**Georgia Tech**

# Naïve Bayes and Logistic Regression

Nakul Gopalan
Georgia Tech

These slides are adopted based on slides from Le Song, Eric Eaton, Mahdi Roozbahani and Chao Zhang.

# Regression

$y(x)$



Temp. $y(x)$

Days $x$

Training Data

Test Data

Mean Sq. Error

$$L(\theta) = \frac{1}{n}\sum_n \left(y(x_i) - \hat{y}(x_i)\right)^2$$

Sparse ness $\rightarrow$ LASSO

$0.1 \quad \longrightarrow 0.01$

$0.01 \quad \longrightarrow \quad 0.0001$

## Polynomial Regression

$$y(x) = 1 + a_1 x^1 + a_2 x^2 \cdots a_n x^n$$

$$z_n = x^n$$

$$y(z) = 1 + a_1 z_1 + a_2 z_2 \cdots + a_n z_n$$

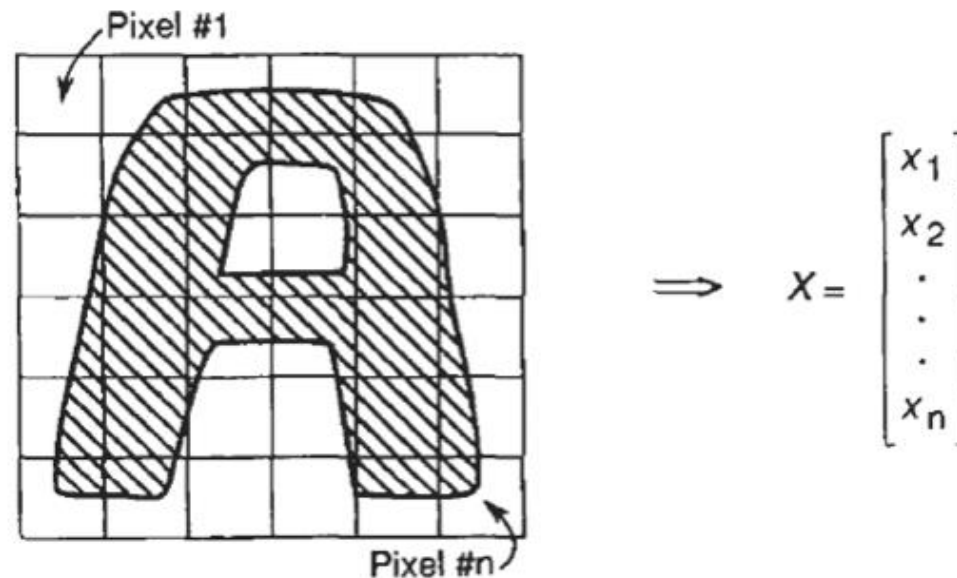Regularization $\rightarrow$

Overfitting & generalization

Lagrange Multiplier :

Ridge $\rightarrow L(\theta) = \frac{1}{2n}\sum_n (y(x_i) - \hat{y}(x_i))^2 + \lambda\|\theta\|^2$

$\rightarrow L(\theta) = \frac{1}{n}\sum (y(x_i) - \hat{y}(x_i))^2 + \lambda\|\theta\|$

# Outline

- Generative and Discriminative Classification ⬅

- The Logistic Regression Model

- Understanding the Objective Function

- Gradient Descent for Parameter Learning
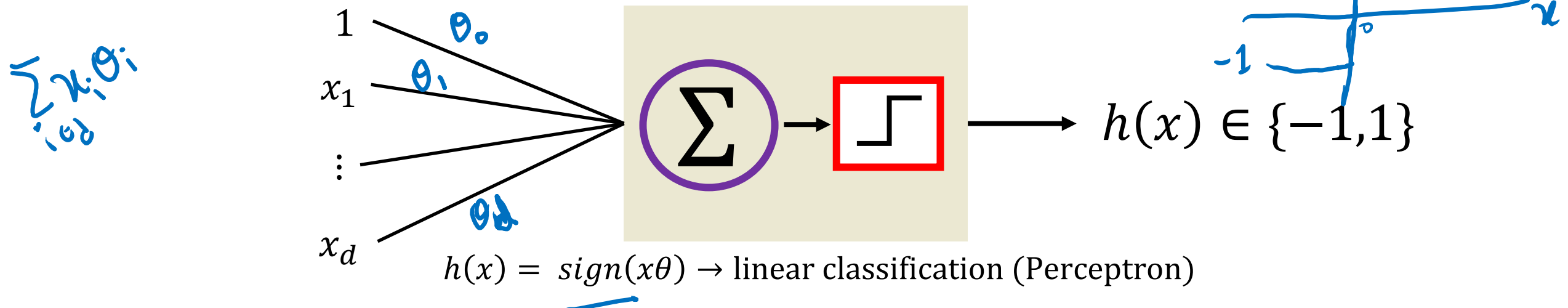
- Multiclass Logistic Regression

# Classification

- Represent the data



Pixel #1

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix}$$

Pixel #n

- A label is provided for each data point, eg., $y \in \{-1, +1\}$

  image of character

- Classifier

$\sum x_i \theta_i$

1   $\theta_0$

$x_1$   $\theta_1$

$\vdots$

$x_d$   $\theta_d$

$\Sigma \rightarrow \boxed{\text{⊓}} \rightarrow h(x) \in \{-1, 1\}$

$h(x) = sign(x\theta) \rightarrow$ linear classification (Perceptron)

