

# CS145: INTRODUCTION TO DATA MINING

## 4: Vector Data: Logistic Regression

---

**Instructor: Yizhou Sun**

[yzsun@cs.ucla.edu](mailto:yzsun@cs.ucla.edu)


October 6, 2021

# Methods to Learn

	Vector Data	Set Data	Sequence Data/Time Series	Text Data	Graph Data
Classification	<b>Logistic Regression;</b> Decision Tree; NN			Naïve Bayes for Text	Label Propagation
Clustering	K-means; Mixture Models			PLSA	Spectral Clustering
Prediction	<b>Linear Regression</b> GLM*		AR Model		
Frequent Pattern Mining		Apriori; FP growth	GSP; PrefixSpan		
Similarity Search			DTW		P-PageRank
Ranking					PageRank

# Vector Data: Logistic Regression

---

- Classification: Basic Concepts 
- Logistic Regression Model
- Generalized Linear Model\*
- Summary

# Supervised vs. Unsupervised Learning

---

- **Supervised learning (classification)**
  - **Supervision:** The training data (observations, measurements, etc.) are accompanied by **labels** indicating the class of the observations
  - New data is classified based on the training set
- **Unsupervised learning (clustering)**
  - The class labels of training data is unknown
  - Given a set of measurements, observations, etc. with the aim of establishing the existence of classes or clusters in the data

# Prediction Problems: Classification vs. Numeric Prediction

---

- **Classification**
  - predicts categorical class labels
  - classifies data (constructs a model) based on the training set and the values (**class labels**) in a classifying attribute and uses it in classifying new data
- **Numeric Prediction**
  - models continuous-valued functions, i.e., predicts unknown or missing values
- **Typical applications**
  - **Medical diagnosis:** if a tumor is cancerous or benign
  - **Fraud detection:** if a transaction is fraudulent
  - **Web page categorization:** which category it is

# Classification—A Two-Step Process (1)

---

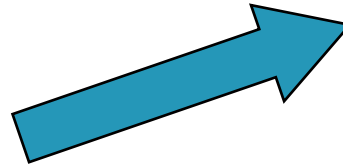
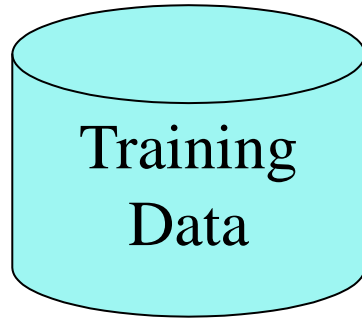
- **Model construction**: describing a set of predetermined classes
  - Each tuple/sample is assumed to belong to a predefined class, as determined by the **class label attribute**
    - For data point  $i$ :  $\langle \mathbf{x}_i, y_i \rangle$
    - Features:  $\mathbf{x}_i$ ; class label:  $y_i$
  - The model is represented as classification rules, decision trees, or mathematical formulae
    - Also called classifier
- The set of tuples used for model construction is **training set**

# Classification—A Two-Step Process (2)

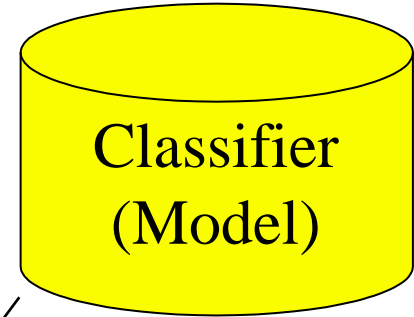
---

- **Model usage**: for classifying future or unknown objects
- **Estimate accuracy of the model**
  - The known label of test sample is compared with the classified result from the model
  - **Test set** is independent of training set (otherwise overfitting)
  - **Accuracy** rate is the percentage of test set samples that are correctly classified by the model
    - Most used for binary classes
- **If the accuracy is acceptable, use the model to classify new data**
- **Note**: If *the test set* is used to select models, it is called **validation (test) set**

