

UVA CS 6316

– Fall 2015 Graduate: Machine Learning

Lecture 15: Logistic Regression / Generative vs. Discriminative

Dr. Yanjun Qi

University of Virginia

Department of
Computer Science

10/21/15

1

Where are we ? →

Five major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models

10/21/15

2

Where are we ? →

Three major sections for classification

- We can divide the large variety of classification approaches into **roughly three major types**

- Discriminative**
 - directly estimate a decision rule/boundary
 - e.g., **logistic regression**, support vector machine, decisionTree
- Generative:**
 - build a generative statistical model
 - e.g., **naïve bayes classifier**, **Bayesian networks**
- Instance based classifiers**
 - Use observation directly (no models)
 - e.g. **K nearest neighbors**

10/21/15

3

X_1	X_2	X_3	C

A Dataset for classification

$$f : X \rightarrow C$$

Output as Discrete
Class Label
 C_1, C_2, \dots, C_L

Generative → $\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X | C)P(C)$

Discriminative → $P(C | \mathbf{X}) \quad C = c_1, \dots, c_L$

- Data**/points/instances/examples/samples/records: [rows]
- Features**/attributes/dimensions/independent variables/covariates/predictors/regressors: [columns, except the last]
- Target**/outcome/response/label/dependent variable: special column to be predicted [last column]

10/21/15

4

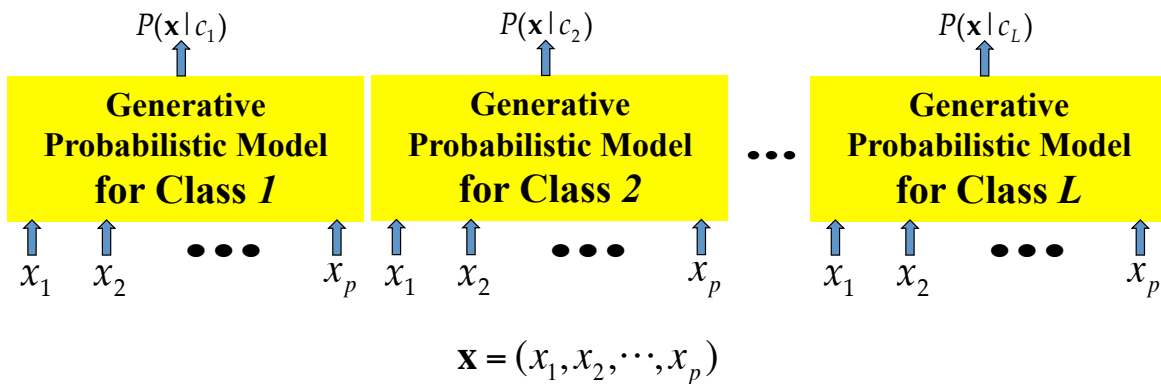
Establishing a probabilistic model for classification (cont.)

Dr. Yanjun Qi / UVA CS 6316 / f15

– **(1) Generative model**

$$\arg \max_c P(C | X) = \arg \max_c P(X, C)$$

$$= \arg \max_c P(X | C) P(C)$$



10/21/15

Adapt from Prof. Ke Chen NB slides

Establishing a probabilistic model for classification

Dr. Yanjun Qi / UVA CS 6316 / f15

– **(2) Discriminative model**

$$P(C | \mathbf{X}) \quad C = c_1, \dots, c_L, \quad \mathbf{X} = (X_1, \dots, X_n)$$

$$P(c_1 | \mathbf{x}) \quad P(c_2 | \mathbf{x}) \quad \dots \quad P(c_L | \mathbf{x})$$



$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

10/21/15

Adapt from Prof. Ke Chen NB slides

Today :

- ✓ Logistic regression
- ✓ Generative vs. Discriminative

10/21/15

7

Multivariate linear regression to Logistic Regression

$$\underline{y} = \underline{\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_i x_i}$$

Dependent

Independent variables

Predicted

Predictor variables

Response variable

Explanatory variables

Outcome variable

Covariables

Logistic regression for
binary classification

$$\ln \left[\frac{P(y|x)}{1 - P(y|x)} \right] = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

