

## Computing Probabilities from Data

- Various probabilities you will need to compute for Naive Bayesian Classifier (using MLE here):

$$\hat{P}(Y=0) = \frac{\text{\#instances in class=0}}{\text{total \#instances}}$$

$$\hat{P}(X_i=0, Y=0) = \frac{\text{\#instances where } X_i=0 \text{ and class}=0}{\text{total \#instances}}$$

$$\hat{P}(X_i=0|Y=0) = \frac{\hat{P}(X_i=0, Y=0)}{\hat{P}(Y=0)} \quad \hat{P}(X_i=0|Y=1) = \frac{\hat{P}(X_i=0, Y=1)}{\hat{P}(Y=1)}$$

$$\hat{P}(X_i=1|Y=0) = 1 - \hat{P}(X_i=0|Y=0)$$

$$\hat{y} = \arg \max_y P(\mathbf{X}|Y)P(Y) = \arg \max_y (\log P(\mathbf{X}|Y)P(Y))$$

$$\log P(\mathbf{X}|Y) = \log P(X_1, X_2, \dots, X_m | Y) = \log \prod_{j=1}^m P(X_j | Y) = \sum_{j=1}^m \log P(X_j | Y)$$

## From Naive Bayes to Logistic Regression

- Recall the Naive Bayes Classifier
  - Predict  $\hat{y} = \arg \max_y P(\mathbf{X}, Y) = \arg \max_y P(\mathbf{X}|Y)P(Y)$
  - Use assumption that  $P(\mathbf{X}|Y) = P(X_1, X_2, \dots, X_m | Y) = \prod_{i=1}^m P(X_i | Y)$
  - We are really modeling joint probability  $P(\mathbf{X}, Y)$
- But for classification, really care about  $P(Y | \mathbf{X})$ 
  - Really want to predict  $\hat{y} = \arg \max_y P(Y | \mathbf{X})$
  - Modeling full joint probability  $P(\mathbf{X}, Y)$  is just proxy for this
- So, how do we model  $P(Y | \mathbf{X})$  directly?
  - Welcome our friend: logistic regression!

## Logistic Regression

- Model *conditional* likelihood  $P(Y | \mathbf{X})$  directly

- Model this probability with *logistic* function:

$$P(Y=1 | \mathbf{X}) = \frac{1}{1 + e^{-z}} \quad \text{where } z = \alpha + \sum_{j=1}^m \beta_j X_j$$

- For simplicity define  $X_0 = 1$  and  $\beta_0 = \alpha$ , so  $z = \sum_{j=0}^m \beta_j X_j$

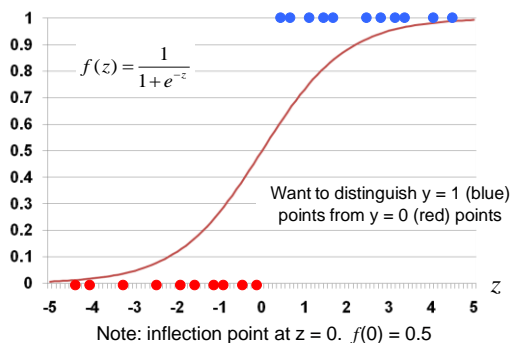
- Since  $P(Y=0 | \mathbf{X}) + P(Y=1 | \mathbf{X}) = 1$ , we obtain:

$$P(Y=0 | \mathbf{X}) = \frac{e^{-z}}{1 + e^{-z}} \quad \text{where } z = \sum_{j=0}^m \beta_j X_j$$

- Note: log-odds is **linear** function of inputs  $X_j$ :

$$\log \frac{P(Y=1 | \mathbf{X})}{P(Y=0 | \mathbf{X})} = \log \frac{1}{e^{-z}} = \log e^z = z = \sum_{j=0}^m \beta_j X_j$$

## The Logistic Function



## Learning Logistic Regression Function

- Can write log-conditional likelihood of data as:

$$LL(\theta) = \sum_{i=1}^n [y_i \log P(Y=1 | \mathbf{X}_i) + (1 - y_i) \log P(Y=0 | \mathbf{X}_i)]$$

where  $P(Y=1 | \mathbf{X}) = \frac{1}{1 + e^{-z}}$ ,  $P(Y=0 | \mathbf{X}) = \frac{e^{-z}}{1 + e^{-z}}$  with  $z = \sum_{j=0}^m \beta_j X_j$

- No analytic derivation of MLE parameters for  $\beta_j$

- But, log-conditional likelihood function is *concave*
- Has a single global maximum (good times)

- Compute gradient of  $LL(\theta)$  w.r.t.  $\beta_j$  where  $0 \leq j \leq m$ :

$$\frac{\partial LL(\theta)}{\partial \beta_j} = x_j (y - \frac{1}{1 + e^{-z}}) = x_j (y - P(Y=1 | \mathbf{X}))$$