# Process (1): Model Construction



Classification Algorithms

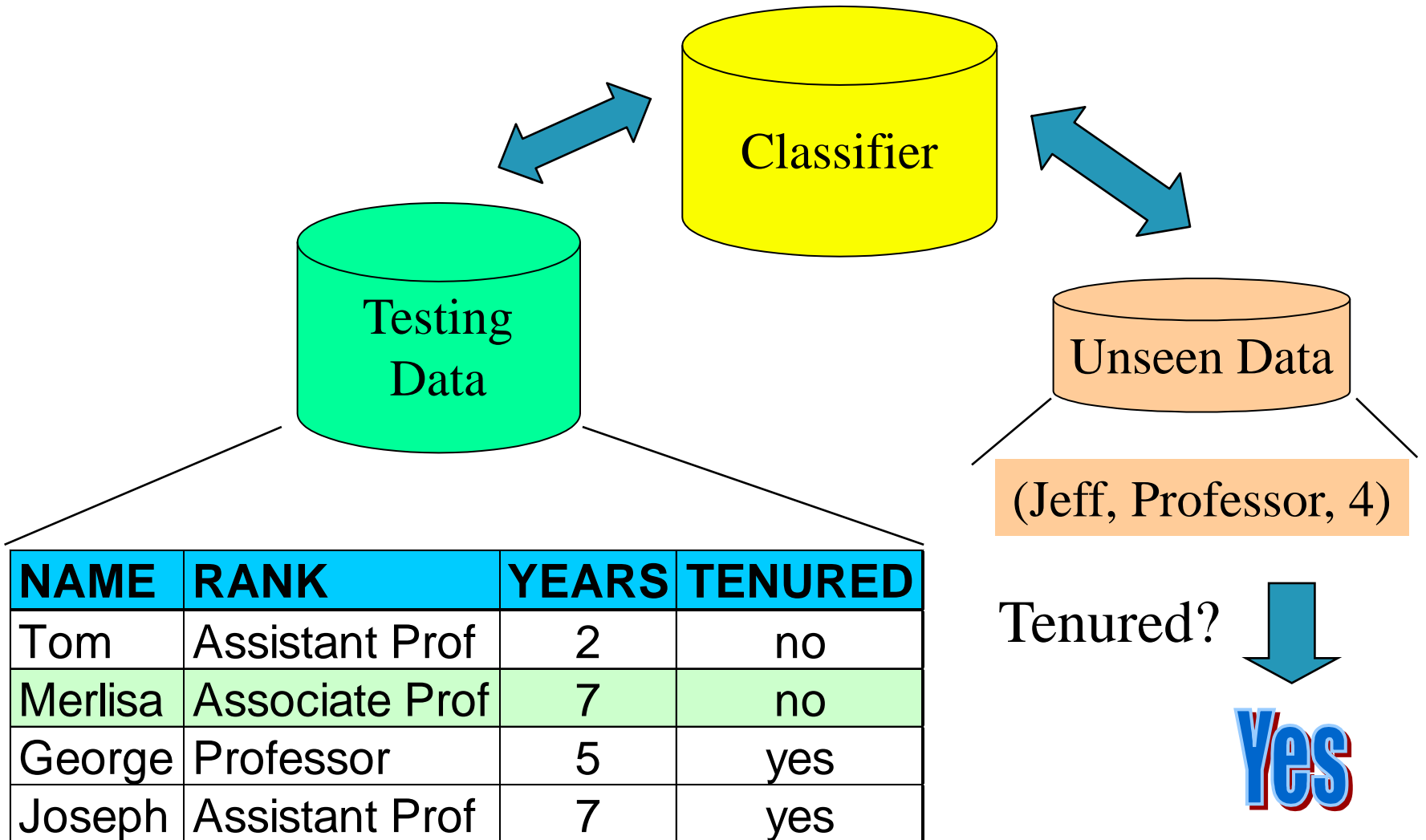


NAME	RANK	YEARS	TENURED
Mike	Assistant Prof	3	no
Mary	Assistant Prof	7	yes
Bill	Professor	2	yes
Jim	Associate Prof	7	yes
Dave	Assistant Prof	6	no
Anne	Associate Prof	3	no

IF rank = 'professor'  
OR years > 6  
THEN tenured = 'yes'