$y$   $+1$   $x$   $-1$
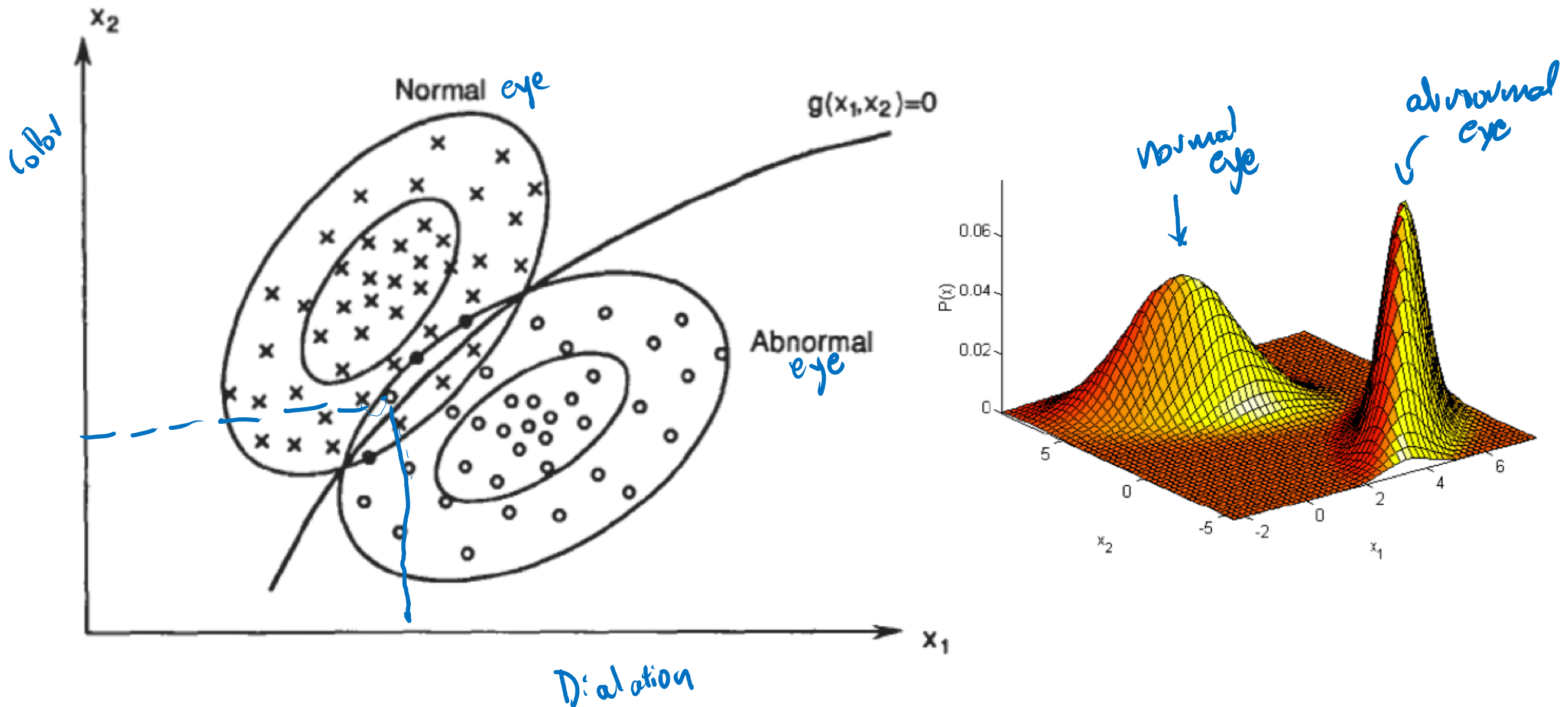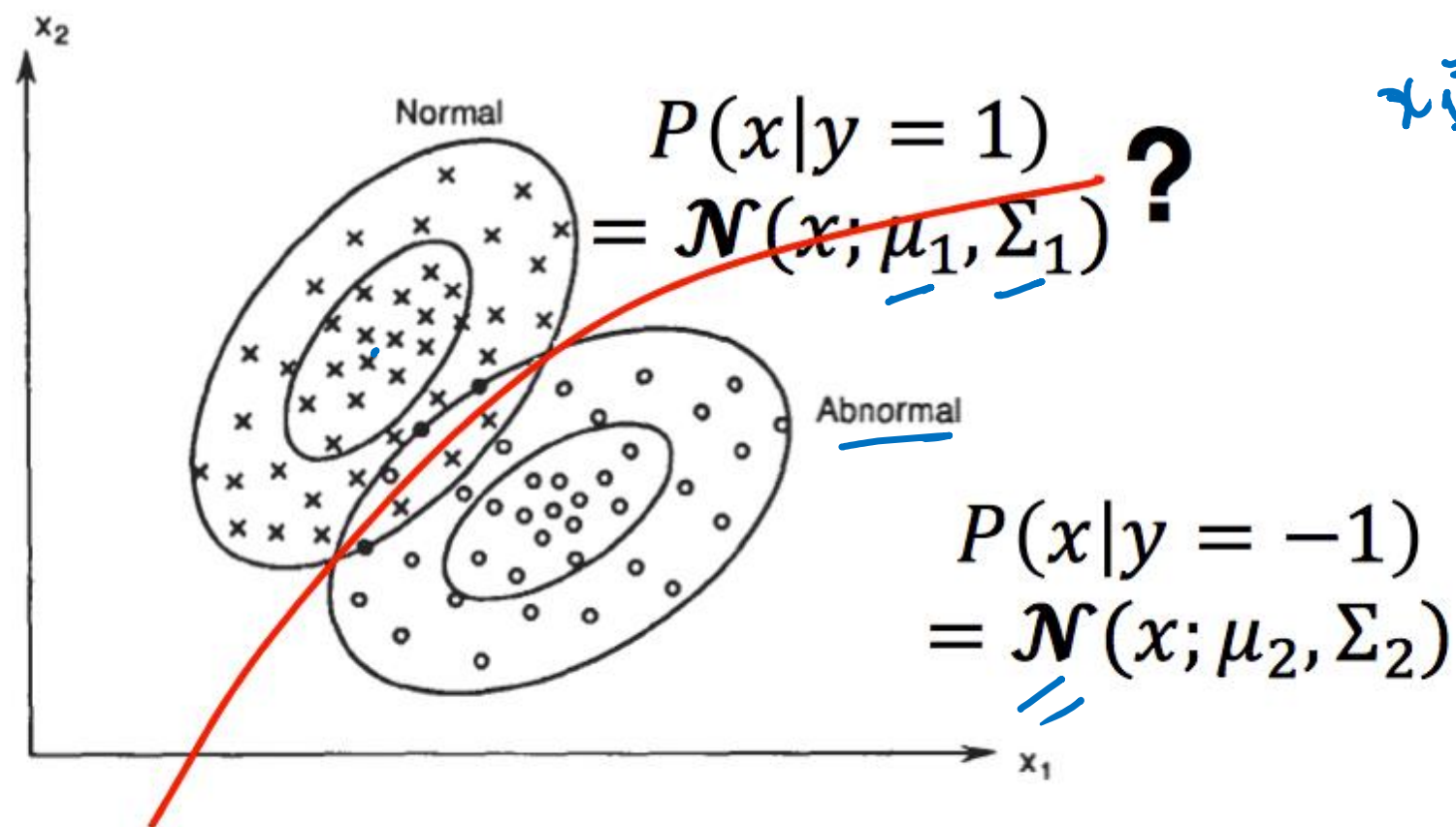
# Decision Making: Dividing the Feature Space

- Distributions of sample from normal (positive class) and abnormal (negative class) tissues

# How to Determine the Decision Boundary?

- Given class conditional distribution: $P(x|y = 1), P(x|y = -1)$, and class prior: $P(y = 1), P(y = -1)$



