

Logistic Regression

Machine Learning – 10701/15781

Carlos Guestrin

Carnegie Mellon University

September 24th, 2007

©Carlos Guestrin 2005-2007

1

Generative v. Discriminative classifiers – Intuition

■ Want to Learn: $h: X \mapsto Y$ $Y \in \{1, 2, 3, \dots, k\}$

- X – features
- Y – target classes

■ Bayes optimal classifier – $P(Y|X)$

■ 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 $P(Y|X=x) = \frac{P(Y, X=x)}{P(X=x)}$
- 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

©Carlos Guestrin 2005-2007

2

generate spm:
sample (or set)
 $P(Y = \text{spam})$

sample words:

$P(X|Y = \text{spam})$

exactly

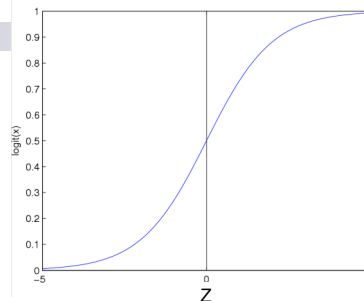
eg. NB:
 $P(X|Y) = \prod P(X_i|Y)$

learn $P(Y, X)$

at classification time:
input x
answer $P(Y|X=x)$

Logistic Regression

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



■ Learn $P(Y|X)$ directly!

- Assume a particular functional form
- Sigmoid applied to a linear function of the data:

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

Features can be discrete or continuous!

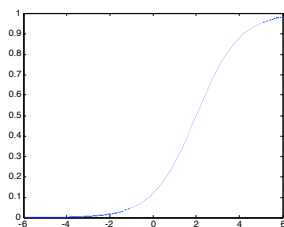
©Carlos Guestrin 2005-2007

3

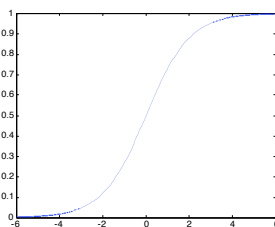
Understanding the sigmoid

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

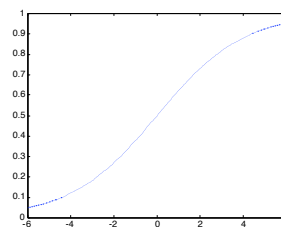
$w_0 = -2, w_1 = -1$



$w_0 = 0, w_1 = -1$



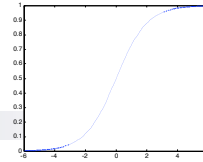
$w_0 = 0, w_1 = -0.5$



©Carlos Guestrin 2005-2007

4

Logistic Regression – a Linear classifier



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

©Carlos Guestrin 2005-2007

5

Very convenient!

$$P(Y = 1|X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

$$P(Y = 0|X = \langle X_1, \dots, X_n \rangle) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

implies

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

implies

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

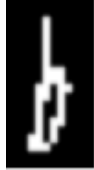
linear
classification
rule!

©Carlos Guestrin 2005-2007

6

What if we have continuous X_i ?

Eg., character recognition: X_i is i^{th} pixel



Gaussian Naïve Bayes (GNB):

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

Sometimes assume variance

- is independent of Y (i.e., σ_i),
- or independent of X_i (i.e., σ_k)
- or both (i.e., σ)

©Carlos Guestrin 2005-2007

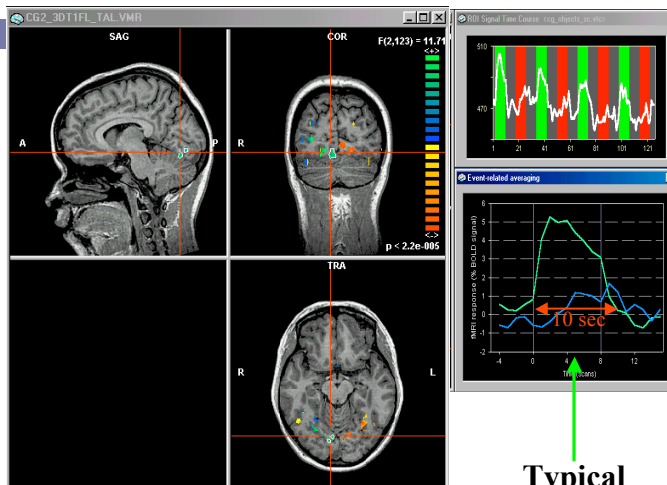
7

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 – Means for $P(\text{BrainActivity} \mid \text{WordCategory})$

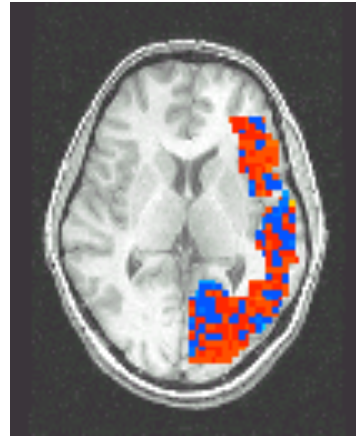
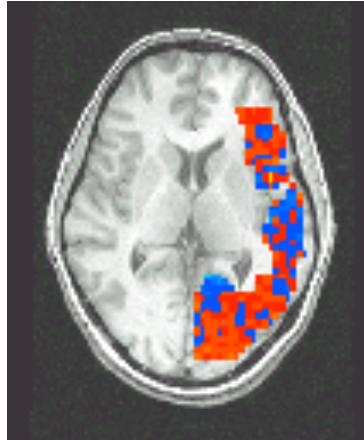
[Mitchell et al.]