# Process (2): Using the Model in Prediction

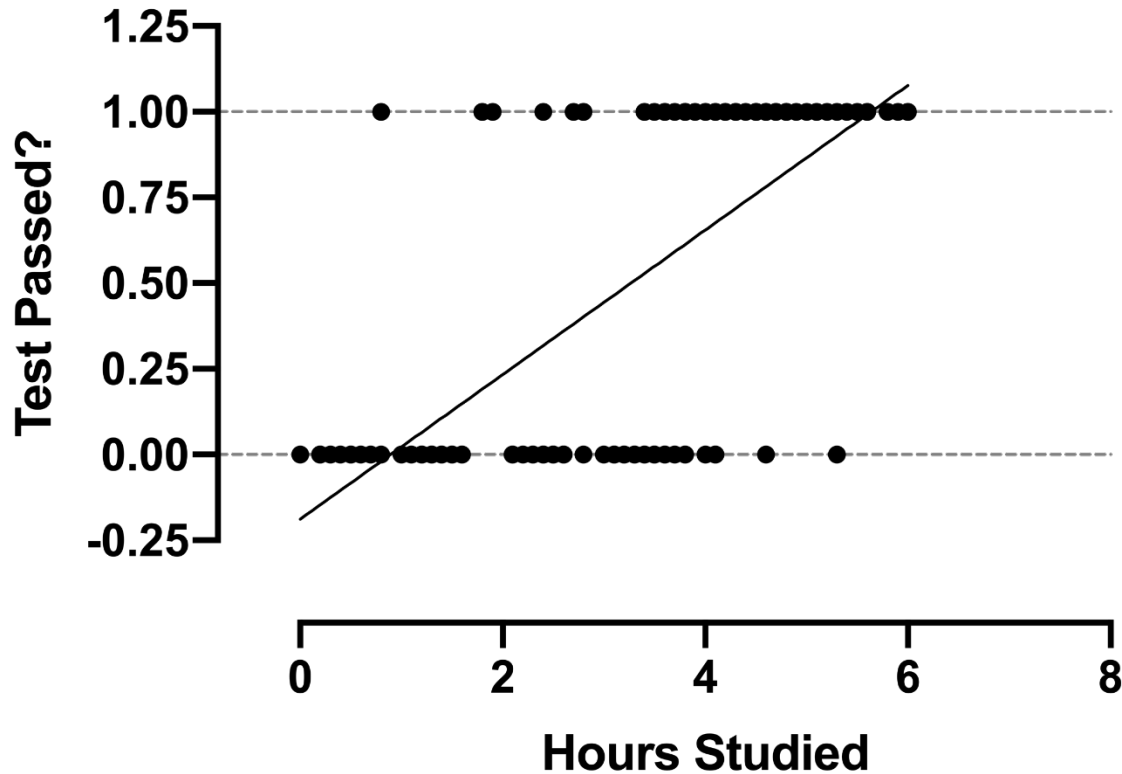


# Vector Data: Logistic Regression

---

- Classification: Basic Concepts
- Logistic Regression Model 
- Generalized Linear Model\*
- Summary

# A toy example



- Can it be solved by linear regression?

[https://www.graphpad.com/guides/prism/latest/curve-fitting/reg\\_simple\\_logistic\\_and\\_linear\\_difference.htm](https://www.graphpad.com/guides/prism/latest/curve-fitting/reg_simple_logistic_and_linear_difference.htm)