$P(x|y = 1) = \mathcal{N}(x; \mu_1, \Sigma_1)$ **?**

$P(x|y = -1) = \mathcal{N}(x; \mu_2, \Sigma_2)$

Normal

Abnormal

$x_3 = x_{Pupil} = 0.05 \text{ mm}$
dilation

$P(x_3 | y = normal)$

$P(x_3 | y = abnormal)$

# Bayes Decision Rule



likelihood — Normal dist. — Prior

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x,y)}{\sum_z P(x,y)}$$

posterior — normalization constant — marginalization

Probability of normal eye given $x_1$

$x_1 = [\phantom{P} 0.005, \phantom{} 0.1]$

$P(y|x)$

| | diolation $\downarrow$ | intensity | $P(y)$ / $P(x)$ |
|---|---|---|---|
| | $x^1$ | $x^2$ | Class |
| 1 | 0.005 | 0.1 | +1 |
| 2 | | 0.5 | -1 |
| 1000 | 0.004 | 0.01 | +1 |

$P(x|y=+1) = \mu_1 = \left[ x_{in}^1 \quad \frac{x_w^2}{2} \right]^T$

$\frac{1}{2} \sum_z \left[ \quad \right]$

Prior: $P(y)$

Likelihood (class conditional distribution : $p(x|y) = \mathcal{N}(x|\mu_y, \Sigma_y)$

Posterior: $P(y|x) = \frac{P(y)\mathcal{N}(x|\mu_y, \Sigma_y)}{\sum_y P(y)\mathcal{N}(x|\mu_y, \Sigma_y)}$

# Bayes Decision Rule

- Learning: prior: $p(y)$, class conditional distribution : $p(x|y)$

- The poster probability of a test point
  *Posterior*

$$q_i(x) := P(y = i|x) = \frac{P(x|y)P(y)}{P(x)}$$

$+1/-1$

*Def.*

$q_i(x) > q_j(x)$

$P(x|y=i)P(y=i) > P(x|y=i)P(i)$

*Prior* $\frac{P(x|y=i)}{P(x)}$

- Bayes decision rule:
  - If $q_i(x) > q_j(x)$, then $y = i$, otherwise $y = j$

$1 > q_j/q_i$

$0 > \ln(q_j/q_i)$

- Alternatively:
  - If ratio $l(x) = \frac{P(x|y=i)}{P(x|y=j)} > \frac{P(y=j)}{P(y=i)}$, then $y = i$, otherwise $y = j$

  - Or look at the log-likelihood ratio $h(x) = -\ln\frac{q_i(x)}{q_j(x)} = 0$ (decision boundary)

# What do People do in Practice?

- Generative models
  - Model prior and likelihood explicitly
  - "Generative" means able to generate synthetic data points
  - Examples: Naive Bayes, Hidden Markov Models

- Discriminative models
  - Directly estimate the posterior probabilities
  - No need to model underlying prior and likelihood distributions
  - Examples: Logistic Regression, SVM, Neural Networks

*Mean, Var.*

*compute $P(x|y)$*

*↳ unseen points*

# Generative Model: Naive Bayes

- Use Bayes decision rule for classification

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

- But assume $p(x|y = 1)$ is fully factorized : Dimensions are independent.

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

$$p\left(x = [x^1, x^2] \,\middle|\, y = 1\right) = \prod_{i=1}^{2}\left(p(x^i|y = 1)\right)$$

- Or the variables corresponding to each dimension of the data are independent given the label

# "Naïve" conditional independence assumption

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)} = \frac{P(x,y)}{P(x)}$$

Joint probability model:

$$P(A,B) = P(A|B)P(B)$$

$$P(x, y_{label=1}) = \mathrm{P}(x_1, \dots, x_d, y_{label=1}) = \mathrm{P}(x_1|x_2, , \dots, x_d, y_{label=1})P(x_2, \dots, x_d, y_{label=1})$$

$$= \mathrm{P}(x_1|x_2, , \dots, x_d, y_{label=1})P(x_2|x_3 \dots, x_d, y_{label=1})P(x_3, \dots, x_d, y_{label=1})$$

$$= \dots$$

too hard

$$= \mathrm{P}(x_1|x_2, , \dots, x_d, y_{label=1})P(x_2|x_3 \dots, x_d, y_{label=1}) \dots P(x_{d-1}|x_d, y_{label=1})P(x_d|y_{label=1})\mathrm{P}(y_{label=1})$$

Naïve Bayes assumption: let's rewrite it as:

$$P(x, y_{label=1}) = P(x_1|y_{label=1})P(x_2|y_{label=1}) \dots P(x_n|y_{label=1})P(y_{label=1}) =$$

$$P(y_{label=1}) \prod_{i=1}^{d} P(x_i|y_{label=1})$$

Gaussian naïve Bayes
A typical assumption

[Example](Example)

# Naïve Bayes cat vs dog!

| $y$ | $x_1$ cat height kg | $x_2$ cat ht. ft | $x_3$ tail length |
|-----|-----|-----|-----|
| Cat | 2 | • 1 | 1 |
| dog | 6 | 2 | 2 |
| dog | 30 | 3 | 0 |
| cat | 1 | 0.5 | 0.5 |

Joint dist. $= P(x=[10,2,0], y=cat) = P(y=cat) P(10 \mid cat)$. $P(ht=2 \mid cat)$. $P(0 \mid cat)$

# Administrative things

- Project team composition due this weekend

- Quiz out today. Let us know if you have problems with it, and take as many as you can, it will only help!!

- Homework due next week. Deadlines will start getting closeby from now.

# Naïve Bayes



For Cat vs dog:

$$P(y = \text{Cat} \mid x) = P(y = \text{dog} \mid x)$$

$x \rightarrow$ features like tail length, height, weight.