10/21/15

8

$$y \in \{0,1\} \quad \ln \left[\frac{P(y=1|x)}{1-P(y=1|x)} \right] = \left[\alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p \right]$$

(1) Linear decision boundary [separate two classes]

$$\ln \frac{P(y=1|x)}{1-P(y=1|x)} = \ln \frac{P(y=1|x)}{P(y=0|x)} = 0$$

$$(2) p(y|x) \Rightarrow \frac{P(y=1|x)}{1-P(y=1|x)} = e^{\beta^T x}$$

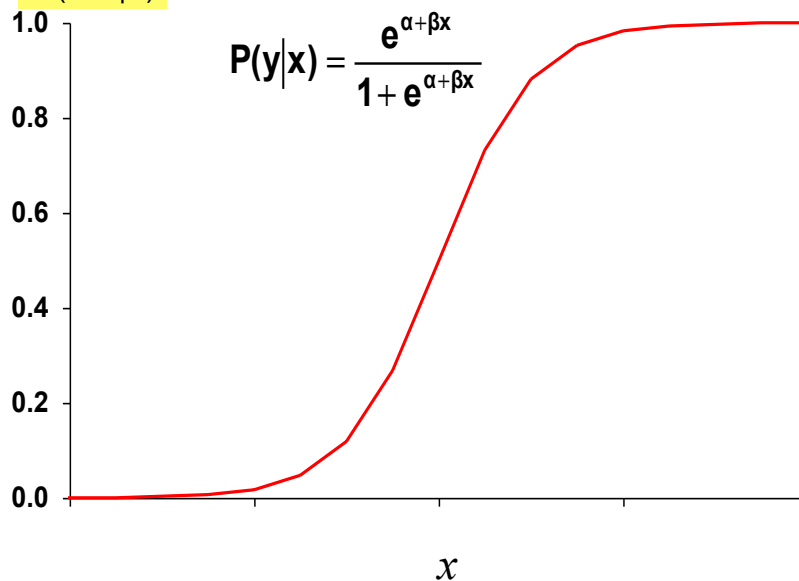
$$\Rightarrow P(y=1|x) = \frac{e^{\beta^T x}}{1 + e^{\beta^T x}}$$

The logistic function (1)

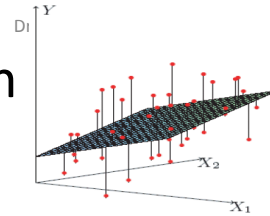
-- is a common "S" shape func

e.g.
Probability of
disease

$P(Y=1|X)$



RECAP: Probabilistic Interpretation of Linear Regression



- Let us assume that the target variable and the inputs are related by the equation:

$$y_i = \theta^T \mathbf{x}_i + \varepsilon_i$$

where ε is an error term of unmodeled effects or random noise

- Now assume that ε follows a Gaussian $N(0, \sigma)$, then we have:

$$p(y_i | x_i; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

- By IID assumption \rightarrow likelihood \rightarrow MLE estimator

$$L(\theta) = \prod_{i=1}^n p(y_i | x_i; \theta) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n \exp\left(-\frac{\sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2}{2\sigma^2}\right)$$

$$l(\theta) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$

10/21/15

11

Dr. Yanjun Qi / UVA CS 6316 / f15

Logistic Regression—when?

Logistic regression models are appropriate for target variable coded as 0/1.

We only observe “0” and “1” for the target variable—but we think of the target variable conceptually as a probability that “1” will occur.