# Linear Regression VS. Logistic Regression

---

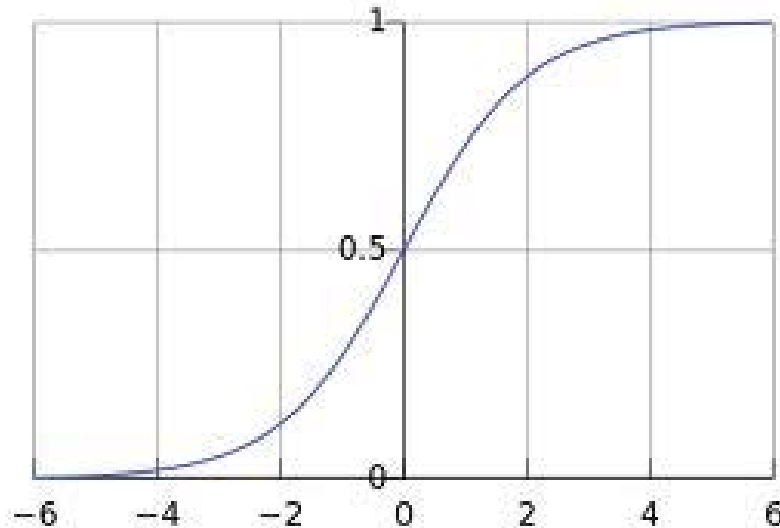
- Linear Regression (prediction)
  - $Y$ : *continuous value*  $(-\infty, +\infty)$ 
    - $Y = \mathbf{x}^T \boldsymbol{\beta} = \beta_0 + x_1\beta_1 + x_2\beta_2 + \dots + x_p\beta_p$
    - $Y|\mathbf{x}, \boldsymbol{\beta} \sim N(\mathbf{x}^T \boldsymbol{\beta}, \sigma^2)$
- Logistic Regression (classification)
  - $Y$ : *discrete value from  $m$  classes*
    - $p(Y = C_j|\mathbf{x}, \boldsymbol{\beta}) \in [0,1]$  and  $\sum_j p(Y = C_j|\mathbf{x}, \boldsymbol{\beta}) = 1$

# Logistic Function

---

