

CS 2750 Machine Learning
Lecture 10

Generative classification models

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

Classification

- **Data:** $D = \{d_1, d_2, \dots, d_n\}$
 $d_i = \langle \mathbf{x}_i, y_i \rangle$
 - y_i **represents a discrete class value**
 - **Goal: learn** $f : X \rightarrow Y$
 - **Binary classification**
 - **A special case when** $Y \in \{0,1\}$
 - **First step:**
 - we need to devise a model of the function f
-

Discriminant functions

- A common way to represent a classifier is by using
 - **Discriminant functions**
- Works for both the binary and multi-way classification
- **Idea:**
 - For every class $i = 0, 1, \dots, k$ define a function $g_i(\mathbf{x})$ mapping $X \rightarrow \mathfrak{R}$
 - When the decision on input \mathbf{x} should be made choose the class with the highest value of $g_i(\mathbf{x})$

$$y^* = \arg \max_i g_i(\mathbf{x})$$

Logistic regression model

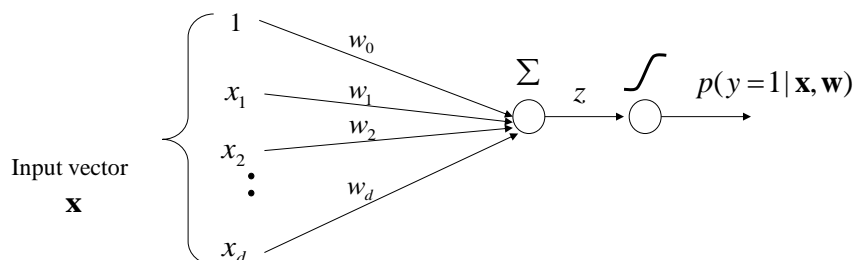
- **Discriminant functions:**

$$g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x}) \quad g_0(\mathbf{x}) = 1 - g(\mathbf{w}^T \mathbf{x})$$

- Values of discriminant functions vary in interval $[0,1]$

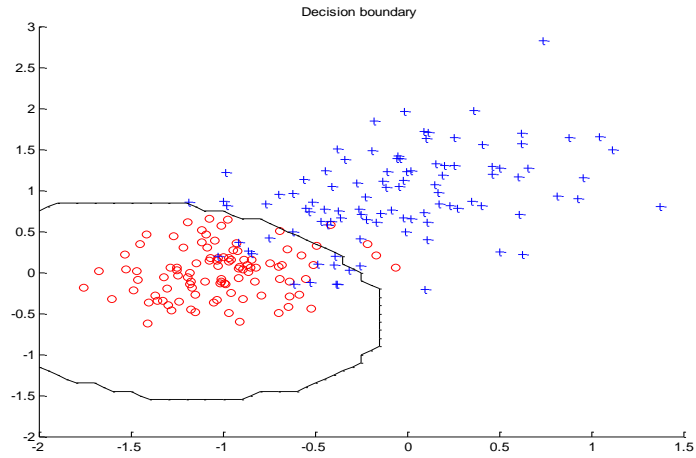
- **Probabilistic interpretation**

$$f(\mathbf{x}, \mathbf{w}) = p(y = 1 | \mathbf{w}, \mathbf{x}) = g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x})$$



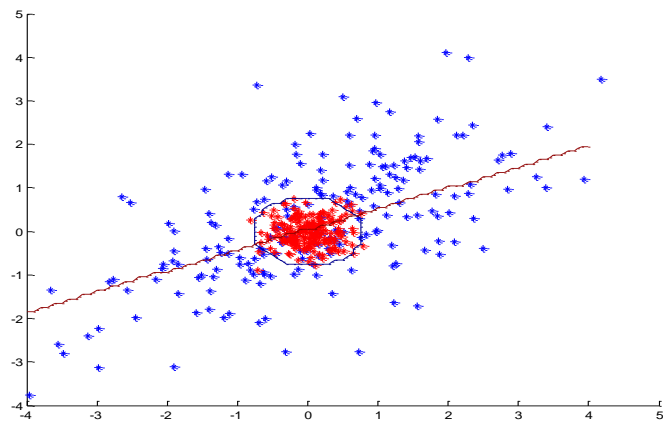
When does the logistic regression fail?

