

بنام خدا

سمینار درس بازشناسی الگو

کاربرد گوسیها و مخلوط آنها در بازشناسی الگو

استاد: دکتر کبیر

ارایه: سامان پروانه

فهرست مطالب:

- دلیل استفاده از مدلسازی در بازشناسی

- تابع گوسی

- تابع چگالی تابع گوسی

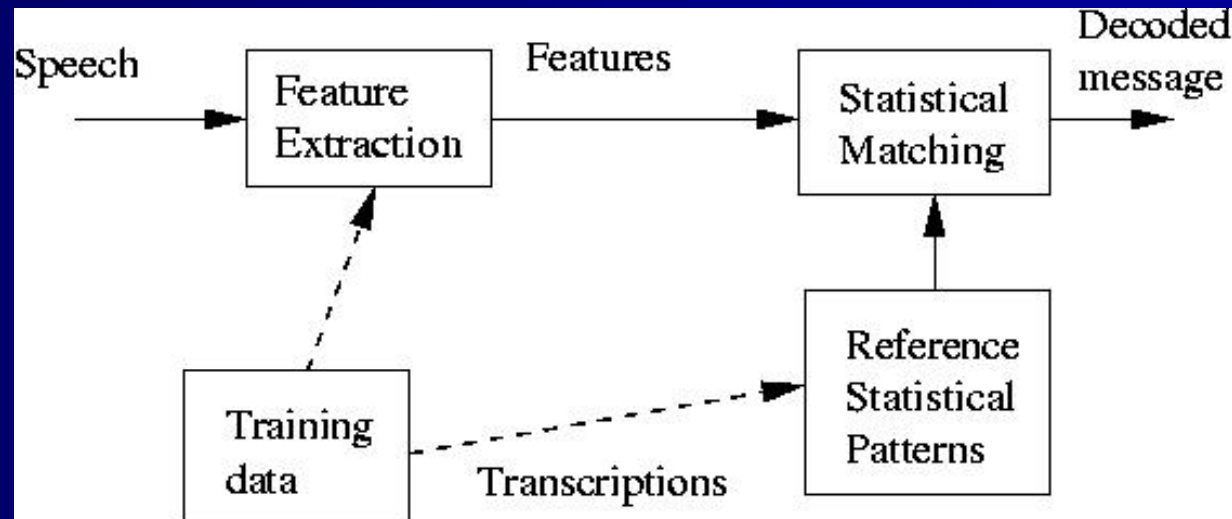
- تخمین پارامترهای تابع گوسی

- مخلوط گوسیها

- تخمین پارامترها

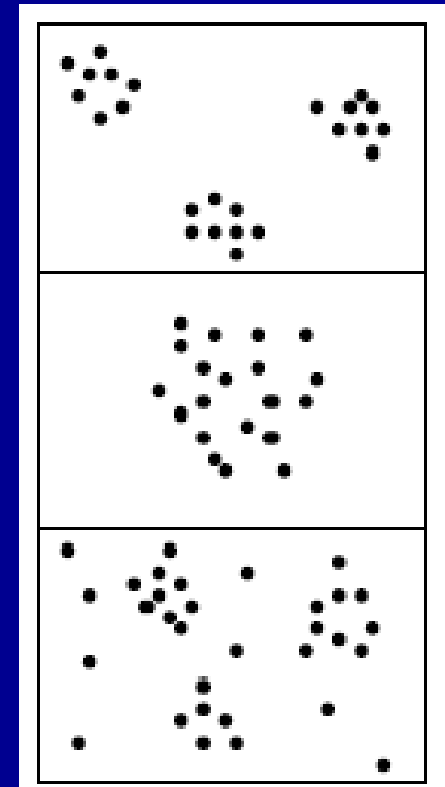
- بررسی همسانی

چرا داده‌ها را مدلسازی می‌کنیم؟



آموزش بدون ناظر:

بعضی اوقات آسان

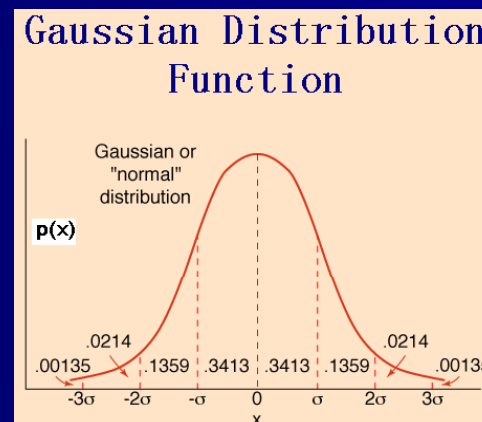


بعضی اوقات غیرممکن

بعضی اوقات در وضعیت بینابینی

تابع گوسی:

- یک توزیع متداول است:
 - بر طبق تئوری حد مرکزی: جمع متغیرهای مستقل به گوسی میل می کند.
 - نویز و خطا معمولاً با توزیع گوسی مدل می شود.
 - معمولاً وقتی اطلاعات دقیق وجود ندارد از این توزیع استفاده می شود برای مثال در بازده بورس، قد مردان و ...
- این توزیع، توزیع نرمال نیز نامیده می شود و دارای شکل زنگی-شکل است.



تابع چگالی احتمال گوسی:

X یک متغیر تصادفی d -بعدی است. $\mathbf{x} \in \mathbb{R}^d$

$$g(\boldsymbol{\mu}, \boldsymbol{\Sigma})(\mathbf{x}) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\boldsymbol{\Sigma})}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

$\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ ← ماتریس کوواریانس

← بردار میانگین

$\boldsymbol{\mu}$ شامل میانگین هر بعد است. $\mu_i = E(x_i)$

$$\boldsymbol{\Sigma} = \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1n} \\ c_{21} & c_{22} & \cdots & c_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \cdots & c_{nn} \end{bmatrix}$$

c_{ii} واریانس

c_{ij} کوواریانس

تابع چگالی احتمال گوسی (ادامه):

$$c_{ij} = E \left((x_i - \mu_i)^T (x_j - \mu_j) \right)$$

$$x_i \text{ and } x_j, i \neq j \longrightarrow c_{ij} = 0 \longrightarrow \text{متعامد}$$