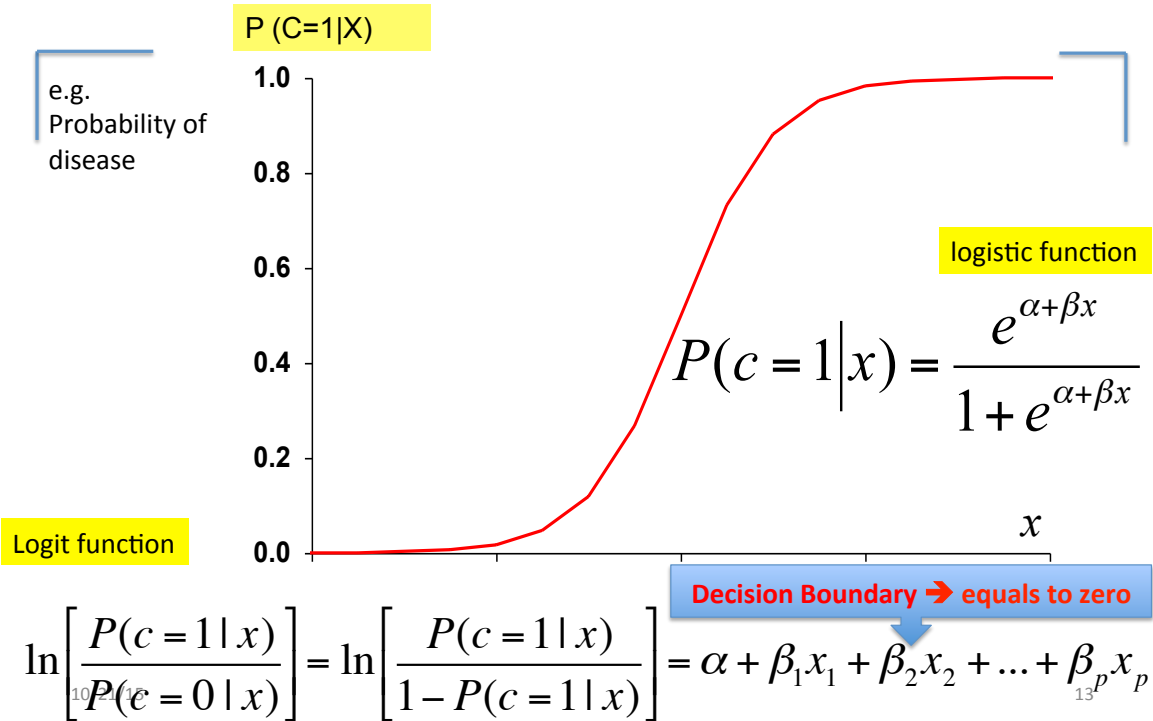
This means we use Bernoulli distribution to model the target variable with its Bernoulli parameter $p(y=1 | x)$ predefined.

The main interest \rightarrow predicting the probability that an event occurs (i.e., the probability that $p(y=1 | x)$).

Binary
 $p(y=1|x)$
logistic

10/21/15

12



The logistic function (2)

$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

logistic

$$\ln \left[\frac{P(y|x)}{1 - P(y|x)} \right] = \alpha + \beta x$$

logit / log-odd



Logit of $P(y|x)$

From probability to logit, i.e. log odds (and back again)

$$z = \log\left(\frac{p}{1-p}\right) \quad \text{logit / log odd function}$$

$$\frac{p}{1-p} = e^z$$

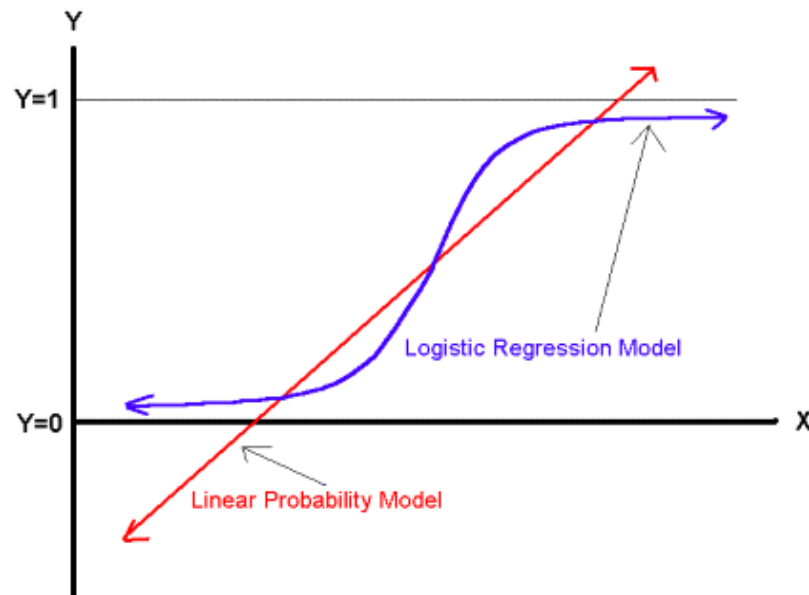
$$p = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}} \quad \text{logistic function}$$