- Nonlinear decision boundary



When does the logistic regression fail?

- Another example of a non-linear decision boundary



Non-linear extension of logistic regression

- use **feature (basis) functions** to model **nonlinearities**
 - the same trick as used for the linear regression

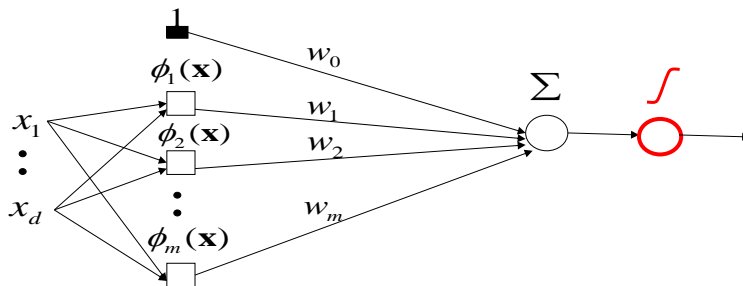
Linear regression

$$f(\mathbf{x}) = w_0 + \sum_{j=1}^m w_j \phi_j(\mathbf{x})$$

Logistic regression

$$p(y = 1 | x) = g(w_0 + \sum_{j=1}^m w_j \phi_j(\mathbf{x}))$$

$\phi_j(\mathbf{x})$ - an arbitrary function of \mathbf{x}

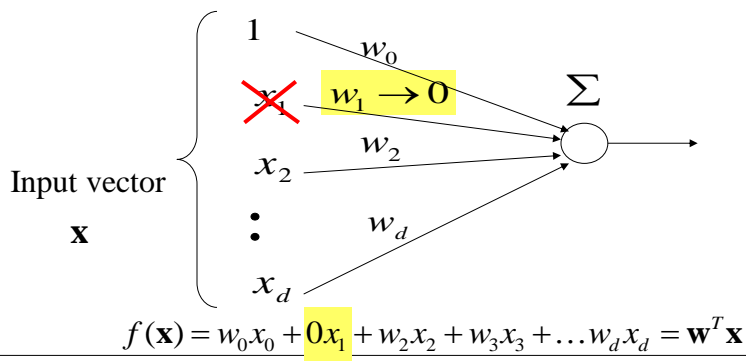


CS 2750 Machine Learning

Regularized logistic regression

- If the model is too complex and can cause overfitting, its prediction accuracy can be improved by **removing some inputs from the model = setting their coefficients to zero**
- **Recall the linear model:**

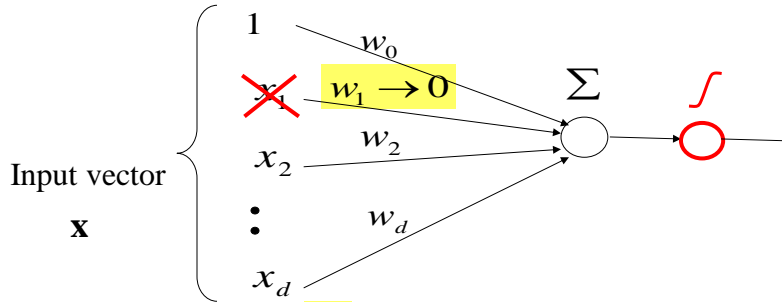
$$f(\mathbf{x}) = w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots w_d x_d = \mathbf{w}^T \mathbf{x}$$



Regularized logistic regression

- If the model is too complex and can cause overfitting, its prediction accuracy can be improved by **removing some inputs from the model = setting their coefficients to zero**
- **We can apply the same idea to the logistic regression:**

$$p(y=1 | \mathbf{x}) = g(\mathbf{w}^T \mathbf{x}) \quad w_0, w_1, \dots, w_k \text{ - parameters (weights)}$$



$$p(y=1 | \mathbf{x}) = g(w_0 x_0 + 0 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_d x_d) = g(\mathbf{w}^T \mathbf{x})$$