چه موقع ماتریس کوواریانس قطری است؟

انحراف استاندارد متغیر تصادفی X است. $\sqrt{\Sigma}$

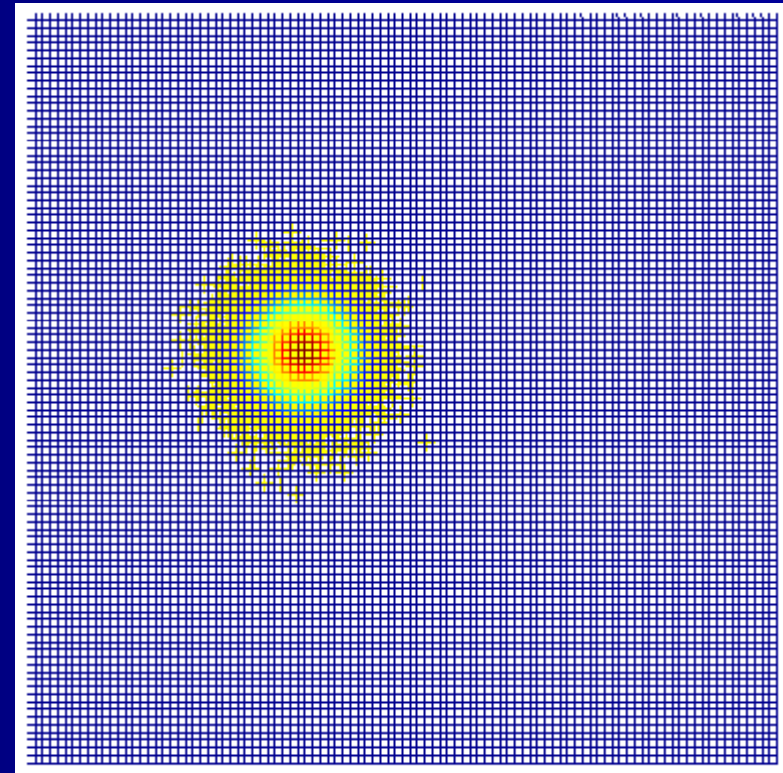
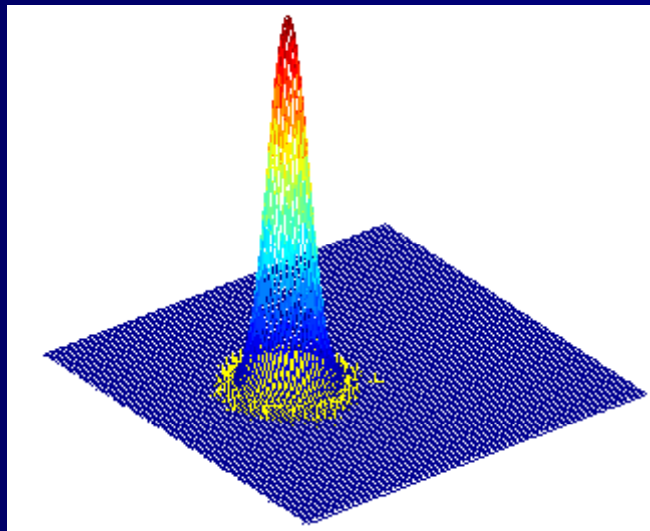
$$x \sim \mathcal{N}(0, \mathbf{I}) \longrightarrow y = \mu + \sqrt{\Sigma} x \longrightarrow y \sim \mathcal{N}(\mu, \Sigma)$$

نمونه فرآیند دایره‌ای:

$$X = \{x_1, x_2, \dots, x_N\}$$

$$\mu = \begin{bmatrix} 730 \\ 1090 \end{bmatrix}$$

$$\Sigma_1 = \begin{bmatrix} 8000 & 0 \\ 0 & 8000 \end{bmatrix}$$

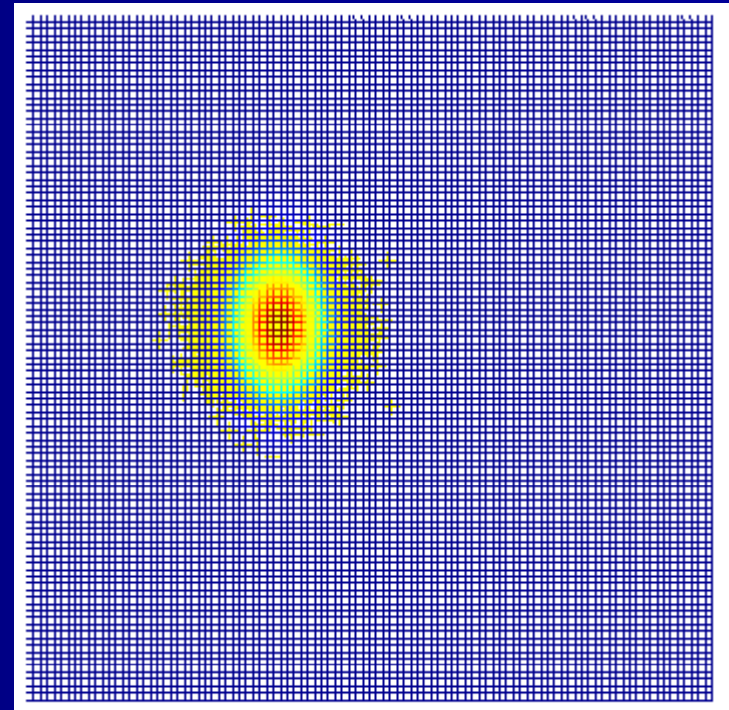
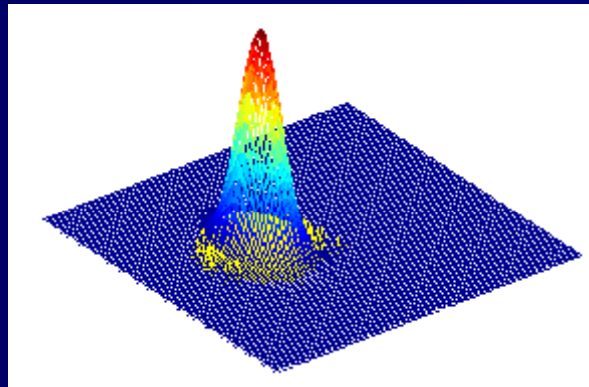


نمونه کوواریانس قطری :