The logistic function (3)

- Advantages of the **logit**
 - Simple transformation of $P(y|x)$
 - Linear relationship with x
 - Can be continuous (Logit between $-\infty$ to $+\infty$)
 - **Directly related to the notion of log odds of target event**

$$\ln\left(\frac{P}{1-P}\right) = \alpha + \beta x \quad \frac{P}{1-P} = e^{\alpha + \beta x}$$

Logistic regression – Binary outcome target variable Y



10/21/15

17

Logistic Regression Assumptions

- Linearity in the logit – the regression equation should have a linear relationship with the logit form of the target variable
- There is no assumption about the feature variables / predictors being linearly related to each other.

10/21/15

18

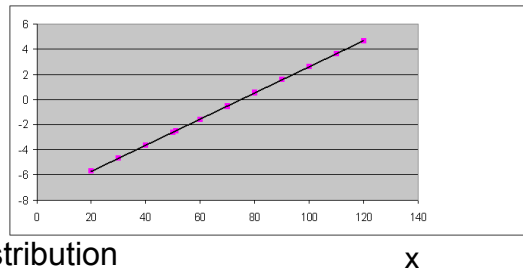
Binary Logistic Regression (K=2)

In summary that the logistic regression tells us two things at once.

- Transformed, the “log odds” (logit) are linear.

$$\ln[p/(1-p)]$$

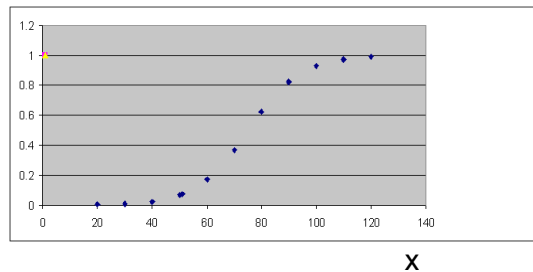
Odds = $p/(1-p)$



This means we use Bernoulli distribution to model the target variable with its Bernoulli parameter $p = p(y=1 | x)$ predefined.

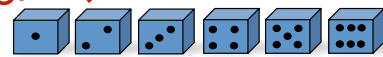
- Logistic Distribution

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



Binary \rightarrow Multinomial Logistic Regression Model

(e.g. $k=6$)



Directly models the posterior probabilities as the output of regression

$$p(y=1|x) = \frac{e^{\beta_1 x}}{1 + e^{\beta_1 x}}$$

$$p(y=0|x) = \frac{1}{1 + e^{\beta_1 x}}$$

$$\Pr(G = k | X = x) = \frac{\exp(\beta_{k0} + \beta_k^T x)}{1 + \sum_{l=1}^{K-1} \exp(\beta_{l0} + \beta_l^T x)}, \quad k = 1, \dots, K-1$$

$$\Pr(G = K | X = x) = \frac{1}{1 + \sum_{l=1}^{K-1} \exp(\beta_{l0} + \beta_l^T x)}$$

$$\ln \frac{P(G=k|x)}{P(G=K|x)} = 0 \Rightarrow \text{linear}$$

$$\beta_{k0} + \beta_k^T x$$

x is p -dimensional input vector

β_k is a p -dimensional vector for each k

Total number of parameters is $(K-1)(p+1)$

Note that the class boundaries are linear

Today :

- ✓ Logistic regression
- ✓ Parameter estimation
- ✓ Generative vs. Discriminative

Parameter Estimation for LR

→ MLE from the data

- **RECAP:** Linear regression → Least squares
- **Logistic regression: → Maximum likelihood estimation**

RECAP: Probabilistic Interpretation of Linear Regression (cont.)

- Hence the log-likelihood is:

$$l(\theta) = n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y_i - \theta^T \mathbf{x}_i)^2$$



MLE

- Do you recognize the last term?

Yes it is:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- Thus under independence assumption, residual square error (RRS) is equivalent to MLE of θ !

10/21/15

23

MLE for Logistic Regression Training

Let's fit the logistic regression model for $K=2$, i.e., number of classes is 2

Training set: (x_i, y_i) , $i=1, \dots, N$

For Bernoulli distribution

$$p(y | x)^y (1 - p)^{1-y}$$

(conditional)
Log-likelihood.