Ridge (L2) penalty

Linear regression – Ridge penalty:

$$J_n(\mathbf{w}) = \underbrace{\frac{1}{n} \sum_{i=1, \dots, n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2}_{\text{Fit to data}} + \underbrace{\lambda \|\mathbf{w}\|_{L2}^2}_{\text{Model complexity penalty}}$$

$$\|\mathbf{w}\|_{L2}^2 = \sum_{i=0}^d w_i^2 = \mathbf{w}^T \mathbf{w} \quad \text{and} \quad \lambda \geq 0$$

Logistic regression:

$$J_n(\mathbf{w}) = \underbrace{-\log P(D | \mathbf{w})}_{\text{Fit to data}} + \underbrace{\lambda \|\mathbf{w}\|_{L2}^2}_{\text{Model complexity penalty}}$$

$$J_n(\mathbf{w}) = - \left[\sum_{i=1}^n y_i \log g(\mathbf{w}^T x_i) + (1 - y_i) \log(1 - g(\mathbf{w}^T x_i)) \right] + \lambda \|\mathbf{w}\|_{L2}^2$$

Fit to data measured using the negative log likelihood

Lasso (L1) penalty

Linear regression – Lasso penalty:

$$J_n(\mathbf{w}) = \underbrace{\frac{1}{n} \sum_{i=1..n} (y_i - \mathbf{w}^T \mathbf{x}_i)^2}_{\text{Fit to data}} + \underbrace{\lambda \|\mathbf{w}\|_{L1}}_{\text{Model complexity penalty}}$$

$$\|\mathbf{w}\|_{L1} = \sum_{i=0}^d |w_i| \quad \text{and} \quad \lambda \geq 0$$

Logistic regression:

$$J_n(\mathbf{w}) = \underbrace{-\log P(D | \mathbf{w})}_{\text{Fit to data}} + \underbrace{\lambda \|\mathbf{w}\|_{L1}}_{\text{Model complexity penalty}}$$

$$J_n(\mathbf{w}) = - \left[\sum_{i=1}^n y_i \log g(\mathbf{w}^T x_i) + (1 - y_i) \log(1 - g(\mathbf{w}^T x_i)) \right] + \lambda \|\mathbf{w}\|_{L1}$$

Fit to data measured using the negative log likelihood

Generative approach to classification

Logistic regression:

- Represents and learns a model of $p(y | \mathbf{x})$
- An example of a **discriminative classification approach**
- Model is unable to sample (generate) data instances (\mathbf{x}, y)

Generative approach:

- Represents and learns the joint distribution $p(\mathbf{x}, y)$
- Model is able to sample (generate) data instances (\mathbf{x}, y)
- The joint model defines probabilistic discriminant functions

How?

$$g_1(\mathbf{x}) = p(y = 1 | \mathbf{x}) = \frac{p(\mathbf{x}, y = 1)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | y = 1) p(y = 1)}{p(\mathbf{x})}$$

$$g_0(\mathbf{x}) = p(y = 0 | \mathbf{x}) = \frac{p(\mathbf{x}, y = 0)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | y = 0) p(y = 0)}{p(\mathbf{x})}$$

$$p(y = 0 | \mathbf{x}) + p(y = 1 | \mathbf{x}) = 1$$

Generative approach to classification

Typical joint model $p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$

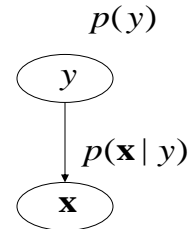
- $p(\mathbf{x} | y) =$ **Class-conditional distributions (densities)**

binary classification: two class-conditional distributions

$$p(\mathbf{x} | y = 0) \quad p(\mathbf{x} | y = 1)$$

- $p(y) =$ **Priors on classes**
 - probability of class y
 - for binary classification: Bernoulli distribution

$$p(y = 0) + p(y = 1) = 1$$



Quadratic discriminant analysis (QDA)