# Naïve Bayes

$$P(x \mid y) \cong \prod_i P(x_i \mid y)$$

<span style="color:red">independence assumption b/w features</span>

$$P(x \mid y = cat) \cong P(\text{tail length} = 1 \mid y = cat) \cdot P(wt = 2 \mid y = cat) \cdot$$

$$P(ht = 1 \mid y = cat)$$

<span style="color:red">learning large conditional dis$^n$.</span>

Actually

<span style="color:red">hard</span>

$$P(x \mid y = cat) = P(\text{tail length} = 1 \mid wt = 2, ht = 1, y = cat) \cdot$$

$$P(wt = 2 \mid ht = 1, y = cat) \cdot P(ht = 1 \mid y = cat)$$

# Discriminative Models

- Directly estimate decision boundary $h(x) = -\ln \frac{q_i(x)}{q_j(x)}$ or posterior distribution $p(y|x)$
  - Logistic regression, Neural networks
  - Do not estimate $p(x|y)$ and $p(y)$

- Why discriminative classifier?
  - Avoid difficult density estimation problem → Generative model
  - Empirically achieve better classification results

# Outline

- Generative and Discriminative Classification

- The Logistic Regression Model

- Understanding the Objective Function

- Gradient Descent for Parameter Learning

- Multiclass Logistic Regression

# Gaussian Naïve Bayes

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

Prior

Naïve ind. assumption

$P(A|B) = \dfrac{P(B|A)\, P(A)}{P(B)}$

$$\prod_{i=1}^{d} p(x_i|y = 1, \mu_{1i}, \sigma_{1i})$$

$$= \prod_{i=1}^{d} \frac{1}{\sqrt{2\pi}\sigma_{1i}} \exp\left(-\frac{1}{2\sigma_{1i}^2}(x_{1i} - \mu_{1i})^2\right)$$

Normal dist.

class    which feature

Cat vs dog

Prior: $p(y = 1) = \pi_1$

$\pi_1$

$\pi_2$

$x_1 = 0.5$

$x_2 = 0.4$

$x_1 = 0.5$

$x_2 = 0.6$

Posterior: $p(y = 1 \mid x, \mu, \sigma, \pi)$

$$= \frac{\pi_1 \prod_{i=1}^{d} \frac{1}{\sqrt{2\pi}\sigma_{1i}} \exp\left(-\frac{1}{2\sigma_{1i}^2}(x_i - \mu_{1i})^2\right)}{\sum_{k=1}^{2} \pi_k \prod_{i=1}^{d} \frac{1}{\sqrt{2\pi}\sigma_{ki}} \exp\left(-\frac{1}{2\sigma_{ki}^2}(x_i - \mu_{ki})^2\right)}$$

Previous slide

Marginalization by summing over y or labels

labels

$P(y_k)$

class    feature i.d.

$P(x) = \sum_{k \in labels} P(x, y_k)$

$= \sum_{k} P(x|y_i) \cdot P(y_k)$

get $\exp(\ln(u))$ of numerator and denominator

$u = \exp(\ln(u))$

$$= \frac{\exp\left(-\sum_{i=1}^{d}\left(\frac{1}{2\sigma_{1i}^2}(x_i - \mu_{1i})^2 + \log\sigma_{1i} + C\right) + \log\pi_1\right)}{\sum_{k=1}^{2}\exp\left(-\sum_{i=1}^{d}\left(\frac{1}{2\sigma_{ki}^2}(x_i - \mu_{ki})^2 + \log\sigma_{ki} + C\right) + \log\pi_k\right)}$$

$$= \frac{\exp\left(-\sum_{i=1}^{d}\left(\frac{1}{2\sigma_i^2}(x_i - \mu_{1i})^2 + \log\sigma_i + C\right) + \log\pi_1\right)}{\sum_{k=1}^{2}\exp\left(-\sum_{i=1}^{d}\left(\frac{1}{2\sigma_i^2}(x_i - \mu_{ki})^2 + \log\sigma_i + C\right) + \log\pi_k\right)}$$

dividing by Nu.

$-\mu_2 - (-\mu_1)$

$\dfrac{c_1}{\sum_{k\in 1,2} c_k} = \dfrac{c_1}{c_1 + c_2} = \dfrac{1}{1 + \frac{c_2}{c_1}}$

$$= \frac{1}{1 + \exp\left(-\sum_{i=1}^{d}\left(x_i \frac{1}{\sigma_i}(\mu_{1i} - \mu_{2i}) + \frac{1}{\sigma_i^2}(\mu_{1i}^2 - \mu_{2i}^2)\right) + \log\frac{\pi_2}{\pi_1}\right)}$$

$\dfrac{\exp(c_2)}{\exp(c_1)} = \exp(c_2 - c_1)$

$\underbrace{\sum_i \theta_i x_i}$

$\underbrace{\theta_0}$

$$P(y = 1|x) = \cfrac{1}{1 + \exp\left(-\sum_{i=1}^{d}\left(x_i \frac{1}{\sigma_i}(\mu_{1i} - \mu_{2i}) + \frac{1}{\sigma_i^2}(\mu_{1i}^2 - \mu_{2i}^2)\right) + \log\frac{\pi_2}{\pi_1}\right)}$$

Number of parameters:

$2d + 1 \rightarrow d$ $mean, d$ $variance, and$ $1$ $for$ $prior$

$$P(y = 1|x) = \frac{1}{1 + \exp[-(\sum_i(\theta_i x_i) + \theta_0)]} = \frac{1}{1 + \exp(-s)}$$

Number of parameters = $d + 1 \rightarrow \theta_0, \theta_1, \theta_2, \dots, \theta_d$

Why not directly learning $P(y = 1|x)$ or $\theta$ parameters?

Gaussian Naïve Bayes is a subset of logistic regression

# Logistic function for posterior probability

Many equations can give us this shape

Let's use the following function:

$$g(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}} \qquad s = x\theta$$



$g(s)$

This formula is called sigmoid function

It is easier to use this function for optimization

$\frac{1}{1+1} \quad \therefore e^0 = 1$

derivative is hard to compute

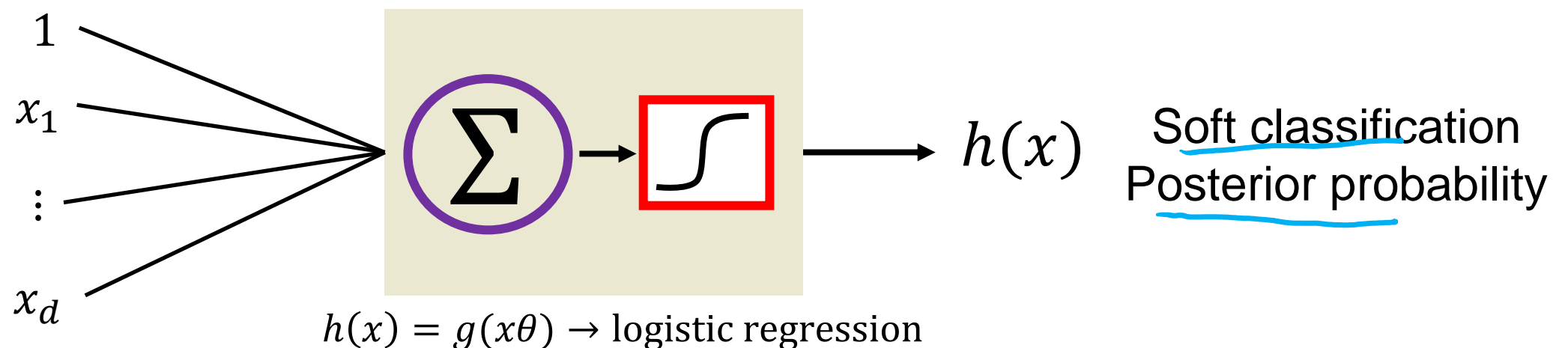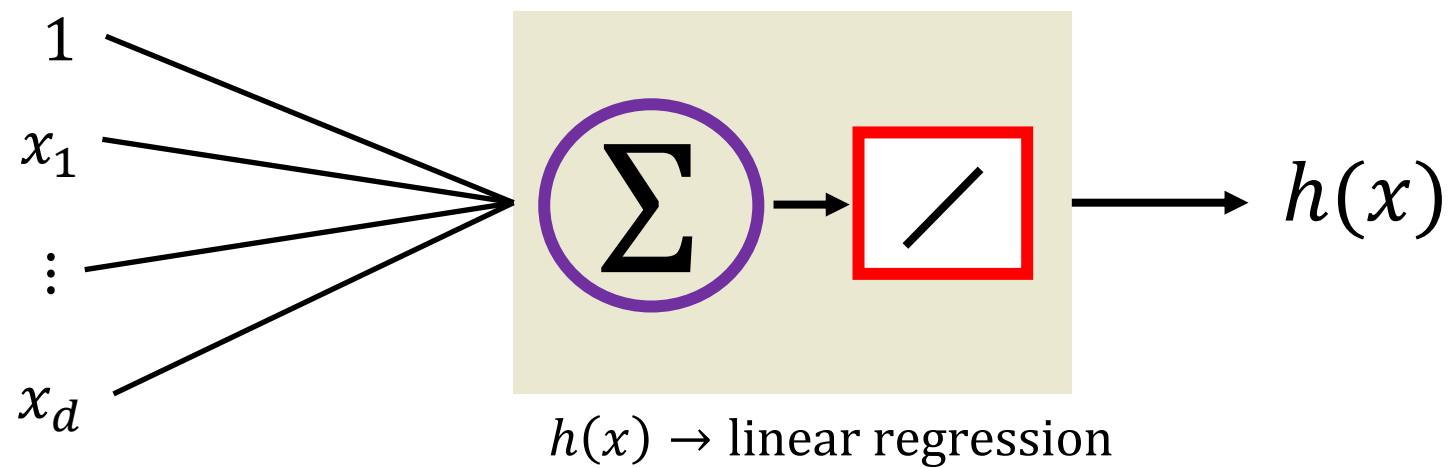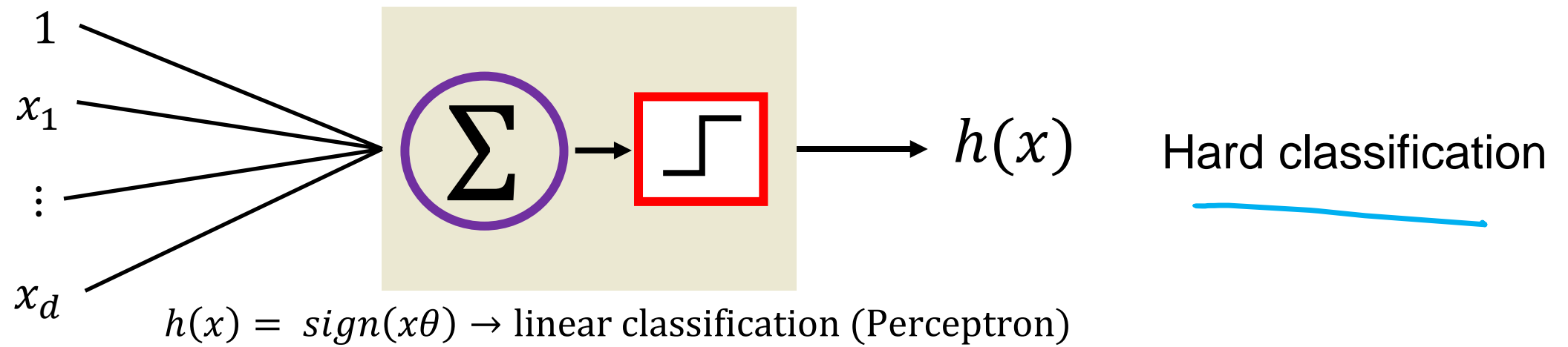$$g(s) = \frac{e^s}{1+e^s} = \frac{1}{1+e^{-s}}$$

## Sigmoid Function

$$s = \sum_{i=0}^{d} x_i \theta_i = \theta_0 + \theta_1 x_1 + \cdots + \theta_d x_d$$



$h(x) = g(x\theta) \rightarrow \text{logistic regression}$

Soft classification
Posterior probability

$$s = \sum_{i=0}^{d} x_i \theta_i = \theta_0 + \theta_1 x_1 + \cdots + \theta_d x_d$$

**Three linear models**



$$h(x) = sign(x\theta) \rightarrow \text{linear classification (Perceptron)}$$

Hard classification

$$h(x) \rightarrow \text{linear regression}$$

$$h(x) = g(x\theta) \rightarrow \text{logistic regression}$$

Soft classification
Posterior probability

$g(s)$ is interpreted as probability

Example: Prediction of heart attacks

Input $x$: cholesterol level, age, weight, finger size, etc.

$g(s)$: probability of heart attack within a certain time

We can't have a hard prediction here

$s = x\theta$    Let's call this risk score

$$g(s) = \frac{1}{e^{-s} + 1}$$

$$h_\theta(x) = p(y|x) = \begin{cases} g(s), & y = 1 \\ 1 - g(s), & y = 0 \end{cases}$$

Using posterior probability directly

# Logistic regression model

$$p(y|x) = \begin{cases} \dfrac{1}{1 + \exp(-x\theta)} & y = 1 \\[4mm] 1 - \dfrac{1}{1 + \exp(-x\theta)} = \dfrac{\exp(-x\theta)}{1 + \exp(-x\theta)} & y = 0 \end{cases}$$

$= g(s) \quad = \dfrac{1}{1 + e^{-s}}$

We need to find $\theta$ parameters, let's set up log-likelihood for **n** datapoints

$$l(\theta) := \log \prod_{i=1}^{n} p(y_i, |x_i, \theta)$$

→ data points

↳ $0^{th}$ class ↳ $i^{th}$ feature

$$= \sum_{i} \theta^T x_i^T (y_i - 1) - \log(1 + \exp(-x_i\theta))$$

This form is concave, negative of this form is convex

# The gradient of $l(\theta)$

$$l(\theta) := \log \prod_{i=1}^{n} p(y_i, |x_i, \theta)$$

$$= \sum_i \theta^T x_i^T (y_i - 1) - \log(1 + \exp(-x_i\theta))$$

$\frac{d}{d\theta} \log(f(\theta)) = \frac{1}{f(\theta)} \cdot \frac{d}{d\theta} f(\theta)$