How?

$$\begin{aligned} l(\beta) &= \sum_{i=1}^N \{\log \Pr(Y = y_i | X = x_i)\} \\ &= \sum_{i=1}^N y_i \log(\Pr(Y = 1 | X = x_i)) + (1 - y_i) \log(\Pr(Y = 0 | X = x_i)) \\ &= \sum_{i=1}^N \left(y_i \log \frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)} + (1 - y_i) \log \frac{1}{1 + \exp(\beta^T x_i)} \right) \\ &= \sum_{i=1}^N (y_i \beta^T x_i - \log(1 + \exp(\beta^T x_i))) \end{aligned}$$

$p(y_i | x_i)$

x_i are $(p+1)$ -dimensional input vector with leading entry 1
 β is a $(p+1)$ -dimensional vector

10/21/15 We want to maximize the log-likelihood in order to estimate β

24

$$l(\beta) = \sum_{i=1}^N \{\log \Pr(Y = y_i | X = x_i)\}$$

$$\begin{aligned} & \log \left\{ \Pr(Y = y_i | X = x_i) = \mathcal{P}(y_i | x_i) \right\} \Rightarrow \begin{matrix} y_i = 1 \\ y_i = 0 \end{matrix} \\ & = \log \left\{ \mathcal{P}(y_i = 1 | x) ^{y_i} (1 - \mathcal{P}(y_i = 1 | x_i))^{1 - y_i} \right\} \\ & = y_i \log \mathcal{P}(y_i = 1 | x) + (1 - y_i) \log (1 - \mathcal{P}(y_i = 1 | x)) \end{aligned}$$

Newton-Raphson for LR (optional)

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^N \left(y_i - \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)} \right) x_i = 0$$

($p+1$) Non-linear equations to solve for ($p+1$) unknowns β

Solve by Newton-Raphson method:

$$\beta^{new} \leftarrow \beta^{old} - \left[\left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} \right) \right]^{-1} \frac{\partial l(\beta)}{\partial \beta},$$

$$\text{where, } \left(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T} \right) = - \sum_{i=1}^N x_i x_i^T \left(\frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)} \right) \left(\frac{1}{1 + \exp(\beta^T x_i)} \right)$$

minimizes a quadratic approximation to the function we are really interested in.

$$\theta_{k+1} = \theta_k - \mathbf{H}_K^{-1} \mathbf{g}_k$$

$p(x_i; \beta)$

$1 - p(x_i; \beta)$

Newton-Raphson for LR...

$$\frac{\partial l(\beta)}{\partial \beta} = \sum_{i=1}^N (y_i - \frac{\exp(\beta^T x)}{1 + \exp(\beta^T x)}) x_i = X^T (y - p)$$

$\rightarrow p(y=1|x) = \frac{e^{\beta^T x}}{1 + e^{\beta^T x}}$

$$(\frac{\partial^2 l(\beta)}{\partial \beta \partial \beta^T}) = -X^T W X$$

So, NR rule becomes:

$$\beta^{new} \leftarrow \beta^{old} + (X^T W X)^{-1} X^T (y - p),$$

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix}_{N \times (p+1)}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}_{N \times 1}, \quad p = \begin{bmatrix} \exp(\beta^T x_1) / (1 + \exp(\beta^T x_1)) \\ \exp(\beta^T x_2) / (1 + \exp(\beta^T x_2)) \\ \vdots \\ \exp(\beta^T x_N) / (1 + \exp(\beta^T x_N)) \end{bmatrix}_{N \times 1}$$

X : $N \times (p+1)$ matrix of x_i

y : $N \times 1$ matrix of y_i

p : $N \times 1$ matrix of $p(x_i; \beta^{old})$

W : $N \times N$ diagonal matrix of $p(x_i; \beta^{old})(1 - p(x_i; \beta^{old}))$

$$\left(\frac{\exp(\beta^T x_i)}{1 + \exp(\beta^T x_i)} \right) \left(1 - \frac{1}{1 + \exp(\beta^T x_i)} \right)$$

Newton-Raphson for LR...