Model:

- **Class-conditional distributions are**
 - **multivariate normal distributions**

$$\mathbf{x} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad \text{for } y = 0$$

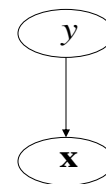
$$\mathbf{x} \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \quad \text{for } y = 1$$

Multivariate normal $\mathbf{x} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- **Priors on classes (class 0,1)** $y \sim$ *Bernoulli*
 - **Bernoulli distribution**

$$p(y, \theta) = \theta^y (1 - \theta)^{1-y} \quad y \in \{0, 1\}$$



Learning of parameters of the QDA model

Density estimation in statistics

- We see examples – we do not know the parameters of Gaussians (class-conditional densities)

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

- **ML estimate of parameters** of a multivariate normal $N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for a set of n examples of \mathbf{x}

Optimize log-likelihood: $l(D, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \log \prod_{i=1}^n p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \hat{\boldsymbol{\Sigma}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}})(\mathbf{x}_i - \hat{\boldsymbol{\mu}})^T$$

- How about **class priors**?

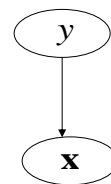
Learning Quadratic discriminant analysis (QDA)

- **Learning Class-conditional distributions**

- Learn parameters of 2 multivariate normal distributions

$$\mathbf{x} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0) \quad \text{for } y = 0$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \quad \text{for } y = 1$$



- Use the density estimation methods

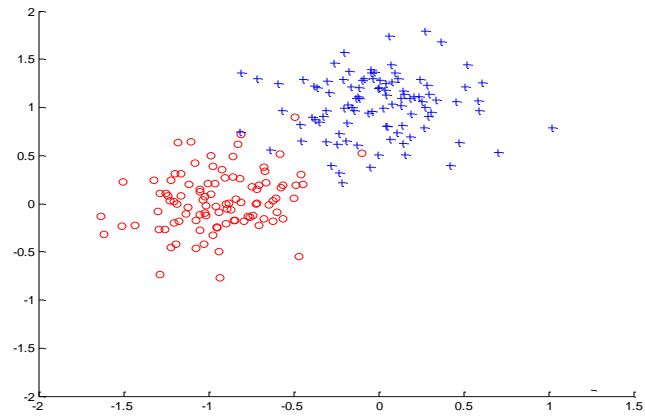
- **Learning Priors on classes (class 0,1)** $y \sim \text{Bernoulli}$

- Learn the parameter of the Bernoulli distribution

- Again use the density estimation methods

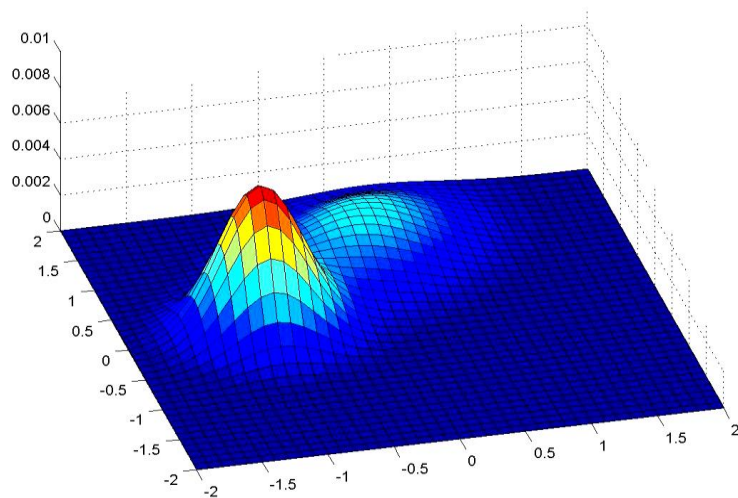
$$p(y, \theta) = \theta^y (1 - \theta)^{1-y} \quad y \in \{0, 1\}$$