- Gradient

$$\frac{\partial l(\theta)}{\partial \theta} = \sum_i x_i^T (y_i - 1) + x_i^T \frac{\exp(-x_i\theta)}{1 + \exp(-x_i\theta)}$$
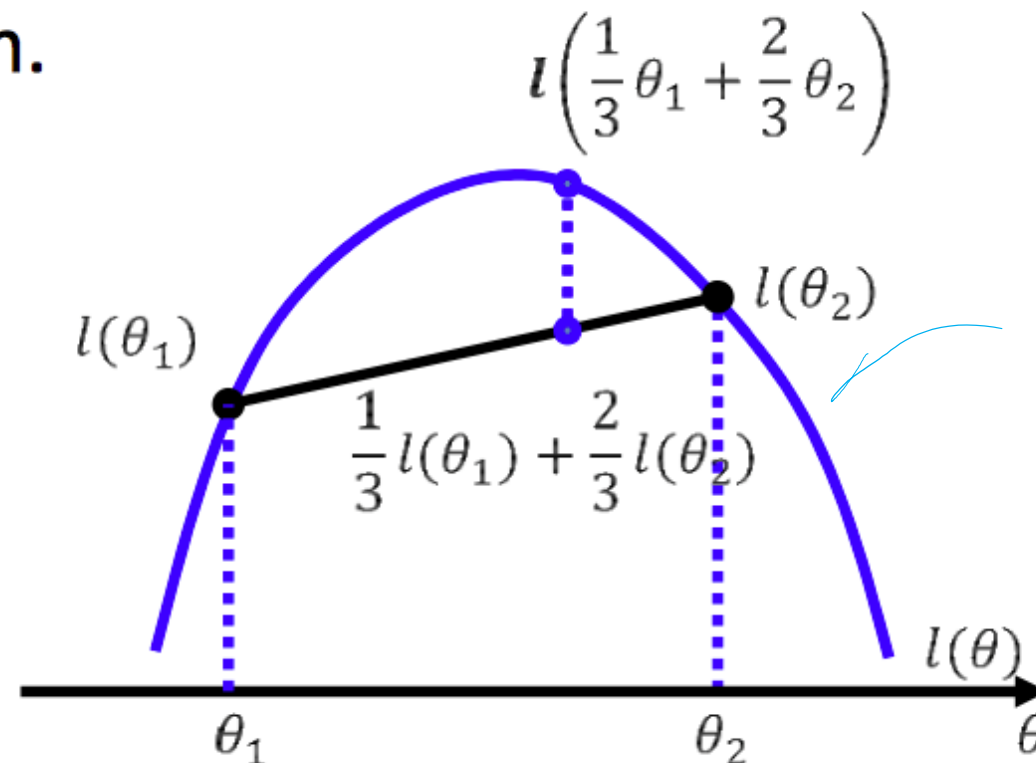
- Setting it to 0 does not lead to closed form solution

# The Objective Function

- Find $\theta$, such that the conditional likelihood of the labels is maximized

$$\max_{\theta} l(\theta) := log \prod_{i=1}^{n} p(y_i, |x_i, \theta)$$

- Good news: $l(\theta)$ is concave function of $\theta$, and there is a single global optimum.

$$l\left(\frac{1}{3}\theta_1 + \frac{2}{3}\theta_2\right)$$

$a f(\theta_1) + b f(\theta_2) \leq$

$f(a\theta_1 + b\theta_2)$

$l(\theta_2)$

$l(\theta_1)$

$$\frac{1}{3}l(\theta_1) + \frac{2}{3}l(\theta_2)$$

$l(\theta)$

$\theta_1$   $\theta_2$   $\theta$

- Bad new: no closed form solution (resort to numerical method)

# Gradient Descent

- One way to solve an *unconstrained* optimization problem is gradient descent

- Given an initial guess, we *iteratively* refine the guess by taking the direction of the negative gradient

- Think about going down a hill by taking the steepest direction at each step

$$\theta_{k+1} = \theta_k - \gamma_k \nabla f(\theta_k)$$

- Update rule

$$x_{k+1} = x_k - \gamma_k \nabla f(x_k)$$

$\gamma_k$ is called the step size or learning rate

# Gradient Ascent(concave)/Descent(convex) algorithm

- Initialize parameter $\theta^0$

- Do

$$\theta^{t+1} \leftarrow \theta^t + \eta \sum_i x_i^T(y_i - 1) + x_i^T \frac{\exp(-x_i\theta)}{1 + \exp(-x_i\theta)}$$

*learning rate*

$\frac{d\ \ell(\theta)}{d\theta)}$ → log liklihood of seeing n data points

*params. for decision boundary.*

- While the $||\theta^{t+1} - \theta^t|| > \epsilon$

we do not need:
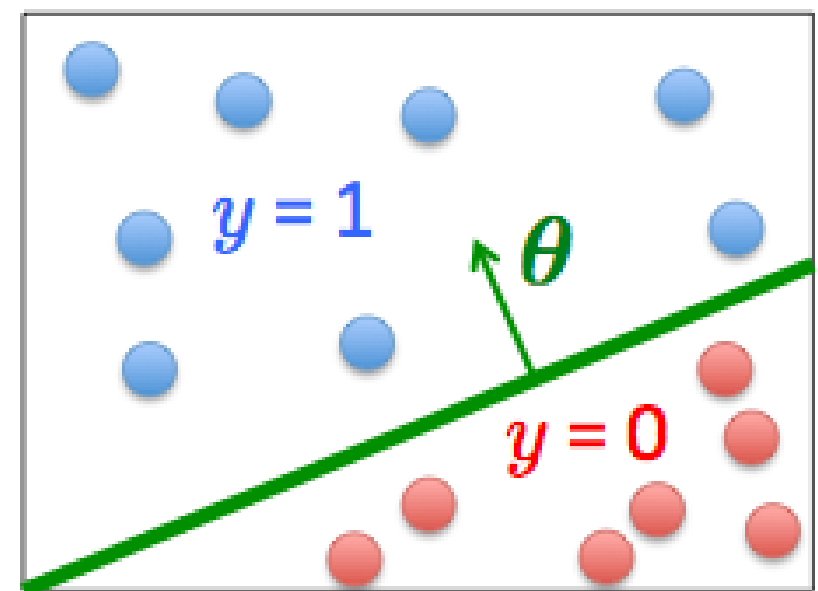$P(x), P(y), P(x|y)$
$P(y|x)$

