example: GNB for classifying mental states

[Mitchell et al.]



**~1 mm resolution**

**~2 images per sec.**

**15,000 voxels/image**

**non-invasive, safe**

**measures Blood Oxygen Level Dependent (BOLD) response**

**Typical impulse response**

**8**

©Carlos Guestrin 2005-2007

## Learned Bayes Models – <u>Means</u> for P(BrainActivity | WordCategory)

[Mitchell et al.]

Pairwise classification accuracy: <u>85%</u>
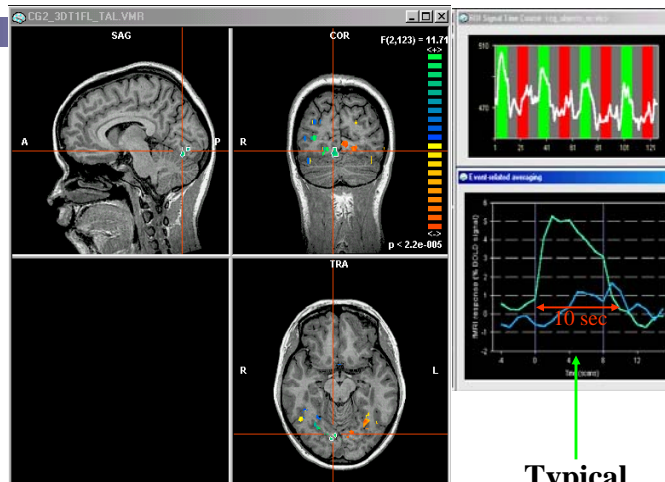
<u>People words</u>          Animal words

9

---

# Logistic regression v. Naïve Bayes

- Consider learning f: <u>X → Y</u>, where
  - □ X <u>is a vector of real-valued features, < X$_1$ … X$_n$ ></u>
  - □ Y is boolean
- Could use a <u>Gaussian Naïve Bayes classifier</u>
  - □ assume all X$_i$ are conditionally independent given Y
  - □ model P(X$_i$ | Y = y$_k$) as Gaussian N($\mu_{ik}$,$\sigma_i$)    *variance only dependson x$_i$ on pixel; not on class*
  - □ <u>model P(Y) as Bernoulli($\theta$,1-$\theta$)</u>
- What does that imply about the form of P(Y|X)?

$$P(Y = 1|X =< X_1, ... X_n >) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

**Cool!!!!**

10

5

# Derive form for P(Y|X) for continuous $X_i$

$e^{\ln x} = x$  $\quad P(y=1|X) = \frac{1}{1+e^{w_0 + \sum w_i x_i}}$

Bayes rule

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

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

$$= \frac{1}{1 + \exp\left(\ln\frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}\right)}$$

$$= \frac{1}{1 + \exp\left(\left(\ln\frac{1-\theta}{\theta}\right) + \boxed{\sum_i \ln\frac{P(X_i|Y=0)}{P(X_i|Y=1)}}\right)}$$

expand

"looks like $w_0$"
* independent of $x_i$

**11**

©Carlos Guestrin 2005-2007

---

# Ratio of class-conditional probabilities

$\ln\frac{1}{e^{-x}} = x$  $\qquad P(y=1|x) = \frac{1}{1+e^{w_0 + \sum w_i x_i}}$

$i$ indexes over features

$$\ln\frac{P(X_i|Y=0)}{P(X_i|Y=1)} = $$

$$P(X_i = x_i | Y = y_k) = \frac{1}{\sigma_i\sqrt{2\pi}} \, e^{\frac{-(x_i - \mu_{ik})^2}{2\sigma_i^2}}$$

var doesn't depend on class $k$

$$\ln \frac{\frac{1}{\sigma_i\sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}}}{\frac{1}{\sigma_i\sqrt{2\pi}} e^{-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}}} =$$

$w_i$  $\qquad$ also part of $w_0$

$$\frac{-(x_i - \mu_{i0})^2}{2\sigma_i^2} + \frac{(x_i - \mu_{i1})^2}{2\sigma_i^2} = \frac{(\mu_{i0} - \mu_{i1})x_i}{\sigma_i^2} + \frac{(\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2}$$

$$= \frac{-x_i^2 + 2x_i\mu_{i0} - \mu_{i0}^2 + x_i^2 - 2\mu_{i1}x_i + \mu_{i1}^2}{2\sigma_i^2}$$

**12**

©Carlos Guestrin 2005-2007

6

# Derive form for $P(Y|X)$ for continuous $X_i$

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

$$= \frac{1}{1 + \exp(\ (\ln \frac{1-\theta}{\theta}) + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)})}$$

$$\sum_i \left( \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right)$$

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

$$w_0 = \ln \frac{1-\theta}{\theta} + \sum_i \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2}$$