QDA



2 Gaussian class-conditional densities

Class conditional densities



QDA: Making class decision

Basically we need to design discriminant functions

- **Posterior of a class** – choose the class with better posterior probability

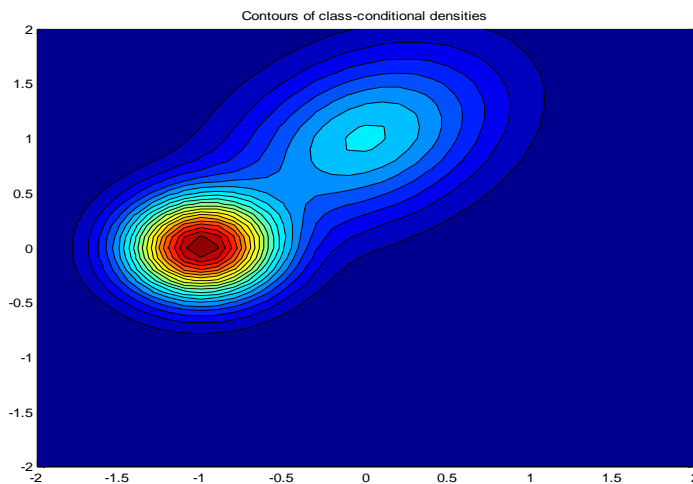
$$\underbrace{p(y=1|\mathbf{x})}_{g_1(\mathbf{x})} > \underbrace{p(y=0|\mathbf{x})}_{g_0(\mathbf{x})} \quad \longrightarrow \quad \begin{array}{l} \text{then } y=1 \\ \text{else } y=0 \end{array}$$

$$p(y=1|\mathbf{x}) = \frac{p(\mathbf{x}|\mu_1, \Sigma_1)p(y=1)}{p(\mathbf{x}|\mu_0, \Sigma_0)p(y=0) + p(\mathbf{x}|\mu_1, \Sigma_1)p(y=1)}$$

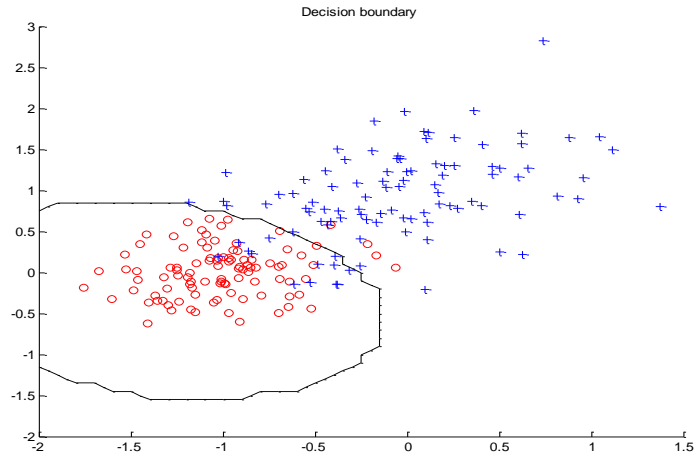
- **Notice it is sufficient to compare:**