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

$$\boldsymbol{\mu} = \begin{bmatrix} 730 \\ 1090 \end{bmatrix}$$

$$\boldsymbol{\Sigma}_2 = \begin{bmatrix} 8000 & 0 \\ 0 & 18500 \end{bmatrix}$$

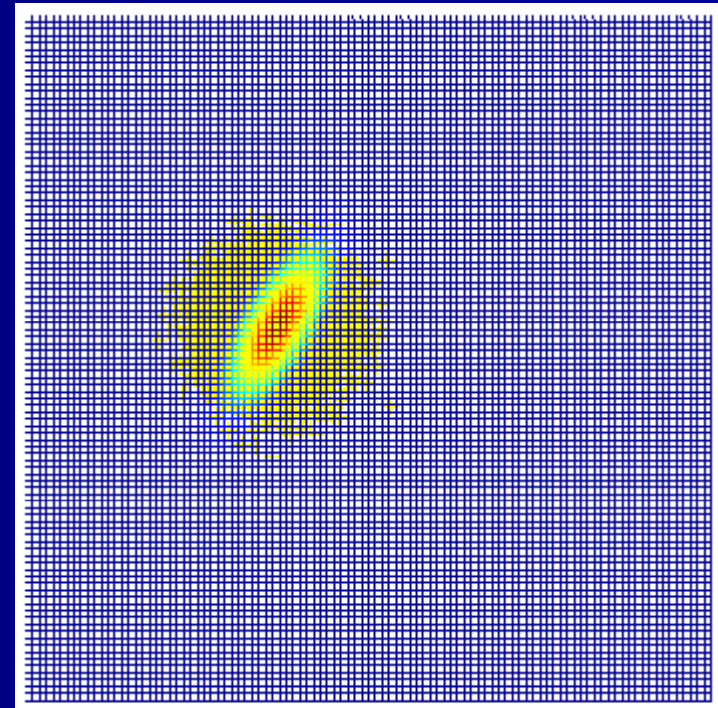
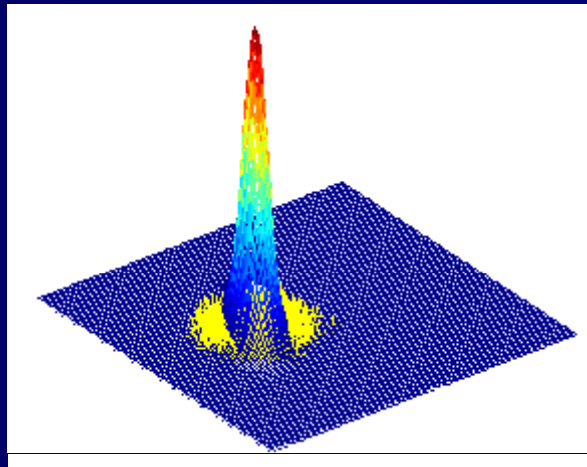


نمونه کوواریانس کامل:

$$X = \{x_1, x_2, \dots, x_N\}$$

$$\mu = \begin{bmatrix} 730 \\ 1090 \end{bmatrix}$$

$$\Sigma_s = \begin{bmatrix} 8000 & 8400 \\ 8400 & 18500 \end{bmatrix}$$



تخمین پارامترهای تابع گوسی از داده:

$$\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

تخمین گر میانگین:

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$$

→ mean(\mathcal{X})

تخمین گر بدون بایاس کوواریانس:

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

→ cov(\mathcal{X})

همسانی یک نمونه نسبت به مدل گوسی:

همسانی (Likelihood) برای طبقه‌بندی استفاده می‌شود.

- با معلوم بودن یک مدل تولید داده (معلوم بودن پارامترها) همسانی برابر است با:

$$p(\mathbf{x}_i | \Theta)$$

- برای مدل گوسی معادل با ارزیابی رابطه زیر است:

$$g_{(\mu, \Sigma)}(\mathbf{x}) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)}$$

- همسانی همزمان (Joint Likelihood): برای یک مجموعه از نمونه‌هایی که بصورت یکنواخت توزیع شده‌اند، حاصلضرب همسانی هر کدام از نمونه‌ها است.

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \rightarrow p(X | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \Theta) = \prod_{i=1}^N p(\mathbf{x}_i | \mu, \Sigma) = \prod_{i=1}^N g_{(\mu, \Sigma)}(\mathbf{x}_i)$$

Log-likelihood بجای Likelihood:

۱- حاصلضرب را به حاصلجمع تبدیل می کند.

$$p(X|\Theta) = \prod_{i=1}^N p(x_i|\Theta) \Leftrightarrow \log p(X|\Theta) = \log \prod_{i=1}^N p(x_i|\Theta) = \sum_{i=1}^N \log p(x_i|\Theta)$$

۲- در حالت گوسی، از محاسبات نمایی جلوگیری می کند.

$$p(x|\Theta) = \frac{1}{\sqrt{2\pi}^d \sqrt{\det(\Sigma)}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$
$$\log p(x|\Theta) = \frac{1}{2} [-d \log(2\pi) - \log(\det(\Sigma)) - (x-\mu)^T \Sigma^{-1}(x-\mu)]$$

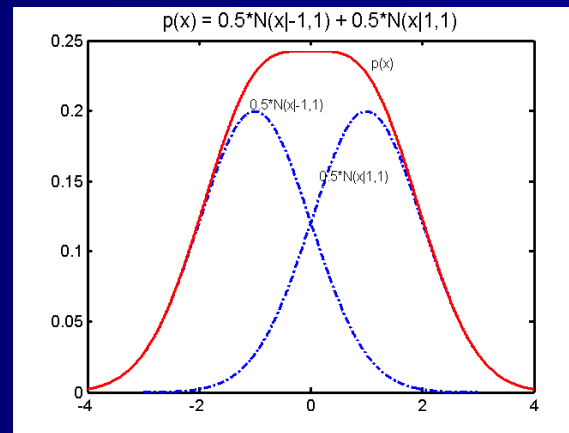
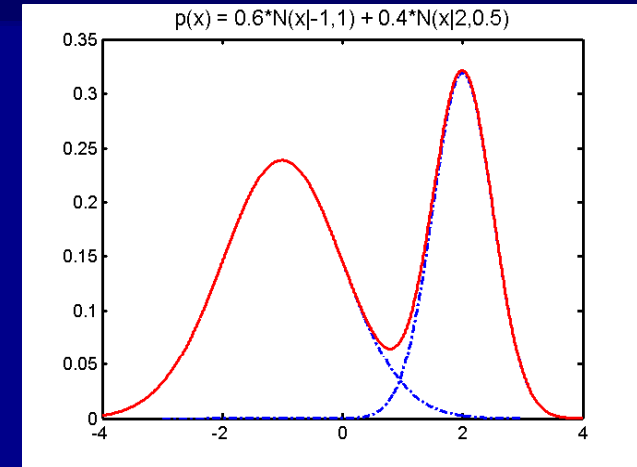
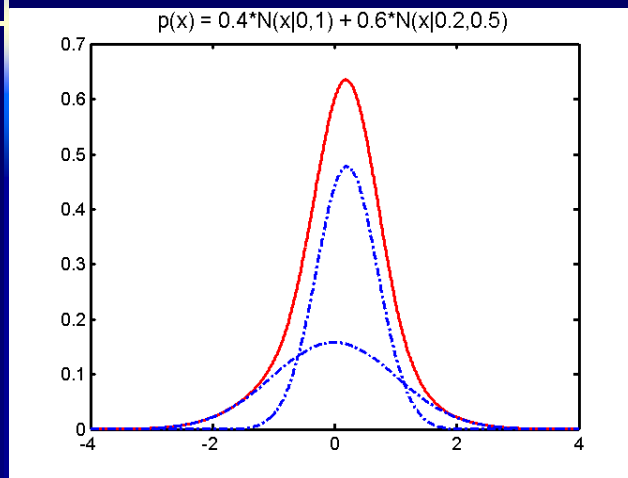
چون Log-likelihood، صعودی یکنوا است در نتیجه روابط مشابه Likelihood است و می تواند بصورت مستقیم برای طبقه بندی استفاده بشود.

$$p(x|\Theta_1) > p(x|\Theta_2) \Leftrightarrow \log p(x|\Theta_1) > \log p(x|\Theta_2)$$