Pairwise classification accuracy: 85%

People words



Animal words



9

©Carlos Guestrin 2005-2007

Logistic regression v. Naïve Bayes

- Consider learning $f: X \rightarrow Y$, where
 - X is a vector of real-valued features, $\langle X_1 \dots X_n \rangle$
 - Y is boolean
- Could use a Gaussian Naïve Bayes classifier
 - assume all X_i are conditionally independent given Y
 - model $P(X_i \mid Y = y_k)$ as Gaussian $N(\mu_{ik}, \sigma_i)$
 - model $P(Y)$ as Bernoulli($\theta, 1-\theta$)

- What does that imply about the form of $P(Y|X)$?

$$P(Y = 1 \mid X = \langle X_1, \dots, X_n \rangle) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Cool!!!!

10

©Carlos Guestrin 2005-2007

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

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

©Carlos Guestrin 2005-2007

11

Ratio of class-conditional probabilities

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

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

©Carlos Guestrin 2005-2007

12

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

$$\begin{aligned}
 P(Y = 1|X) &= \frac{P(Y = 1)P(X|Y = 1)}{P(Y = 1)P(X|Y = 1) + P(Y = 0)P(X|Y = 0)} \\
 &= \frac{1}{1 + \exp\left(\ln \frac{1-\theta}{\theta} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)} \\
 &\quad \sum_i \left(\frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} X_i + \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \right) \\
 P(Y = 1|X) &= \frac{1}{1 + \exp(w_0 + \sum_{i=1}^n w_i X_i)}
 \end{aligned}$$

©Carlos Guestrin 2005-2007

13

Gaussian Naïve Bayes v. Logistic Regression

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

**Set of Logistic
Regression parameters**

- 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!!!**
- **Loss function!!!**
 - Optimize different functions → Obtain different solutions

©Carlos Guestrin 2005-2007

14

Logistic regression for more than 2 classes

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

Logistic regression more generally

- Logistic regression in more general case, where $Y \in \{Y_1 \dots Y_R\}$: 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!

Loss functions: Likelihood v. Conditional Likelihood

- Generative (Naïve Bayes) Loss function:

Data likelihood

$$\begin{aligned}\ln P(\mathcal{D} | \mathbf{w}) &= \sum_{j=1}^N \ln P(\mathbf{x}^j, y^j | \mathbf{w}) \\ &= \sum_{j=1}^N \ln P(y^j | \mathbf{x}^j, \mathbf{w}) + \sum_{j=1}^N \ln P(\mathbf{x}^j | \mathbf{w})\end{aligned}$$

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

Conditional Data Likelihood

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

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

17

©Carlos Guestrin 2005-2007

Expressing Conditional Log Likelihood

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

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

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

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

18

©Carlos Guestrin 2005-2007

Maximizing Conditional Log Likelihood



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

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

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

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

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

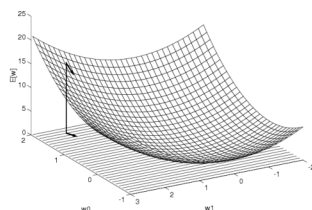
©Carlos Guestrin 2005-2007

19

Optimizing concave function – Gradient ascent



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



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

Learning rate, $\eta > 0$

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

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

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

©Carlos Guestrin 2005-2007

20

Maximize Conditional Log Likelihood: Gradient ascent

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

©Carlos Guestrin 2005-2007

21

Gradient Descent for LR

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

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

For $i = 1 \dots n$,

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

repeat

©Carlos Guestrin 2005-2007

22

That's all M(C)LE. How about MAP?

$$p(\mathbf{w} \mid Y, \mathbf{X}) \propto P(Y \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$$

- One common approach is to define priors on \mathbf{w}
 - Normal distribution, zero mean, identity covariance
 - “Pushes” parameters towards zero
- Corresponds to **Regularization**
 - Helps avoid very large weights and overfitting
 - More on this later in the semester
- MAP estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

©Carlos Guestrin 2005-2007

23

M(C)AP as Regularization

$$\ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$$

$$p(\mathbf{w}) = \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{\frac{-w_i^2}{2\kappa^2}}$$