- Maximize  $LL(\theta)$ : iteratively update  $\beta_j$  using gradient:

$$\beta_j^{\text{new}} = \beta_j^{\text{old}} + c x_j (y - \frac{1}{1 + e^{-z}}) \quad \text{where } z = \sum_{j=0}^m \beta_j^{\text{old}} x_j \text{ and } c = \text{constant}$$

## Wanna See How We Computed Gradient?

- Log-conditional likelihood of data point  $i$ ,  $LL_i(\theta)$ :

$$LL_i(\theta) = y_i \log P(Y=1 | \mathbf{X}_i) + (1 - y_i) \log P(Y=0 | \mathbf{X}_i)$$

- Rearrange terms:

$$LL_i(\theta) = y_i \log \frac{P(Y=1 | \mathbf{X}_i)}{P(Y=0 | \mathbf{X}_i)} + \log P(Y=0 | \mathbf{X}_i)$$

- Substitute values for  $P(Y | \mathbf{X})$  and simplify:

$$\begin{aligned} LL_i(\theta) &= y_i \sum_{j=0}^m \beta_j X_{ij} + \log e^{-z} - \log(1 + e^{-z}) \\ &= y_i \sum_{j=0}^m \beta_j X_{ij} - \sum_{j=0}^m \beta_j X_{ij} - \log(1 + e^{-z}) \end{aligned}$$

- Compute gradient of  $LL(\theta)$  w.r.t.  $\beta_j$  where  $0 \leq j \leq m$ :

$$\frac{\partial LL_i(\theta)}{\partial \beta_j} = y x_{ij} - x_{ij} + \frac{x_j e^{-z}}{1 + e^{-z}} = x_j (y - \frac{1}{1 + e^{-z}} + \frac{e^{-z}}{1 + e^{-z}}) = x_j (y - \frac{1}{1 + e^{-z}})$$

## “Batch” Logistic Regression Algorithm

```

Initialize:  $\beta_j = 0$  for all  $0 \leq j \leq m$ 
// "epochs" = number of passes over data during learning
for (i = 0; i < epochs; i++) {
  Initialize: gradient[j] = 0 for all  $0 \leq j \leq m$ 
  // Compute "batch" gradient vector
  for each training instance ( $\langle x_1, x_2, \dots, x_m \rangle, y$ ) in data {
    // Add contribution to gradient for each data point
    for (j = 0; j <= m; j++) {
      // Note:  $x_j$  below is j-th input variable and  $x_0 = 1$ .
      gradient[j] +=  $x_j(y - \frac{1}{1+e^{-z}})$  where  $z = \sum_{j=0}^m \beta_j x_j$ 
    }
  }
  // Update all  $\beta_j$ . Note learning rate  $\eta$  is pre-set constant
   $\beta_j += \eta * \text{gradient}[j]$  for all  $0 \leq j \leq m$ 
}

```

## Classification with Logistic Regression

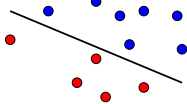
- Training: determine parameters  $\beta_j$  (for all  $0 \leq j \leq m$ )
  - After parameters  $\beta_j$  have been learned, test classifier
- To test classifier, for each new (test) instance  $\mathbf{X}$ :
  - Compute:  $p = P(Y=1|\mathbf{X}) = \frac{1}{1+e^{-z}}$ , where  $z = \sum_{j=0}^m \beta_j X_j$
  - Classify instance as:  $\hat{y} = \begin{cases} 1 & p > 0.5 \\ 0 & \text{otherwise} \end{cases}$
  - Note about evaluation set-up: parameters  $\beta_j$  are **not** updated during “testing” phase
  - In real systems, model parameters are updated as new data becomes available (on a periodic basis)

## Linear Separability

- Recall that log-odds is **linear** function of inputs  $X_j$ :

$$\log \frac{P(Y=1|\mathbf{X})}{P(Y=0|\mathbf{X})} = \sum_{j=0}^m \beta_j X_j$$

- Logistic regression is trying to fit a **line** that separates data instances where  $y=1$  from those where  $y=0$



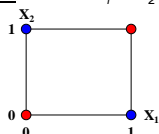
- We call such data (or the functions generating the data) “**linearly separable**”