- Newton-Raphson

$$- \beta^{new} = \beta^{old} + (X^T W X)^{-1} X^T (y - p)$$

$$= (X^T W X)^{-1} X^T W (X \beta^{old} + W^{-1} (y - p))$$

$$= (X^T W X)^{-1} X^T W z$$

Re expressing Newton step as weighted least square step

- Adjusted response

$$z = X \beta^{old} + W^{-1} (y - p)$$

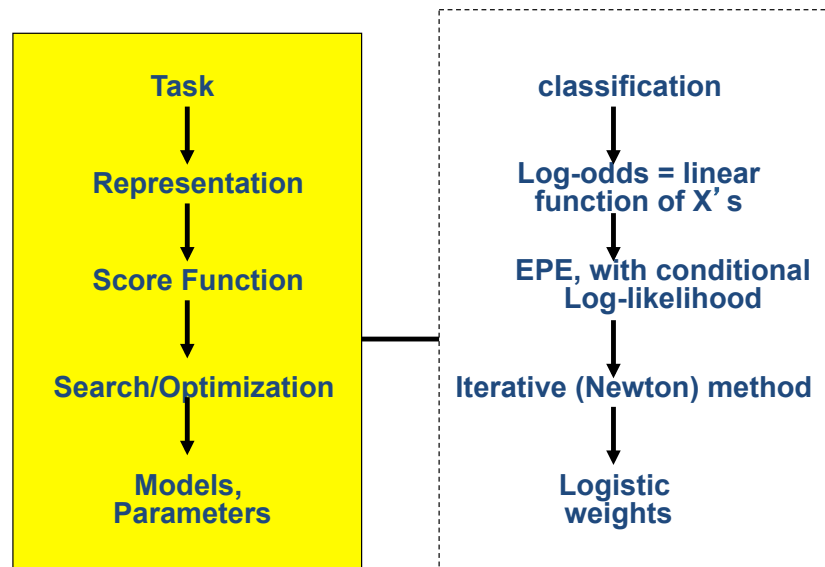
$(X^T W X)^{-1} X^T W z$

- Iteratively reweighted least squares (IRLS)

$$\beta^{new} \leftarrow \arg \min_{\beta} (z - X \beta^T)^T W (z - X \beta^T)$$

$$\leftarrow \arg \min_{\beta} (y - p)^T W^{-1} (y - p)$$

Logistic Regression



$$P(c = 1|x) = \frac{e^{\alpha+\beta x}}{1 + e^{\alpha+\beta x}}$$

10/21/15

29

Today :

- ✓ Logistic regression
- ✓ Generative vs. Discriminative

10/21/15

30

Discriminative vs. Generative

Generative approach

- Model the **joint distribution** $p(X, C)$ using $p(X | C = c_k)$ and $p(C = c_k)$