- Logistic Function / sigmoid function:

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



*Note:*  $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

# Modeling Probabilities of Two Classes

---

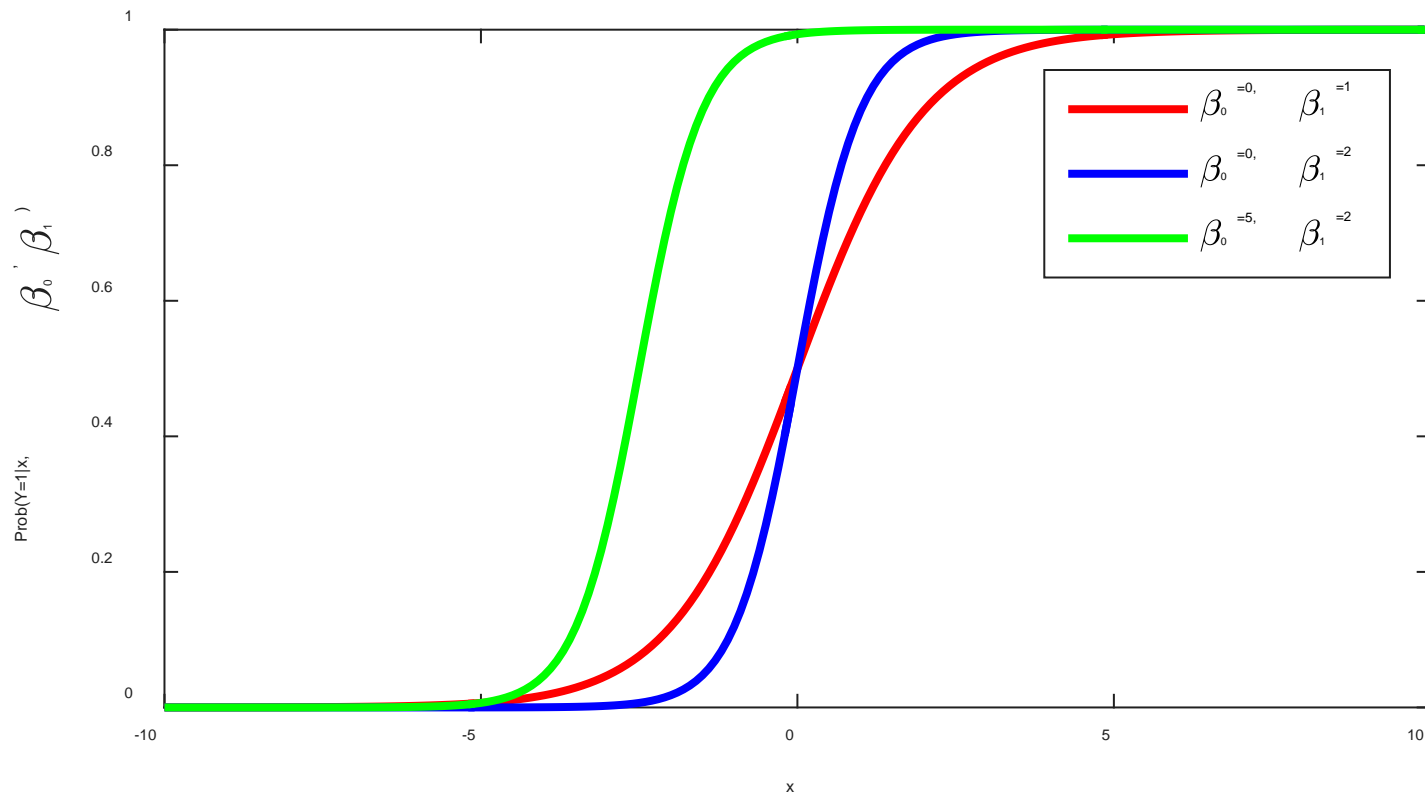
- $P(Y = 1|\mathbf{x}, \beta) = \sigma(\mathbf{x}^T \beta) = \frac{1}{1+\exp\{-\mathbf{x}^T \beta\}} = \frac{\exp\{\mathbf{x}^T \beta\}}{1+\exp\{\mathbf{x}^T \beta\}}$
- $P(Y = 0|\mathbf{x}, \beta) = 1 - \sigma(\mathbf{x}^T \beta) = \frac{\exp\{-\mathbf{x}^T \beta\}}{1+\exp\{-\mathbf{x}^T \beta\}} = \frac{1}{1+\exp\{\mathbf{x}^T \beta\}}$