# Logistic Regression

$$g(s) = \frac{e^s}{1 + e^s} = \frac{1}{1 + e^{-s}}$$

$$s = x\theta$$



$g(s)$

$s$

$x\theta$ should be large <u>negative</u> values for negative instances

$x\theta$ should be large <u>positive</u> values for positive instances

- Assume a threshold and...
  - Predict $y = 1$ if $g(s) \geq 0.5$
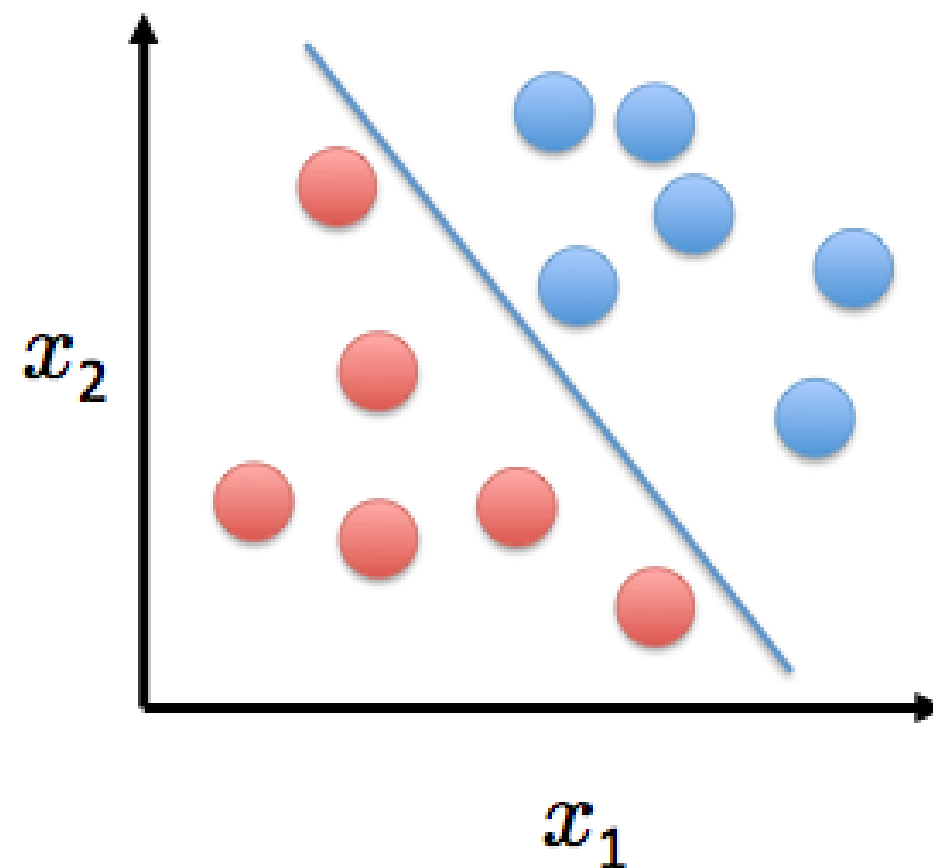  - Predict $y = 0$ if $g(s) < 0.5$



$y = 1$

$\theta$

$y = 0$

# Outline

- Generative and Discriminative Classification

- The Logistic Regression Model

- Understanding the Objective Function

- Gradient Descent for Parameter Learning
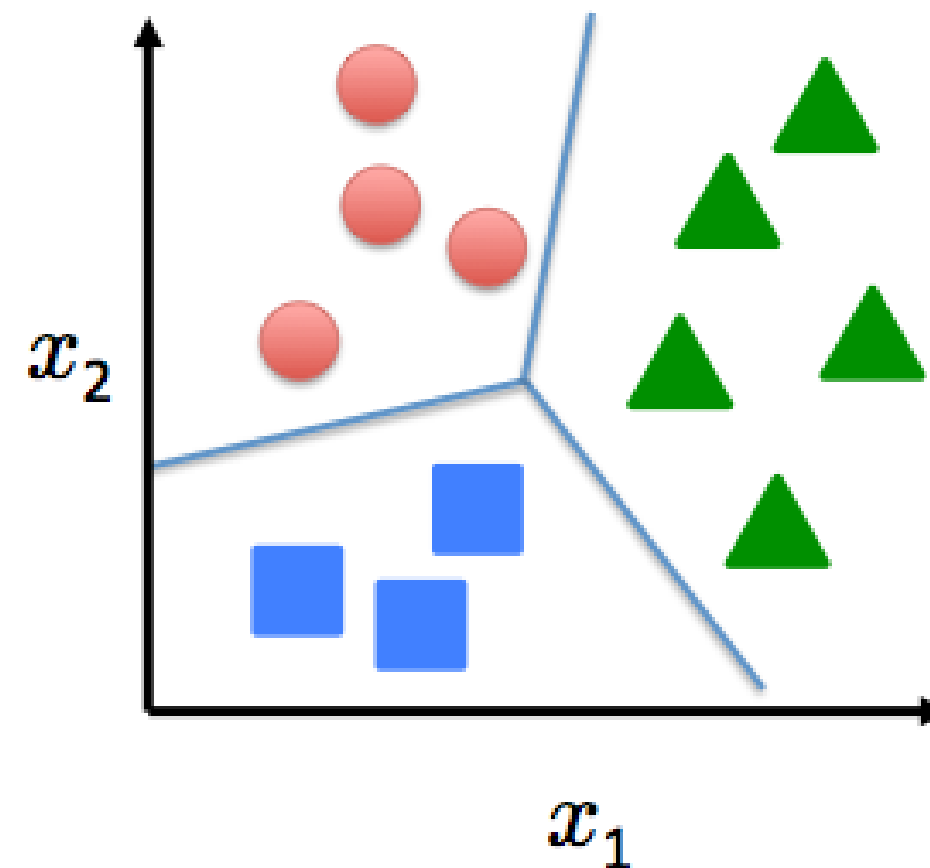
- Multiclass Logistic Regression ⬅