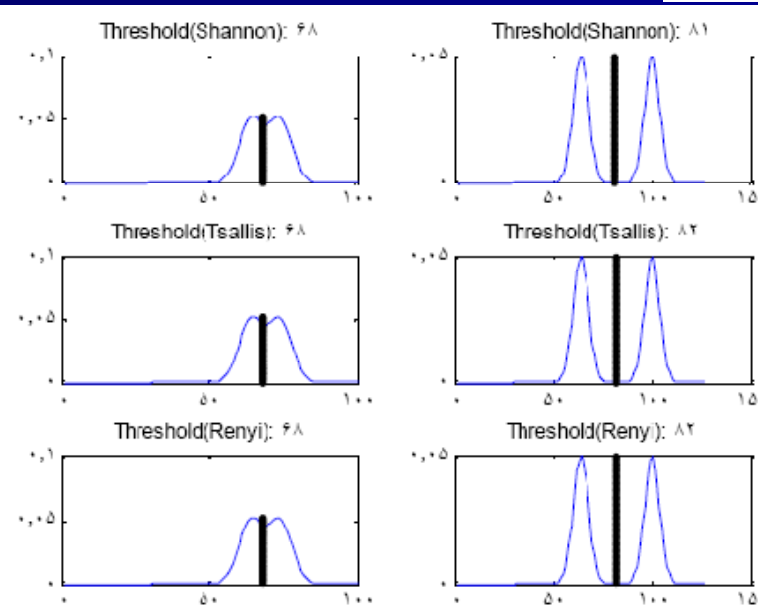
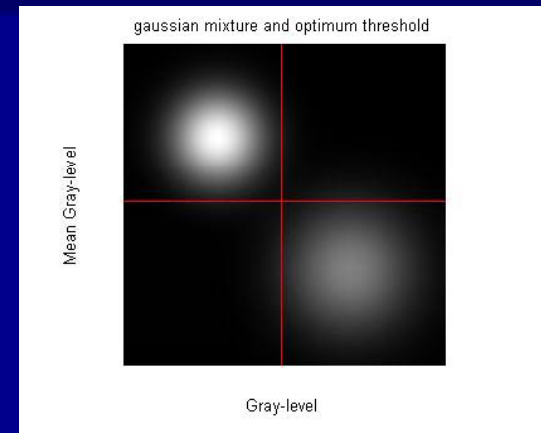
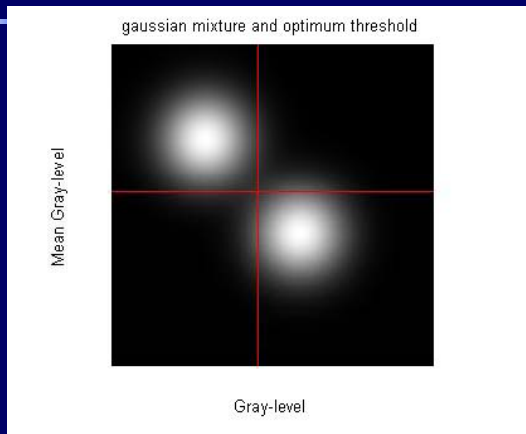
مدل مخلوط گوسی (GMM):

- مدل مخلوط گوسیها
 - بر اساس مدلسازی آماری خوشه‌ها است.
 - خوشه بندی ← یک مسئله تخمین پارامتر
- مشابه K-means که جمع فاصله نمونه‌ها تا مراکز خوشه را مینیمم می‌کرد، از الگوریتم EM استفاده می‌کنیم تا همسانی (Likelihood) مدل مخلوط گوسی را ماکزیمم کنیم.

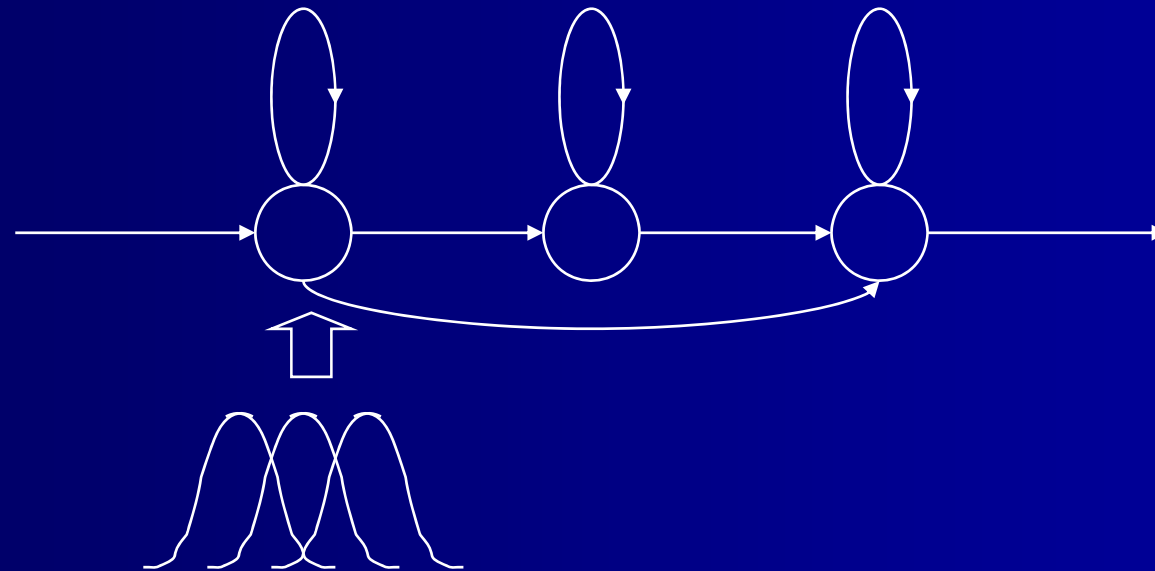
چند مخلوط گوسی نمونه (حالت یک بعدی):



کاربرد مخلوط گوسی در پردازش تصویر:



کاربرد مخلوط گوسی در HMM:

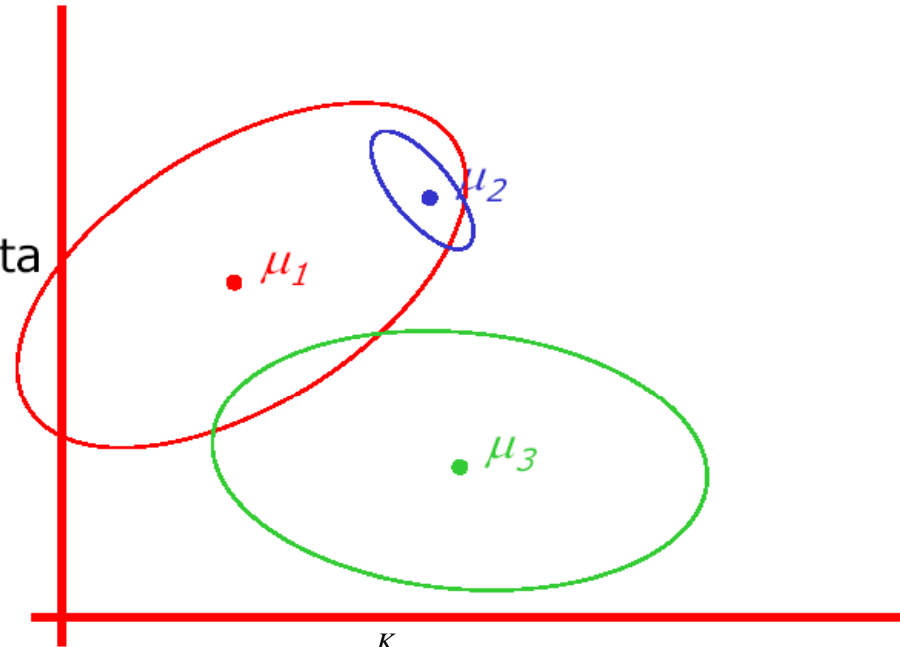


The General GMM assumption

- There are k components. The i 'th component is called ω_i
- Component ω_i has an associated mean vector μ_i
- Each component generates data from a Gaussian with mean μ_i and covariance matrix Σ_i

Assume that each datapoint is generated according to the following recipe:

1. Pick a component at random. Choose component i with probability $P(\omega_i)$.
2. Datapoint $\sim N(\mu_i, \Sigma_i)$

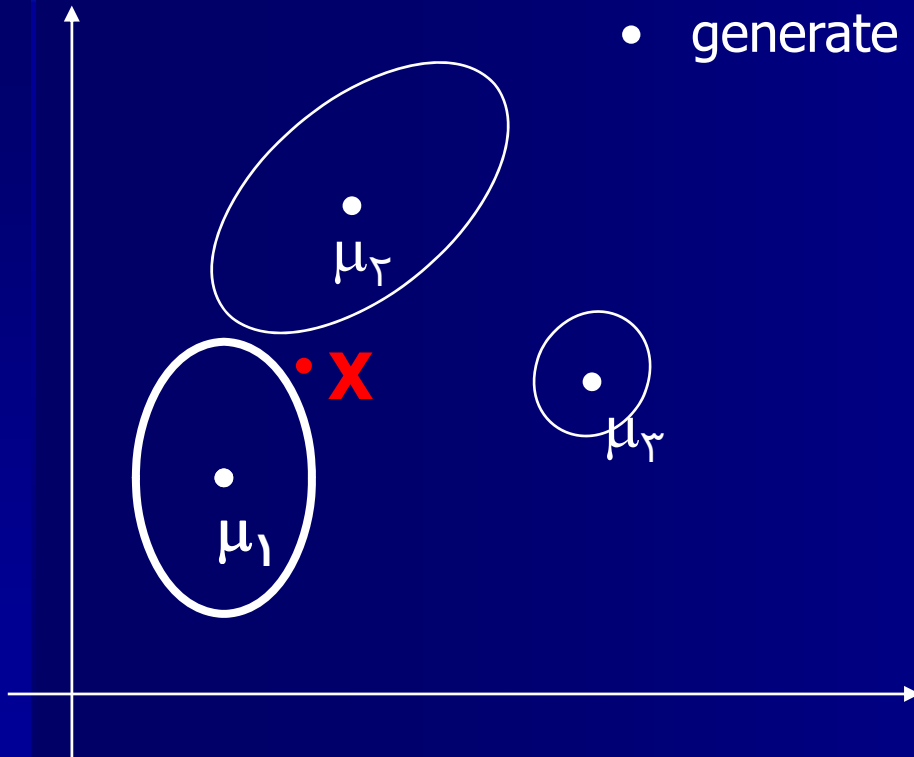


$$p(x) = \sum_{i=1}^K P(\omega_i) N(x | \mu_i, \Sigma_i)$$

$$P(\omega_i) > 0, \quad \sum_{i=1}^K P(\omega_i) = 1$$

مخلوط گوسی (ادامه):

- choose component with probability w_k
- generate $X \sim N(\mu_k, \sigma_k)$



چند نکته دربارهٔ مخلوط گوسی:

- تابع چگالی $p(x)$ جمع (مخلوط) وزن یافتهٔ تعدادی تابع چگالی گوسی است.
- این جمع توابع چگالی با متغیری که جمع وزن یافتهٔ متغیرهای دیگر است متفاوت می‌باشد: $Y = aX_1 + bX_2$

ترکیب مخلوط گوسی:

Class ۱

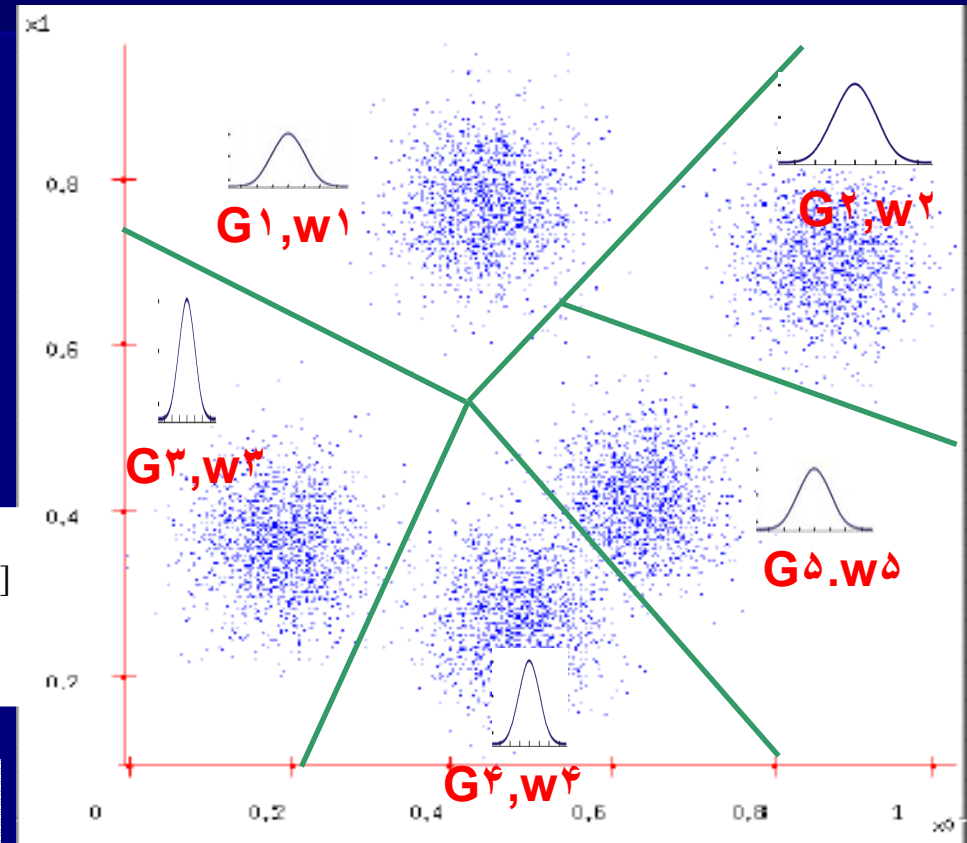
$$P(C_j | X) = P(X | C_j) \cdot \frac{P(C_j)}{P(X)}$$

$$P(X | C_j) = \sum_{k=1}^{Nc} w_k G_k$$

$$G_k \equiv p(X | G_i) = \frac{1}{(\sqrt{\pi})^{d/2} |V_i|^{1/2}} \cdot e^{[-1/2(X-\mu_i)^T V_i^{-1} (X-\mu_i)]}$$

Variables: μ_i , V_i , w_k

We use EM (estimate-maximize) algorithm to approximate this variables.