$$\beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

- In other words
  - $y|\mathbf{x}, \beta \sim \text{Bernoulli}(\sigma(\mathbf{x}^T \beta))$

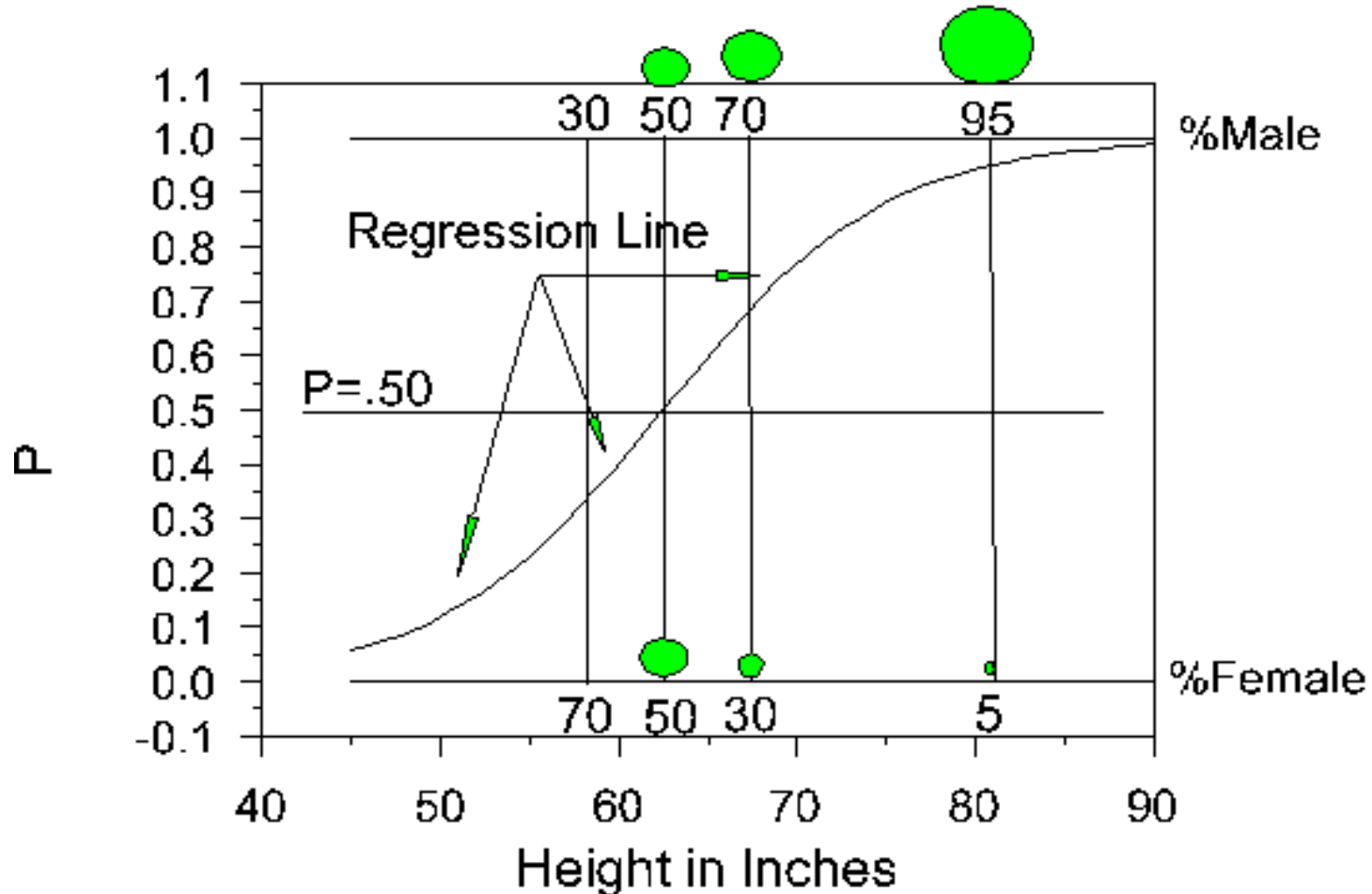
# The 1-d Situation

- $P(Y = 1|x, \beta_0, \beta_1) = \sigma(\beta_1 x + \beta_0)$



# Example

## Regression of Sex on Height



Q: What do we know about  $\beta_0$  here? Positive or negative?



# Parameter Estimation

---

- MLE estimation
  - Given a dataset  $D$ , with  $N$  data points
  - For a single data object with attributes  $\mathbf{x}_i$ , class label  $y_i$ 
    - Let  $p_i = p(y_i = 1 | \mathbf{x}_i, \beta)$ , the prob. of  $i$  in class 1
    - The probability of observing  $y_i$  would be
      - If  $y_i = 1$ , then  $p_i$
      - If  $y_i = 0$ , then  $1 - p_i$
      - Combing the two cases:  $p_i^{y_i}(1 - p_i)^{1-y_i}$

$$L = \prod_i p_i^{y_i} (1 - p_i)^{1-y_i} = \prod_i \left( \frac{\exp\{\mathbf{x}^T \beta\}}{1 + \exp\{\mathbf{x}^T \beta\}} \right)^{y_i} \left( \frac{1}{1 + \exp\{\mathbf{x}^T \beta\}} \right)^{1-y_i}$$

# Optimization

---

- Equivalent to maximize log likelihood
  - $\log L = \sum_i \{y_i \mathbf{x}_i^T \beta - \log(1 + \exp\{\mathbf{x}_i^T \beta\})\}$

- Gradient **ascent** update:

- $$\beta^{new} = \beta^{old} + \eta \frac{\partial \log L(\beta)}{\partial \beta}$$

- Newton-Raphson update

Step size

- $$\beta^{new} = \beta^{old} - \left( \frac{\partial^2 \log L(\beta)}{\partial \beta \partial \beta^T} \right)^{-1} \frac{\partial \log L(\beta)}{\partial \beta}$$