← Class prior

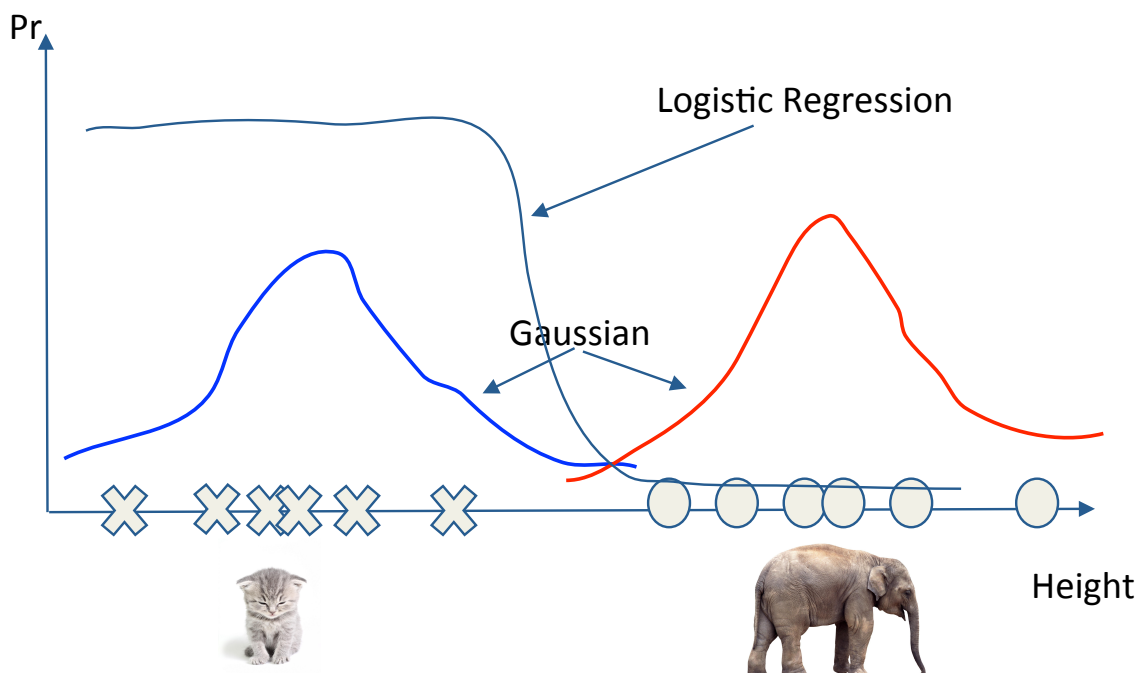
Discriminative approach

- Model the **conditional distribution** $p(c | X)$ directly

e.g.,

$$\frac{1}{1 + e^{-(\beta_0 + \beta_1 * X)}}$$

Discriminative vs. Generative



LDA vs. Logistic Regression

• LDA (Generative model)

- Assumes Gaussian class-conditional densities and a common covariance
- Model parameters are estimated by maximizing the full log likelihood, parameters for each class are estimated independently of other classes, $K_p + \frac{p(p+1)}{2} + (K - 1)$ parameters
- Makes use of marginal density information $\Pr(x)$
- Easier to train, low variance, more efficient if model is correct
- Higher asymptotic error, but converges faster

• Logistic Regression (Discriminative model)

- Assumes class-conditional densities are members of the (same) exponential family distribution
- Model parameters are estimated by maximizing the conditional log likelihood, simultaneous consideration of all other classes, $(K - 1)(p + 1)$ parameters
- Ignores marginal density information $\Pr(x)$
- Harder to train, robust to uncertainty about the data generation process
- Lower asymptotic error, but converges more slowly

10/21/15

33

Discriminative vs. Generative

● Definitions

- h_{gen} and h_{dis} : generative and discriminative classifiers
- $h_{\text{gen, inf}}$ and $h_{\text{dis, inf}}$: same classifiers but trained on the entire population (asymptotic classifiers)
- $n \rightarrow \text{infinity}$, $h_{\text{gen}} \rightarrow h_{\text{gen, inf}}$ and $h_{\text{dis}} \rightarrow h_{\text{dis, inf}}$

Ng, Jordan,. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.

Discriminative vs. Generative

Proposition 1:

$$\epsilon(h_{dis,inf}) \leq \epsilon(h_{gen,inf})$$

Proposition 1 states that asymptotically, the error of the discriminative logistic regression is smaller than that of the generative naive Bayes. This is easily shown

- p : number of dimensions
- n : number of observations
- ϵ : generalization error

Logistic Regression vs. NBC

Discriminative classifier (Logistic Regression)

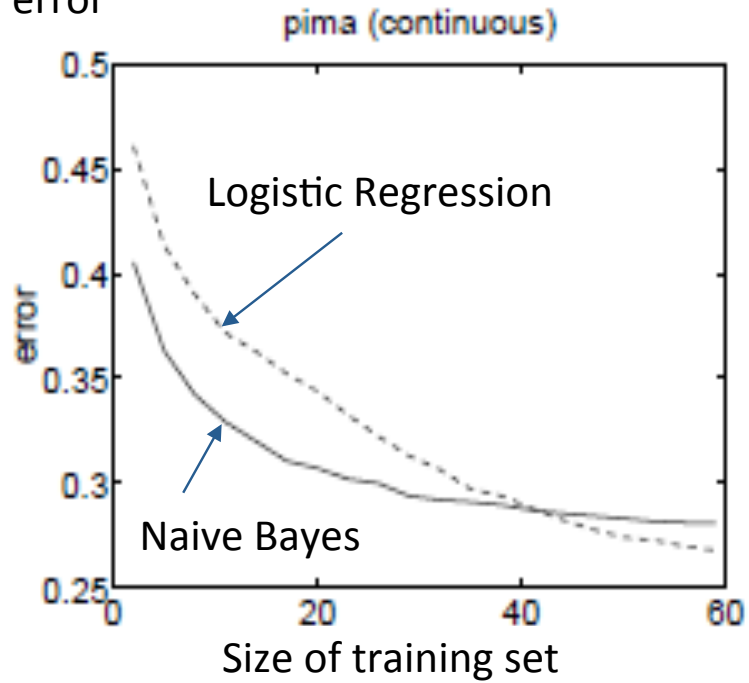
- Smaller asymptotic error
- Slow convergence $\sim O(p)$