تعداد گوسیها چند عدد باشد؟

Problem of scoring models with different complexities

Model too flexible \Rightarrow overfit the data \Rightarrow high variance

Model too restrictive \Rightarrow can't fit the data \Rightarrow high bias

Bias-variance tradeoff: compromise

احتمال مشترک و شرطی:

اگر دو متغیر X و Y را داشته باشیم و بخواهیم احتمالاتی که شامل X و Y است را ترکیب کنیم:

■ احتمال مشترک (Joint Probability): آنرا با $p(x,y)$ نمایش می‌دهیم و احتمال این است که X مقدار x^1 را بگیرد و Y مقدار y^1 را بگیرد. در اصل احتمال این است که هر دوی این حالتها رخ بدهند:
 X برابر x^1 بشود و Y برابر y^1 بشود.

■ احتمال شرطی (Joint Probability): آنرا با $p(x/y)$ نمایش می‌دهیم و بدین مفهوم است که در صورتیکه Y مقدار y را گرفت، X مقدار x را بگیرد (صرفنظر از غیرممکن بودن اینکه Y مقدار y را بگیرد)

■ رابطه: $p(x,y) = p(x|y)p(y)$ ^{۲۳}

اهمیت احتمال شرطی:

- احتمال شرطی در حوزه‌های مختلفی از جمله طبقه‌بندی، تئوری تصمیم‌گیری، پیش‌بینی و موقعیتهای مشابه دارای اهمیت است.
- دلیل این اهمیت در کاربردهای بالا این است که طبقه‌بندی، تصمیم‌گیری، پیش‌بینی و ... را بر اساس یک ملاک انجام می‌دهیم. بنابراین می‌خواهیم احتمال یک نتیجه را با معلوم بودن ملاک بدانیم. $p(r/e)$.
- در طبقه‌بندی، ملاک مقادیر اندازه‌گیریها یا ویژگیهای X است که طبقه‌بندی باید بر اساس آنها باشد. نتایج ممکن نیز، کلاسها (w) هستند.
- مشکل در این است که تخمین احتمال شرطی $p(w_i/x)$ بطور مستقیم از روی داده خیلی سخت است.

تخمین GMM:

برای GMM در حالت کلی log-likelihood برابر است با:

$$L(\Theta) = \sum_{i=1}^N \log p(x_i | \Theta) = \sum_{i=1}^N \log \left[\sum_{j=1}^K P(w_j) N(x_i | \mu_j, \Sigma_j) \right]$$
$$\Theta = \{P(w_1), \dots, P(w_K), \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$$

همسانی را با کمک رابطه $\frac{\partial L}{\partial \Theta} = 0$ ماکزیمم می‌کنیم ولی باید توجه داشت که حل آن سخت است.

راه‌حل: از الگوریتم ساده شده EM (Expectation-Maximization) استفاده می‌کنیم.

- مرحله E (Expectation) کمک می‌کند تا اطلاعات ناکامل را بازیابی کنیم.

- در مرحله M (Maximization)، همسانی را با استفاده از اطلاعات بازیابی شده ماکزیمم می‌کنیم.

الگوریتم EM:

داده آموزش نا کامل است. (\mathcal{X})

مجموعه داده کامل: $\{\mathcal{X}, \gamma\}$ متغیرهای مخفی: γ

همسانی داده کامل: $L(\theta | \mathcal{X}, \gamma)$ که تابع γ است.

مثال ساده (با اطلاعات کامل):

اگر وقایع "نمرات یک کلاس" باشد:

$w_1 = \text{Gets an A}$

$$P(A) = \frac{1}{2}$$

$w_2 = \text{Gets a B}$

$$P(B) = \mu$$

$w_3 = \text{Gets a C}$

$$P(C) = 2\mu$$

$w_4 = \text{Gets a D}$

$$P(D) = \frac{1}{2} - 3\mu$$

$$(\text{دقت: } \mu \leq 1/6)$$

فرض کنید می‌خواهیم μ را از داده تخمین بزنیم. در یک کلاس مفروض داریم:

a A's

b B's

c C's

d D's

تخمین بیشترین همسانی برای μ با معلوم بودن a, b, c, d چقدر است؟

با ML براحتی حل می شود!

ML = Maximum Likelihood

$$P(A) = \frac{1}{2} \quad P(B) = \mu \quad P(C) = \gamma\mu \quad P(D) = \frac{1}{2} - \gamma\mu$$

$$P(a, b, c, d | \mu) = K \left(\frac{1}{2}\right)^a (\mu)^b (\gamma\mu)^c \left(\frac{1}{2} - \gamma\mu\right)^d$$

$$\log P(a, b, c, d | \mu) = \log K + a \log \frac{1}{2} + b \log \mu + c \log \gamma\mu + d \log \left(\frac{1}{2} - \gamma\mu\right)$$

FOR MLE μ , SET $\frac{\partial \text{LogP}}{\partial \mu} = 0$

$$\frac{\partial \text{LogP}}{\partial \mu} = \frac{b}{\mu} + \frac{\gamma c}{\gamma\mu} - \frac{\gamma d}{1/2 - \gamma\mu} = 0$$

$$\text{Gives max like } \mu = \frac{b + c}{\gamma(b + c + d)}$$

A	B	C	D
۱۴	۶	۹	۱۰

So if class got

$$\text{MLE } \mu = \frac{1}{10}$$

مثال ساده (با اطلاعات مخفی):

Someone tells us that

Number of High grades (A's + B's) = h

Number of C's = c

Number of D's = d

What is the max. like estimate of μ now?

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = \frac{1}{2}\mu$$

$$P(D) = \frac{1}{2} - \frac{1}{2}\mu$$

EXPECTATION

If we know the value of μ we could compute the expected value of a and b

Since the ratio $a:b$ should be the same as the ratio $\frac{1}{2} : \mu$

$$a = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h \quad b = \frac{\mu}{\frac{1}{2} + \mu} h$$

MAXIMIZATION

If we know the expected values of a and b we could compute the maximum likelihood value of μ

$$\mu = \frac{b + c}{2(b + c + d)}$$

کاربرد EM برای مثال بدیهی:

We begin with a guess for μ

We iterate between EXPECTATION and MAXIMALIZATION to improve our estimates of μ and a and b .