# Multiclass Logistic Regression
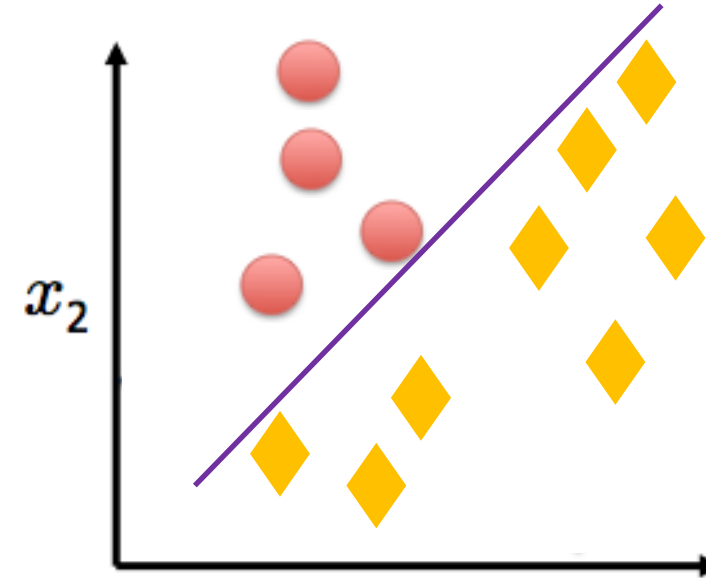


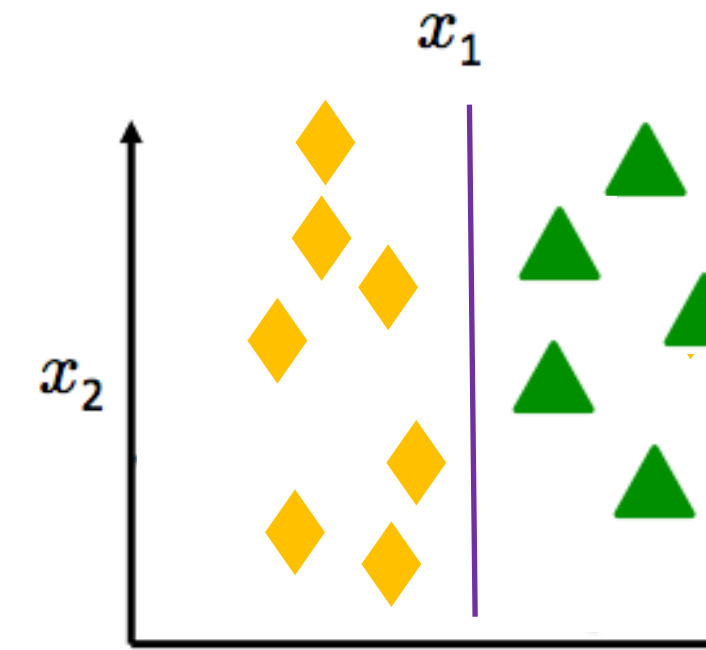Binary classification:

Multi-class classification:

Disease diagnosis:     healthy / cold / flu / pneumonia

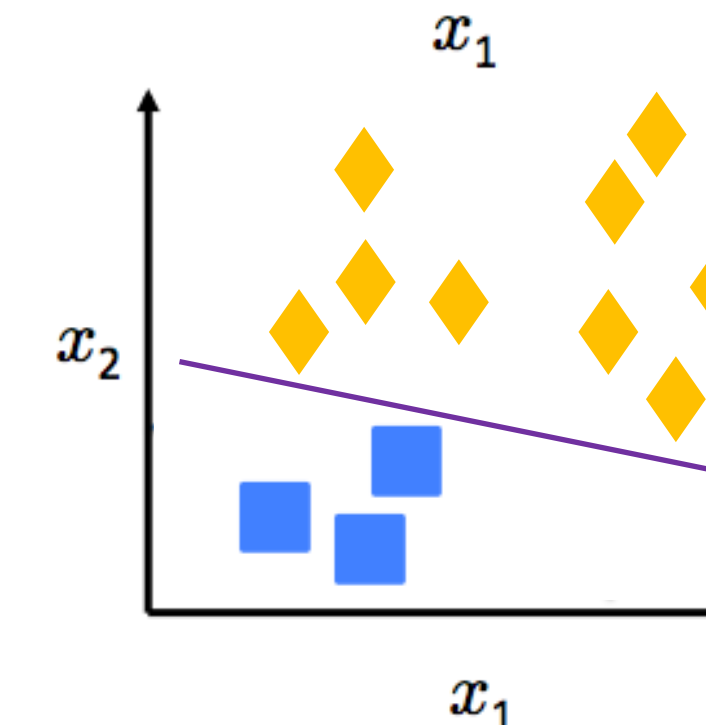Object classification:  desk / chair / monitor / bookcase
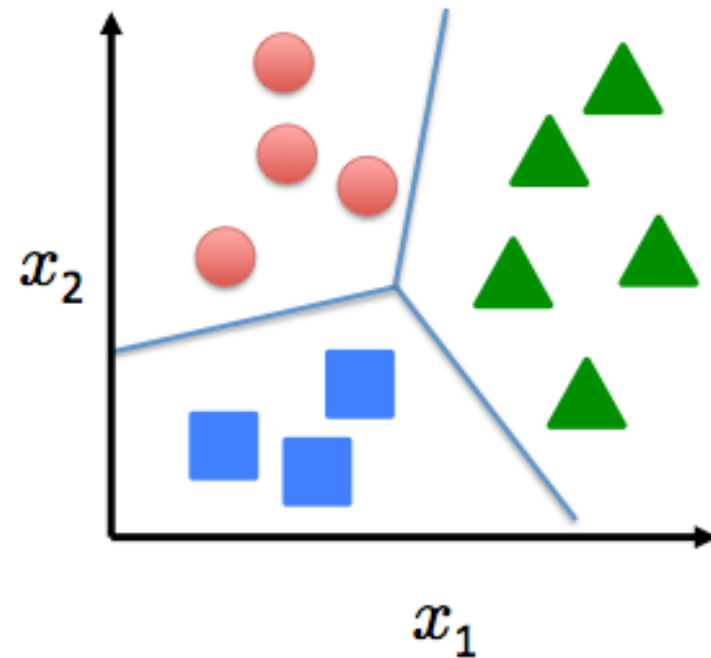
# One-vs-all (one-vs-rest)

$h_\theta^1(x)$

Multi-class classification:

$h_\theta^2(x)$

$h_\theta^{(i)}(x) = p(y = 1 | x, \theta) \ (i = 1,2,3)$

$= g(s) = \dfrac{1}{1 + e^{-s}} \quad s = \theta \cdot x$

$h_\theta^3(x)$

**One-vs-all (one-vs-rest)**

Train a logistic regression $h_\theta^{(i)}(x)$ for each class $i$

To predict the label of a new input $x$, pick class $i$ that maximizes:

$$\max_i h_\theta^{(i)}(x)$$

# Take-Home Messages

- Generative and Discriminative Classification

- The Logistic Regression Model

- Understanding the Objective Function

- Gradient Descent for Parameter Learning

- Multiclass Logistic Regression