$$p(\mathbf{x}|\mu_1, \Sigma_1)p(y=1) > p(\mathbf{x}|\mu_0, \Sigma_0)p(y=0)$$

QDA: Quadratic decision boundary

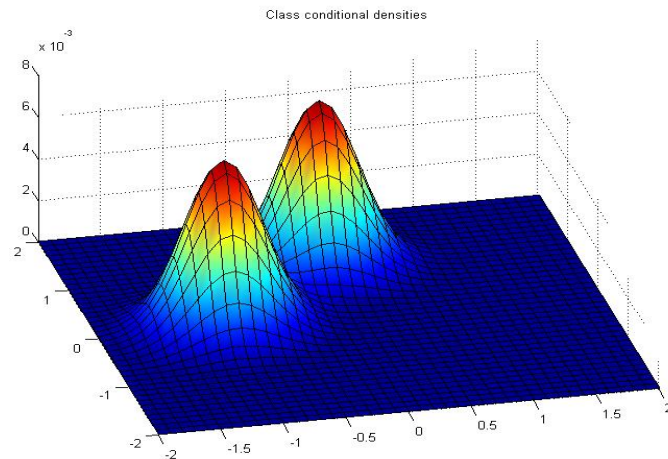


QDA: Quadratic decision boundary

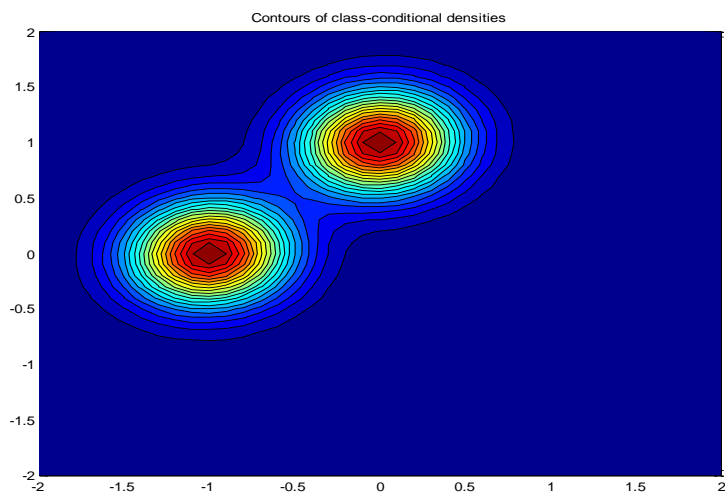


Linear discriminant analysis (LDA)

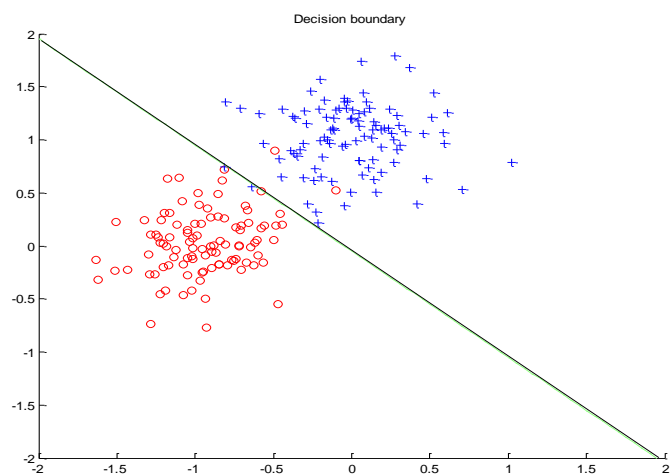
- Assumes covariances are the same $\mathbf{x} \sim N(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}), y = 0$
 $\mathbf{x} \sim N(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}), y = 1$



LDA: Linear decision boundary



LDA: linear decision boundary



Generative classification models

Idea:

1. Represent and learn the distribution $p(\mathbf{x}, y)$
2. Model is able to sample (generate) data instances (\mathbf{x}, y)
3. The model is used to get **probabilistic discriminant functions** $g_0(\mathbf{x}) = p(y = 0 | \mathbf{x})$ $g_1(\mathbf{x}) = p(y = 1 | \mathbf{x})$

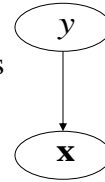
Typical model $p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$

- $p(\mathbf{x} | y) =$ **Class-conditional distributions (densities)**
binary classification: two class-conditional distributions

$$p(\mathbf{x} | y = 0) \quad p(\mathbf{x} | y = 1)$$

- $p(y) =$ **Priors on classes** - probability of class y
binary classification: Bernoulli distribution

$$p(y = 0) + p(y = 1) = 1$$



Naïve Bayes classifier

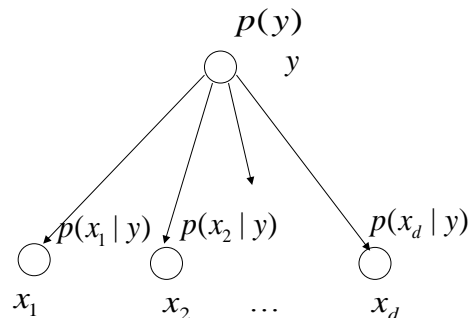
A generative classifier model with an additional simplifying assumption:

- All input attributes are conditionally independent of each other given the class.
- One of the basic ML classification models (often performs very well in practice)

So we have:

$$p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$$

$$p(\mathbf{x} | y) = \prod_{i=1}^d p(x_i | y)$$



Learning parameters of the model

Much simpler density estimation problems

- We need to learn:
 $p(\mathbf{x} | y = 0)$ and $p(\mathbf{x} | y = 1)$ and $p(y)$
- Because of the assumption of the conditional independence we need to learn:
for every input variable i : $p(x_i | y = 0)$ and $p(x_i | y = 1)$
- **Much easier if the number of input attributes is large**
- **Also, the model gives us a flexibility to represent input attributes of different forms !!!**
- E.g. one attribute can be modeled using the Bernoulli, the other using Gaussian density, or a Poisson distribution

Making a class decision for the Naïve Bayes

Discriminant functions

- **Posterior of a class** – choose the class with better posterior probability

$$p(y = 1 | \mathbf{x}) > p(y = 0 | \mathbf{x}) \quad \text{then } y=1 \\ \text{else } y=0$$

$$p(y = 1 | \mathbf{x}) = \frac{\left(\prod_{i=1}^d p(x_i | \Theta_{1,i}) \right) p(y = 1)}{\left(\prod_{i=1}^d p(x_i | \Theta_{1,i}) \right) p(y = 0) + \left(\prod_{i=1}^d p(x_i | \Theta_{2,i}) \right) p(y = 1)}$$

Next: two interesting questions

(1) Two models with linear decision boundaries:

- Logistic regression
- LDA model (2 Gaussians with the same covariance matrices)
 $x \sim N(\mu_0, \Sigma)$ for $y = 0$
 $x \sim N(\mu_1, \Sigma)$ for $y = 1$

- **Question: Is there any relation between the two models?**

(2) Two models with the same gradient:

- Linear model for regression
- Logistic regression model for classification

have the same gradient update

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \sum_{i=1}^n (y_i - f(\mathbf{x}_i)) \mathbf{x}_i$$

- **Question: Why is the gradient the same?**

Logistic regression and generative models

• Two models with linear decision boundaries:

- Logistic regression
- Generative model with 2 Gaussians with the same covariance matrices
 $x \sim N(\mu_0, \Sigma)$ for $y = 0$
 $x \sim N(\mu_1, \Sigma)$ for $y = 1$

- **Question: Is there any relation between the two models?**

Answer: Yes, the two models are related !!!

- When we have **2 Gaussians with the same covariance matrix** the probability of y given \mathbf{x} has the form of a logistic regression model !!!

$$p(y = 1 | \mathbf{x}, \mu_0, \mu_1, \Sigma) = g(\mathbf{w}^T \mathbf{x})$$

Logistic regression and generative models

- Members of the exponential family can be often more naturally described as

$$f(\mathbf{x} | \boldsymbol{\theta}, \boldsymbol{\varphi}) = h(x, \boldsymbol{\varphi}) \exp \left\{ \frac{\boldsymbol{\theta}^T \mathbf{x} - A(\boldsymbol{\theta})}{a(\boldsymbol{\varphi})} \right\}$$

$\boldsymbol{\theta}$ - A location parameter $\boldsymbol{\varphi}$ - A scale parameter

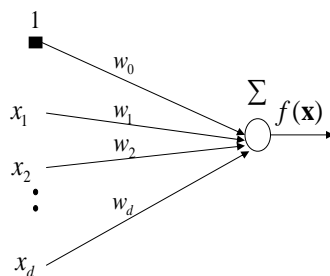
- Claim:** A **logistic regression** is a correct model when class conditional densities are from the same distribution in the exponential family and have **the same scale factor** $\boldsymbol{\varphi}$
- Very powerful result !!!!**
 - We can represent posteriors of many distributions with the same small logistic regression model

CS 2750 Machine Learning

The gradient puzzle ...