Define $\mu(t)$ the estimate of μ on the t 'th iteration

$b(t)$ the estimate of b on t 'th iteration

$\mu(0)$ = initial guess

$$b(t) = \frac{\mu(t)h}{\frac{1}{2} + \mu(t)} = E[b | \mu(t)]$$

$$\mu(t+1) = \frac{b(t) + c}{6(b(t) + c + d)}$$

= max like est of μ given $b(t)$



E-step



M-step

REMEMBER

$$P(A) = \frac{1}{2}$$

$$P(B) = \mu$$

$$P(C) = 2\mu$$

$$P(D) = \frac{1}{2} - 3\mu$$

Continue iterating until converged.

Good news: Converging to local optimum is assured.

Bad news: I said "local" optimum.

همگرایی EM

- Convergence proof based on fact that $\text{Prob}(\text{data} \mid \mu)$ must increase or remain same between each iteration [NOT OBVIOUS]
 - But it can never exceed λ [OBVIOUS]
- So it must therefore converge [OBVIOUS]

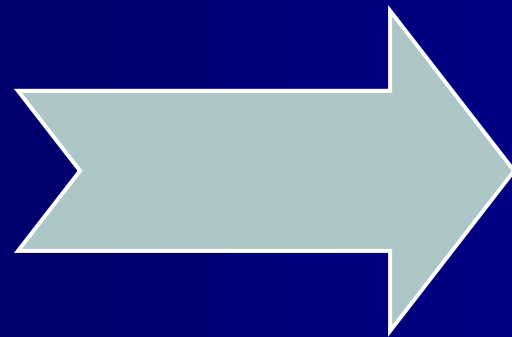
In our example,
suppose we had

$$h = 2.0$$

$$c = 1.0$$

$$d = 1.0$$

$$\mu(\cdot) = \cdot$$



Convergence is generally linear: error decreases by a constant factor each time step.

t	$\mu(t)$	b(t)
0	0	0
1	0.0833	2.857
2	0.0937	3.158
3	0.0947	3.185
4	0.0948	3.187
5	0.0948	3.187
6	0.0948	3.187 ⁽³⁾

الگوریتم EM برای تخمین GMM:

- مرحله E: هر مشاهده بر حسب همسانی هر کدام از گوسیهای متناظر، یک وزن را برای هر خوشه انتساب می‌دهد.

$$P(w_j | x_i, \Theta^{(t)}) = \frac{p(x_i | \mu_j^{(t)}, \Sigma_j^{(t)}) P(w_j)}{\sum_{v=1}^K p(x_i | \mu_v^{(t)}, \Sigma_v^{(t)}) P(w_v)}$$

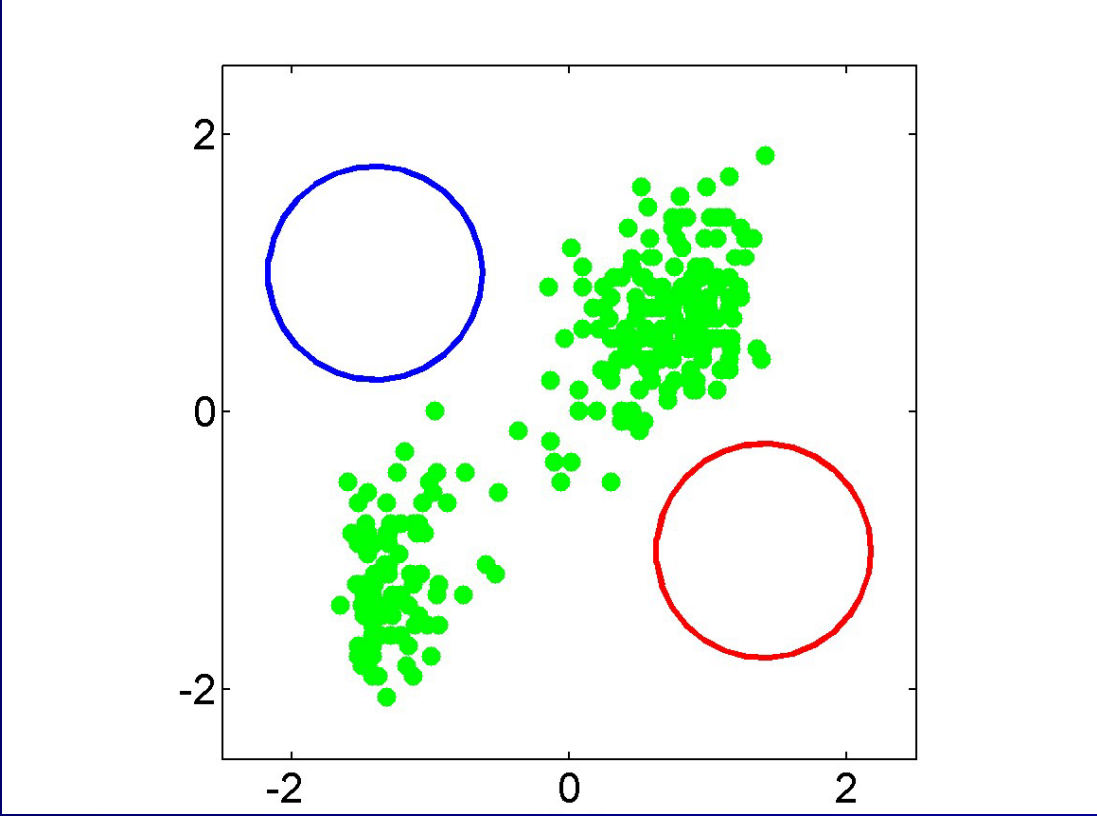
- مرحله M: هر مشاهده با استفاده از ML در میانگینها و کوواریانسهای هر خوشه مشارکت می‌کند.