## Wouldn’t You Like to Be Linear Too?

- Logistic Regression as a linear function
  - As mentioned  $\log \frac{P(Y=1|\mathbf{X})}{P(Y=0|\mathbf{X})} = \sum_{j=0}^m \beta_j X_j$  is linear in inputs  $X_j$
  - That is, each  $X_j$  is multiplied by separate  $\beta_j$
  - No direct interaction (multiplication) of multiple  $X_j$
- Such linearity is also true for Naïve Bayes
  - You will show it on the next problem set
    - It’s pretty straight-forward (pay attention to the top of this slide)
  - So, Logistic Regression and Naïve Bayes have same functional form!
    - Each is just maximizing a different objective function:
      - Naïve Bayes:  $P(\mathbf{X}, Y)$
      - Logistic Regression:  $P(Y|\mathbf{X})$

## Data Often Not Linearly Separable

- Many data sets/functions are not linearly separable
  - Consider function:  $y = x_1 \text{ XOR } x_2$
  - Note:  $y=1$  iff **one** of either  $x_1$  or  $x_2 = 1$



- Not possible to draw a line that successfully separates all the  $y=1$  points (blue) from the  $y=0$  points (red)
- Despite this fact, logistic regression and Naïve Bayes still often work well in practice

## Logistic Regression vs. Naïve Bayes

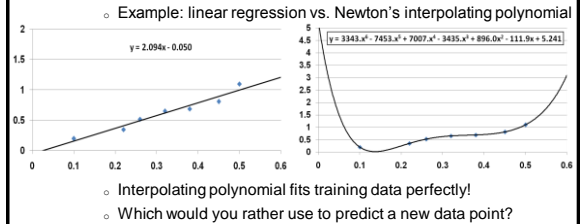
- Compare Naïve Bayes and Logistic Regression
  - Recall that Naïve Bayes models  $P(\mathbf{X}, Y) = P(\mathbf{X}|Y)P(Y)$
  - Logistic Regression directly models  $P(Y|\mathbf{X})$
  - We call Naïve Bayes a “generative model”
    - Tries to model joint distribution of how data is “generated”
    - I.e., could use  $P(\mathbf{X}, Y)$  to generate new data points if we wanted
    - But lots of effort to model something that may not be needed
  - We call Logistic Regression a “discriminative model”
    - Just tries to model way to discriminate  $y=0$  vs.  $y=1$  cases
    - Cannot use model to generate new data points (no  $P(\mathbf{X}, Y)$ )
    - Note: Logistic Regression can be generalized to more than two output values for  $y$  (have multiple sets of parameters  $\beta_j$ )

## Choosing an Algorithm

- Many trade-offs in choosing learning algorithm
  - Continuous input variables
    - Logistic Regression easily deals with continuous inputs
    - Naive Bayes needs to use some parametric form for continuous inputs (e.g., Gaussian) or "discretize" continuous values into ranges (e.g., temperature in range: <50, 50-60, 60-70, >70)
  - Discrete input variables
    - Naive Bayes naturally handles multi-valued discrete data by using multinomial distribution for  $P(X_i | Y)$
    - Logistic Regression requires some sort of representation of multi-valued discrete data (e.g., multiple binary features)
    - Say  $X_i \in \{A, B, C\}$ . Not necessarily a good idea to encode  $X_i$  as taking on input values 1, 2, or 3 corresponding to A, B, or C.

## Good Machine Learning = Generalization

- Goal of machine learning: build models that **generalize** well to predicting new data
  - "Overfitting": fitting the training data too well, so we lose generality of model

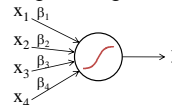


## Issues with Logistic Regression

- Logistic Regression can more easily overfit training data than Naive Bayes
  - Logistic Regression is not modeling whole distribution, it is just optimizing prediction of Y
  - Overfitting can be especially problematic if distributions of training data and testing data differ a bit
  - There are methods to mitigate overfitting in Logistic Regression
    - Use Bayesian priors on parameters  $\beta$ , rather than just maximizing conditional likelihood
    - Intuitively, analogous to using priors in Naive Bayes to get MAP estimates of probabilities
    - But optimization process is more complex in Logistic Regression since conditional likelihood has no analytic solution

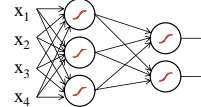
## Logistic Regression and Neural Networks

- Consider logistic regression as:



Logistic regression is same as a one node neural network

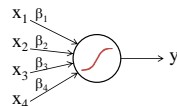
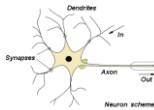
- Neural network



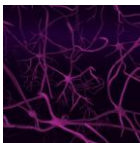
- Neural network uses "back-propagation" algorithm
  - Like Logistic Regression optimization on "steroids"
  - With enough nodes, can approximate *any* function

## Biological Basis for Neural Networks

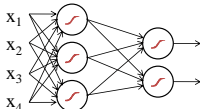
- A neuron



- Your brain



Actually, it's probably someone else's brain



## Now You Too Can Build Terminators!

- Be careful! ☺



"My CPU is a neural net processor, a learning computer. The more contact I have with humans, the more I learn."