**13**

©Carlos Guestrin 2005-2007

---

# Gaussian Naïve Bayes v. Logistic Regression

*transform into parameterization of LR*

**Set of Gaussian Naïve Bayes parameters (feature variance independent of class label)**

**Set of Logistic Regression parameters**

*transform to NB ∃ some $w$'s, but not all $w$'s*

- Representation equivalence
  - **But only in a special case!!!** (GNB with class-independent variances)
- But what's the difference???
- **LR makes no assumptions about** P(**X**|Y) **in learning**!!!  *does not assume independence*
- **Loss function!!!**  *assume form for $P(Y|X)$*
  - Optimize different functions → Obtain different solutions

**14**

©Carlos Guestrin 2005-2007

7

# Logistic regression for more than 2 classes

- Logistic regression in more general case, where $Y \in \{Y_1 ... Y_R\}$ : learn $R$-1 sets of weights

4 class: 3 sets of params

$$P(Y_1 = 1 | x, w_1) \propto e^{w_{10} + \sum_i w_{1i} x_i}$$

$$P(Y = 2 | x, w_2) \propto e^{w_{20} + \sum_i w_{2i} x_i}$$

$$\vdots$$

$$P(Y = R-1 | x, w_{R-1}) \propto e^{w_{R-1,0} + \sum_i w_{R-1,i} x_i}$$

$$P(Y = R | x) = 1 - \sum_{j=1}^{R-1} P(Y = j | x) \propto 1 - \sum_{j=1}^{R-1} e^{w_{j0} + \sum_i w_{ji} x_i}$$

**15**

---

# Logistic regression more generally

- Logistic regression in more general case, where $Y \in \{Y_1 ... Y_R\}$ : 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 $k=R$ (normalization, so no weights for this class)

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

**Features can be discrete or continuous!**

$$\begin{cases} A = 1 \\ B = 2 \\ C = 3 \end{cases}$$

$X_i = $ grade in 10701

**16**

# Loss functions: Likelihood v. Conditional Likelihood

*[handwritten top right: $P(x,y|w)$ = $P(y|x,w) \cdot P(x|w)$]*

- Generative (Naïve Bayes) Loss function:
  **Data likelihood**

*[handwritten: $D = \langle x^j, y^j \rangle_{j:1\ldots N}$]*

$$\ln P(\mathcal{D} \mid \mathbf{w}) = \sum_{j=1}^{N} \ln P(\mathbf{x}^j, y^j \mid \mathbf{w})$$

$$= \sum_{j=1}^{N} \ln P(y^j \mid \mathbf{x}^j, \mathbf{w}) + \sum_{j=1}^{N} \ln P(\mathbf{x}^j \mid \mathbf{w})$$

*[handwritten: classification]*  *[handwritten: for generating data, not important for classification]*

- Discriminative models cannot compute $P(\mathbf{x}|\mathbf{w})$!
- But, discriminative (logistic regression) loss function:
  **Conditional Data Likelihood**

*[handwritten: discriminative likelihood]*

$$\ln P(\mathcal{D}_Y \mid \mathcal{D}_{\mathbf{X}}, \mathbf{w}) = \sum_{j=1}^{N} \ln P(y^j \mid \mathbf{x}^j, \mathbf{w})$$

*[handwritten right: $j$ = training exam; $y^j = 1$ if Spam, = 0 if not Spam; $x^j$ = list of words in $j$-th ex]*

  □ Doesn't waste effort learning $P(X)$ – focuses on $P(Y|\mathbf{X})$ all that matters for classification

©Carlos Guestrin 2005-2007

17

---

# Expressing Conditional Log Likelihood

$$\max_{\mathbf{w}} \quad l(\mathbf{w}) \equiv \sum_{j} \ln P(y^j | \mathbf{x}^j, \mathbf{w})$$

$$P(Y = 0 | \mathbf{X}, \mathbf{w}) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | \mathbf{X}, \mathbf{w}) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

*[handwritten: $j$th component: $P(Y=Spam|X^j,w)$ if $j$-th was spam; $P(Y = not\ spam|x^j,w)$ if $j$ was not spam]*

$$l(\mathbf{w}) = \sum_{j} \left[ y^j \ln P(y = 1 | \mathbf{x}^j, \mathbf{w}) + (1 - y^j) \ln P(y = 0 | \mathbf{x}^j, \mathbf{w}) \right]$$