- where derivatives are evaluated at  $\beta^{old}$

# First Derivative

- It is a  $(p+1)$  vector, with  $j$ th element as

$$\begin{aligned} \bullet \frac{\partial \log L(\beta)}{\partial \beta_j} &= \sum_{i=1}^N y_i x_{ij} - \sum_{i=1}^N \frac{x_{ij} e^{x_i^T \beta}}{1 + e^{x_i^T \beta}} \\ &= \sum_{i=1}^N y_i x_{ij} - \sum_{i=1}^N p_i(\beta) x_{ij} \\ &= \sum_{i=1}^N x_{ij} (y_i - p_i(\beta)) \end{aligned}$$

$p_i(\beta) = \sigma(\beta^T x_i)$

For  $j = 0, 1, \dots, p$

# Second Derivative

- It is a  $(p+1)$  by  $(p+1)$  matrix, Hessian Matrix, with  $j$ th row and  $n$ th column as

$$\begin{aligned}\frac{\partial \log L(\boldsymbol{\beta})}{\partial \beta_j \partial \beta_n} &= - \sum_{i=1}^N \frac{(1 + e^{\boldsymbol{\beta}^T \mathbf{x}_i}) e^{\boldsymbol{\beta}^T \mathbf{x}_i} x_{ij} x_{in} - (e^{\boldsymbol{\beta}^T \mathbf{x}_i})^2 x_{ij} x_{in}}{(1 + e^{\boldsymbol{\beta}^T \mathbf{x}_i})^2} \\ &= - \sum_{i=1}^N x_{ij} x_{in} p_i(\boldsymbol{\beta}) - \sum_{i=1}^N x_{ij} x_{in} (p_i(\boldsymbol{\beta}))^2 \\ &= - \sum_{i=1}^N x_{ij} x_{in} p_i(\boldsymbol{\beta}) (1 - p_i(\boldsymbol{\beta}))\end{aligned}$$

**Matrix form:**

$$\begin{aligned}\frac{\partial \log L(\boldsymbol{\beta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= - \sum_{i=1}^N \mathbf{x}_i p_i(\boldsymbol{\beta}) (1 - p_i(\boldsymbol{\beta})) \mathbf{x}_i^T \\ &= - \mathbf{X}^T \begin{bmatrix} p_1(\boldsymbol{\beta})(1 - p_1(\boldsymbol{\beta})) & & \\ & \ddots & \\ & & p_N(\boldsymbol{\beta})(1 - p_N(\boldsymbol{\beta})) \end{bmatrix} \mathbf{X}\end{aligned}$$

where  $\mathbf{X}$  is the  $N \times (p + 1)$  feature matrix. Note  $\mathbf{X}^T = [x_1 \ x_2 \ \dots \ x_N]$ .

# An Alternative View of the Objective Function

---

- Cross entropy loss
  - Measure the difference from the predicted distribution ( $p$ ) to the ground truth distribution ( $q$ )
    - Cross entropy from  $q$  to  $p$ :  $H(q, p) = -\sum_k q_k \log(p_k)$
  - In the classification setting
    - $q_0 = 1$  and  $q_1 = 0$ , if  $y = 0$ ;  $q_0 = 0$  and  $q_1 = 1$ , if  $y = 1$
    - $p_0 = \frac{1}{1+\exp\{\mathbf{x}^T \boldsymbol{\beta}\}}$  and  $p_1 = \frac{\exp\{\mathbf{x}^T \boldsymbol{\beta}\}}{1+\exp\{\mathbf{x}^T \boldsymbol{\beta}\}}$

# An Alternative View of the Objective Function (Cont.)

---

- If  $y = 0$ 
    - $H(q, p) = \log(1 + \exp\{\mathbf{x}^T \beta\})$
  - If  $y = 1$ 
    - $H(q, p) = -\mathbf{x}^T \beta + \log(1 + \exp\{\mathbf{x}^T \beta\})$
  - Putting together
    - $H(q, p) = -y\mathbf{x}^T \beta + \log(1 + \exp\{\mathbf{x}^T \beta\})$
- Recall Log likelihood:  $\log L = \sum_i \{y_i \mathbf{x}_i^T \beta - \log(1 + \exp\{\mathbf{x}_i^T \beta\})\}$
- The goal: **minimize** the mean cross entropy loss over all the data points