Linear regression

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$



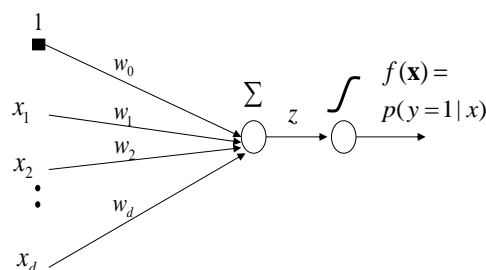
Gradient update:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \sum_{i=1}^n (y_i - f(\mathbf{x}_i)) \mathbf{x}_i$$

$$\text{Online: } \mathbf{w} \leftarrow \mathbf{w} + \alpha (y - f(\mathbf{x})) \mathbf{x}$$

Logistic regression

$$f(\mathbf{x}) = p(y=1 | \mathbf{x}, \mathbf{w}) = g(\mathbf{w}^T \mathbf{x})$$



Gradient update:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \sum_{i=1}^n (y_i - f(\mathbf{x}_i)) \mathbf{x}_i$$

$$\text{Online: } \mathbf{w} \leftarrow \mathbf{w} + \alpha (y - f(\mathbf{x})) \mathbf{x}$$

The same



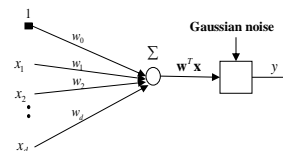
CS 2750 Machine Learning

The gradient puzzle ...

- The **same simple gradient update rule** derived for both the linear and logistic regression models
- Where the magic comes from?
- Under the **log-likelihood** measure the function models and the models for the output selection fit together:

- **Linear model + Gaussian noise**

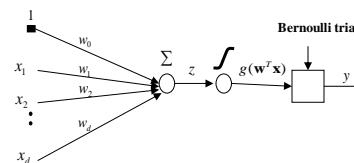
$$y = \mathbf{w}^T \mathbf{x} + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$



- **Logistic + Bernoulli**

$$y = \text{Bernoulli}(\theta)$$

$$\theta = p(y = 1 | \mathbf{x}) = g(\mathbf{w}^T \mathbf{x})$$



Generalized linear models (GLIMs)

Assumptions:

- The conditional mean (expectation) is:
 - Where $f(\cdot)$ is a **response function**
- Output y is characterized by an exponential family distribution with a conditional mean μ

Examples:

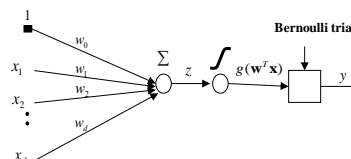
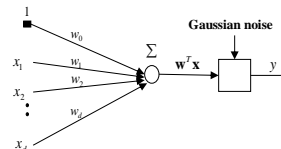
- **Linear model + Gaussian noise**

$$y = \mathbf{w}^T \mathbf{x} + \varepsilon \quad \varepsilon \sim N(0, \sigma^2)$$

- **Logistic + Bernoulli**

$$y \approx \text{Bernoulli}(\theta)$$

$$\theta = g(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$$



Generalized linear models (GLIMs)

- A canonical response functions $f(\cdot)$:
 - encoded in the sampling distribution

$$p(\mathbf{x} | \boldsymbol{\theta}, \boldsymbol{\varphi}) = h(x, \boldsymbol{\varphi}) \exp \left\{ \frac{\boldsymbol{\theta}^T \mathbf{x} - A(\boldsymbol{\theta})}{a(\boldsymbol{\varphi})} \right\}$$

- Leads to a simple gradient form
- Example: Bernoulli distribution

$$p(x | \mu) = \mu^x (1 - \mu)^{1-x} = \exp \left\{ \log \left(\frac{\mu}{1 - \mu} \right) x + \log(1 - \mu) \right\}$$
$$\theta = \log \left(\frac{\mu}{1 - \mu} \right) \quad \mu = \frac{1}{1 + e^{-\theta}}$$

- Logistic function matches the Bernoulli
-