*[handwritten:
if $y^j = 1$ : $\ln P(y=1|x^j w)$ + 0
if $y^j = 0$ : 0 + $\ln P(y=0|x^j,w)$]*

*[handwritten: $l(w) = \sum_j \{ y_j [w_0 + \sum_i w_i x_i - \ln(1 + e^{w_0 + \sum_i w_i x_i})] + (1 - y^j)[- \ln(1 + e^{w_0 + \sum_i w_i x_i})] \}$]*

©Carlos Guestrin 2005-2007

18

9

# Maximizing Conditional Log Likelihood

$$P(Y = 0|X, W) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|X, W) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$\max_{w} \quad l(\mathbf{w}) \equiv \ln \prod_j P(y^j|\mathbf{x}^j, \mathbf{w})$$

linear part

$$= \sum_j \left[ y^j(w_0 + \sum_i^n w_i x_i^j) - \ln(1 + exp(w_0 + \sum_i^n w_i x_i^j)) \right]$$

**Good news**: $l(\mathbf{w})$ is concave function of $\mathbf{w} \rightarrow$ no locally optimal solutions

**Bad news**: no closed-form solution to maximize $l(\mathbf{w})$

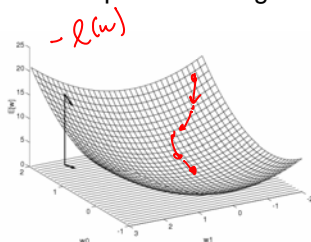**Good news**: concave functions easy to optimize

19

---

# Optimizing concave function – Gradient ascent

( Conjugate G.D.)
better.

- Conditional likelihood for Logistic Regression is concave $\rightarrow$ Find optimum with gradient ascent

$-l(w)$

**Gradient:**
$$\nabla_\mathbf{w} l(\mathbf{w}) = [\frac{\partial l(\mathbf{w})}{\partial w_0}, \ldots, \frac{\partial l(\mathbf{w})}{\partial w_n}]'$$

step size

**Learning rate, $\eta > 0$**

0.01

**Update rule:**
$$\triangle \mathbf{w} = \eta \nabla_\mathbf{w} l(\mathbf{w})$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}$$

- Gradient ascent is simplest of optimization approaches
  - e.g., Conjugate gradient ascent much better (see reading)

20

10

# Maximize Conditional Log Likelihood: Gradient ascent

$$\frac{\partial \ln f(w)}{\partial w} = \frac{\frac{\partial f(w)}{\partial w}}{f(w)}$$

$$P(Y = 0 | X, W) = \frac{1}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1 | X, W) = \frac{exp(w_0 + \sum_i w_i X_i)}{1 + exp(w_0 + \sum_i w_i X_i)}$$

$$\frac{\partial e^{f(w)}}{\partial w} = \frac{\partial f(w)}{\partial w} e^{f(w)} \qquad \frac{\partial}{\partial w_1} = x_i^j$$

$$l(\mathbf{w}) = \sum_j y^j (w_0 + \sum_i^n w_i x_i^j) - \ln(1 + exp(w_0 + \sum_i^n w_i x_i^j))$$

$$\frac{\partial l(w)}{\partial w_1} = \sum_j \left[ y^j x_i^j - \frac{\partial}{\partial w_1} \ln\left( 1 + e^{w_0 + \sum w_i x_i^j} \right) \right]$$

$$= \sum_j \left[ y^j x_i^j - \frac{x_i^j e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}} \right]$$

$$P(Y=1|XW)$$

$$= \sum_j x_i^j \left[ y^j - P(Y=1|X,w) \right]$$

if $j^{th}$ example is positive: if $x_i^j$ is positive want to make

if $j^{th}$ '' is negative: if $x_i^j$ is positive $w_i$ large want to make $w_i$ small

©Carlos Guestrin 2005-2007

**21**

---

# Gradient Descent for LR

$w_0$

$w_1$

Iterations

Gradient ascent algorithm: iterate until change $< \varepsilon$

$w^{(t)}:$ $w$ at $t'th$ iteration

$$w_0^{(t+1)} \leftarrow w_0^{(t)} + \eta \sum_j [y^j - \tilde{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w})^{(t)}]$$

no $x_0^j$, $x_0^j = 1$

For $i = 1 \ldots n,$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \sum_j x_i^j [y^j - \tilde{P}(Y^j = 1 | \mathbf{x}^j, \mathbf{w})^{(t)}]$$

$$\frac{e^{w_0 + \sum_i w_i x_i^j}}{1 + e^{w_0 + \sum_i w_i x_i^j}}$$

repeat

**22**

©Carlos Guestrin 2005-2007