# What about Multiclass Classification?

---

- It is easy to handle under logistic regression, say  $M$  classes, using softmax function
  - $P(Y = j|x) = \frac{\exp\{x^T \beta_j\}}{1 + \sum_{m=1}^{M-1} \exp\{x^T \beta_m\}}$ , for  $j = 1, \dots, M - 1$
  - $P(Y = M|x) = \frac{1}{1 + \sum_{m=1}^{M-1} \exp\{x^T \beta_m\}}$
- Loss function
  - Cross entropy loss from observed class distribution (e.g.,  $(0,0,1,0,0)$ ) to  $p$

slido




**Which of the following  
statement(s) are correct?**

ⓘ Start presenting to display the poll results on this slide.



# Vector Data: Logistic Regression

---

- Classification: Basic Concepts
- Logistic Regression Model
- Generalized Linear Model\* 
- Summary

# Recall Linear Regression and Logistic Regression

---

- Linear Regression
  - $y|\mathbf{x}, \beta \sim N(\mathbf{x}^T \beta, \sigma^2)$
- Logistic Regression
  - $y|\mathbf{x}, \beta \sim \text{Bernoulli}(\sigma(\mathbf{x}^T \beta))$
- How about other distributions?
  - Yes, as long as they belong to exponential family

# Exponential Family

---

- Canonical Form
  - $p(\mathbf{y}; \boldsymbol{\eta}) = b(\mathbf{y}) \exp(\boldsymbol{\eta}^T T(\mathbf{y}) - a(\boldsymbol{\eta}))$
  - $\boldsymbol{\eta}$ : natural parameter
  - $T(\mathbf{y})$ : sufficient statistic
  - $a(\boldsymbol{\eta})$ : log partition function for normalization
  - $b(\mathbf{y})$ : function that only dependent on  $\mathbf{y}$

# Examples of Exponential Family

- Many:

$$p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta))$$

- Gaussian, Bernoulli, Poisson, beta, Dirichlet, categorical, ...

- For Gaussian (not interested in  $\sigma$ )

$$\begin{aligned} p(y; \mu) &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(y - \mu)^2\right) \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}y^2\right) \cdot \exp\left(\mu y - \frac{1}{2}\mu^2\right) \end{aligned}$$

$$\begin{aligned} \eta &= \mu \\ T(y) &= y \\ a(\eta) &= \mu^2/2 \\ &= \eta^2/2 \\ b(y) &= (1/\sqrt{2\pi}) \exp(-y^2/2) \end{aligned}$$

- For Bernoulli

$$\begin{aligned} p(y; \phi) &= \phi^y (1 - \phi)^{1-y} \\ &= \exp(y \log \phi + (1 - y) \log(1 - \phi)) \\ &= \exp\left(\left(\log\left(\frac{\phi}{1 - \phi}\right)\right) y + \log(1 - \phi)\right) \end{aligned}$$

$\eta$

$$\begin{aligned} T(y) &= y \\ a(\eta) &= -\log(1 - \phi) \\ &= \log(1 + e^\eta) \\ b(y) &= 1 \end{aligned}$$


# Recipe of GLMs

---

- Determines a distribution for  $y$ 
  - E.g., Gaussian, Bernoulli, Poisson
- Form the linear predictor for  $\eta$ 
  - $\eta = \mathbf{x}^T \boldsymbol{\beta}$
- Determines a link function:  $\mu = g^{-1}(\eta)$ 
  - Connects the linear predictor to the mean of the distribution
  - E.g.,  $\mu = \eta$  for Gaussian,  $\mu = \sigma(\eta)$  for Bernoulli,  $\mu = \exp(\eta)$  for Poisson

# Vector Data: Logistic Regression

---

- Classification: Basic Concepts
- Logistic Regression Model
- Generalized Linear Model\*
- Summary 

# Summary

---

- What is classification
  - Supervised learning vs. unsupervised learning, classification vs. prediction
- Logistic regression
  - Sigmoid function, multiclass classification
- Generalized linear model\*
  - Exponential family, link function