Penalizes high weights, also applicable in linear regression

©Carlos Guestrin 2005-2007

24

Gradient of M(C)AP

$$\frac{\partial}{\partial w_i} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

$$p(\mathbf{w}) = \prod_i \frac{1}{\kappa \sqrt{2\pi}} e^{\frac{-w_i^2}{2\kappa^2}}$$

25

©Carlos Guestrin 2005-2007

MLE vs MAP

- Maximum conditional likelihood estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[\prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

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

- Maximum conditional a posteriori estimate

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[p(\mathbf{w}) \prod_{j=1}^N P(y^j | \mathbf{x}^j, \mathbf{w}) \right]$$

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

26

©Carlos Guestrin 2005-2007

Naïve Bayes vs Logistic Regression

Consider Y boolean, X_i continuous, $X = \langle X_1 \dots X_n \rangle$

Number of parameters:

- NB: $4n + 1$
- LR: $n + 1$

Estimation method:

- NB parameter estimates are uncoupled
- LR parameter estimates are coupled

27

©Carlos Guestrin 2005-2007

G. Naïve Bayes vs. Logistic Regression 1

[Ng & Jordan, 2002]

- Generative and Discriminative classifiers
- Asymptotic comparison (# training examples \rightarrow infinity)
 - when model correct
 - GNB, LR produce identical classifiers
 - when model incorrect
 - LR is less biased – does not assume conditional independence
 - **therefore LR expected to outperform GNB**

28

©Carlos Guestrin 2005-2007

G. Naïve Bayes vs. Logistic Regression 2

[Ng & Jordan, 2002]

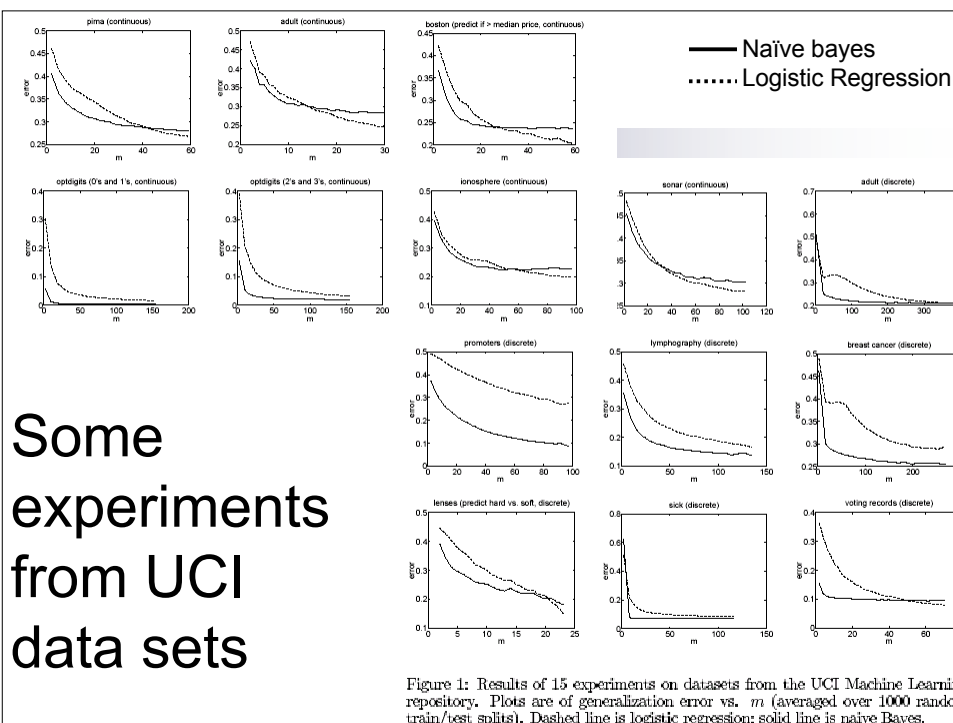
■ Generative and Discriminative classifiers

■ Non-asymptotic analysis

- convergence rate of parameter estimates, $n = \#$ of attributes in X
 - Size of training data to get close to infinite data solution
 - GNB needs $O(\log n)$ samples
 - LR needs $O(n)$ samples
- **GNB converges more quickly to its (perhaps less helpful) asymptotic estimates**

29

©Carlos Guestrin 2005-2007



What you should know about Logistic Regression (LR)

- Gaussian Naïve Bayes with class-independent variances representationally equivalent to LR
 - Solution differs because of objective (loss) function
- In general, NB and LR make different assumptions
 - NB: Features independent given class \rightarrow assumption on $P(\mathbf{X}|Y)$
 - LR: Functional form of $P(Y|\mathbf{X})$, no assumption on $P(\mathbf{X}|Y)$
- LR is a linear classifier
 - decision rule is a hyperplane
- LR optimized by conditional likelihood
 - no closed-form solution
 - concave \rightarrow global optimum with gradient ascent
 - Maximum conditional a posteriori corresponds to regularization
- Convergence rates
 - GNB (usually) needs less data
 - LR (usually) gets to better solutions in the limit

31

©Carlos Guestrin 2005-2007