Generative classifier (Naive Bayes)

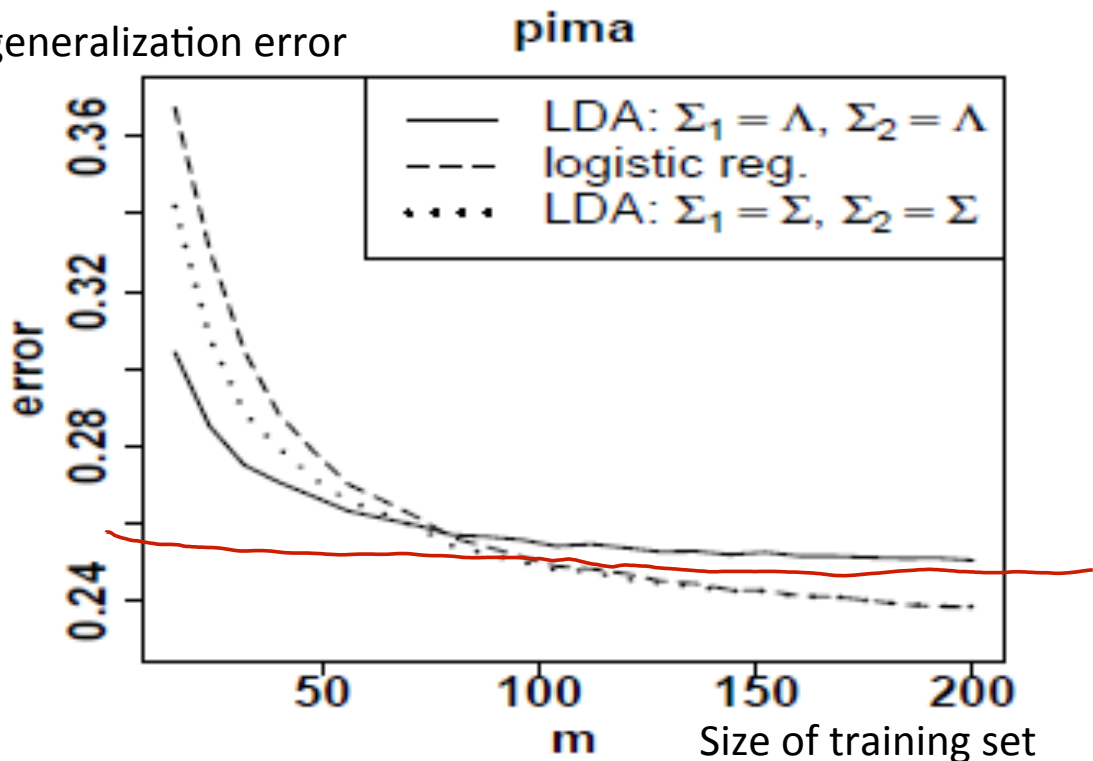
- Larger asymptotic error
- Can handle missing data (EM)
- Fast convergence $\sim O(\lg(p))$

Ng, Jordan,. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.

generalization error



generalization error



Xue, Jing-Hao, and D. Michael Titterton. "Comment on "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes". " *Neural processing letters* 28.3 (2008): 169-187.

Discriminative vs. Generative

- Empirically, **generative** classifiers approach their asymptotic error faster than discriminative ones
 - Good for small training set
 - Handle missing data well (EM)
- Empirically, **discriminative** classifiers have lower asymptotic error than generative ones
 - Good for larger training set

Yanjun Qi / UVA CS 4501-01-6501-07

References

- Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide
- Prof. Andrew Moore's slides
- Prof. Eric Xing's slides
- Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.