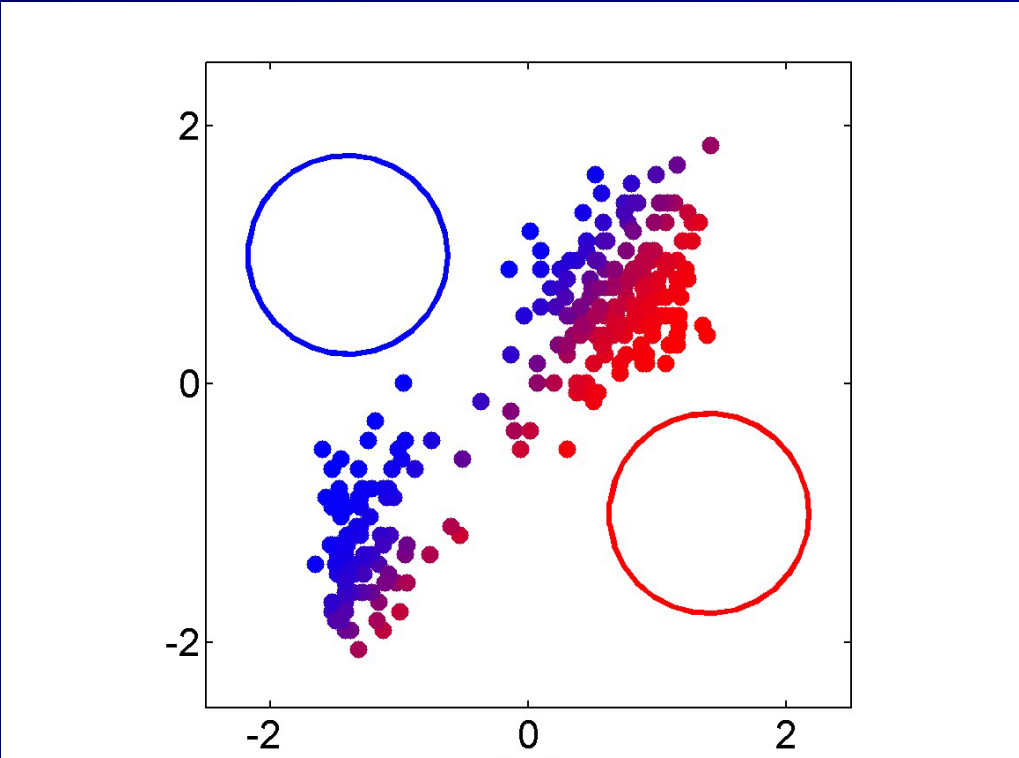
$$\hat{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^N P(w_j | x_i, \Theta^{(t)}) x_i}{\sum_{i=1}^N P(w_j | x_i, \Theta^{(t)})}$$

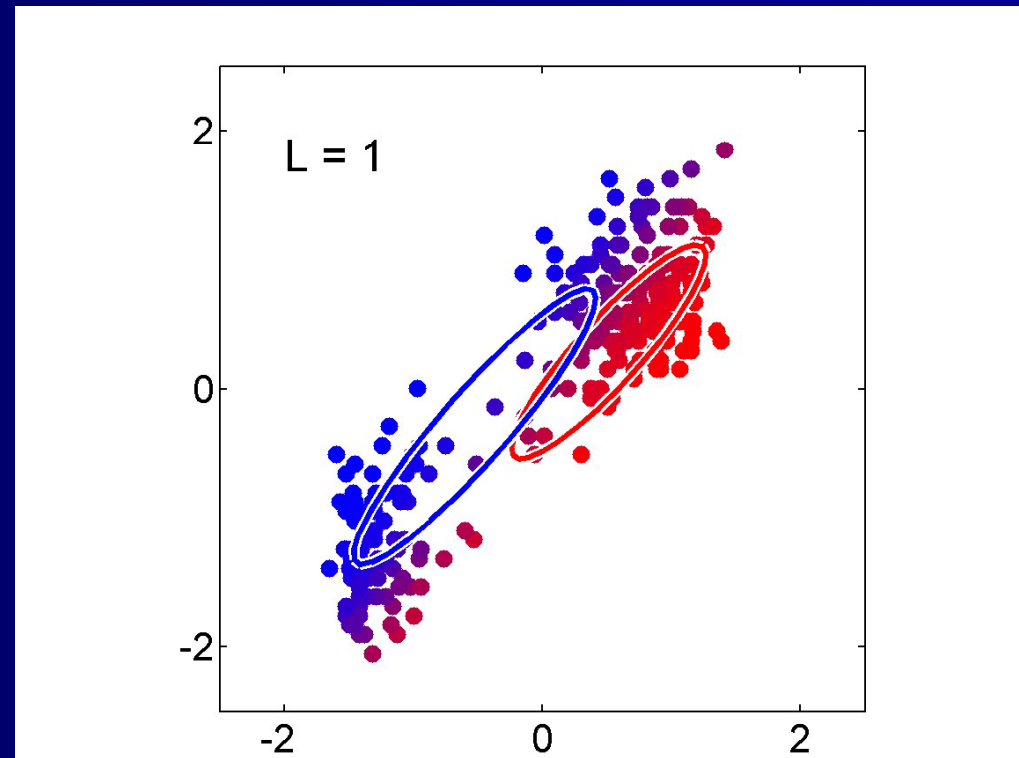
$$\hat{\Sigma}_j^{(t+1)} = \frac{\sum_{i=1}^N P(w_j | x_i, \Theta^{(t)}) (x_i - \hat{\mu}_j^{(t+1)})(x_i - \hat{\mu}_j^{(t+1)})^T}{\sum_{i=1}^N P(w_j | x_i, \Theta^{(t)})}$$

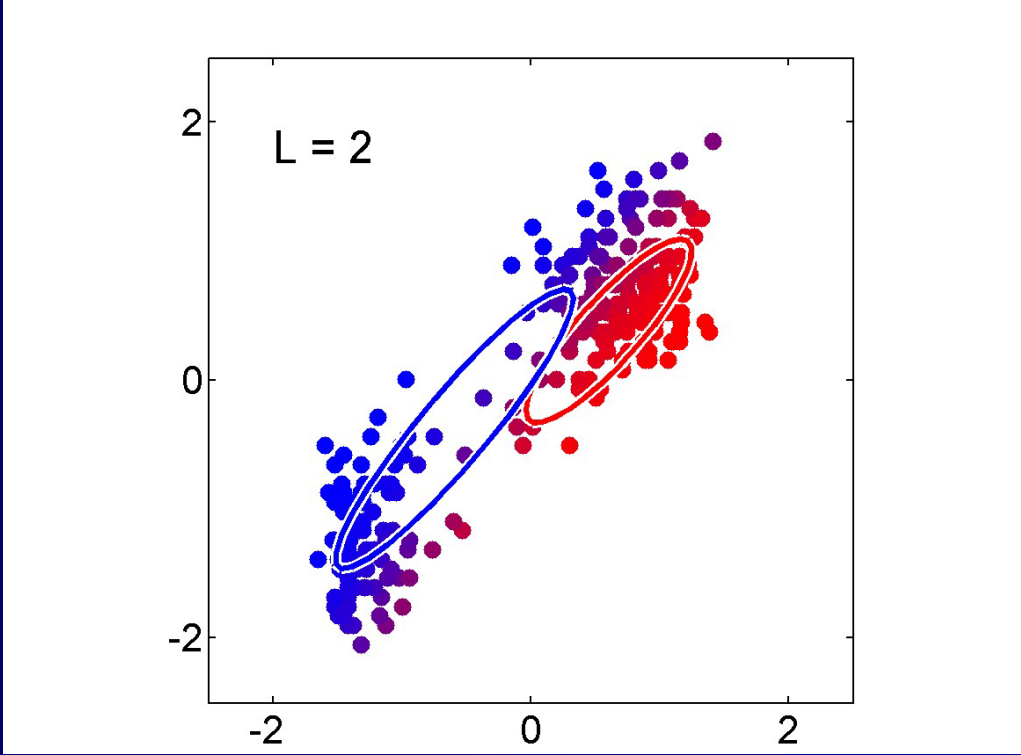
$$P(w_j)^{(t+1)} = \frac{\sum_{i=1}^N P(w_j | x_i, \Theta^{(t)})}{N}$$

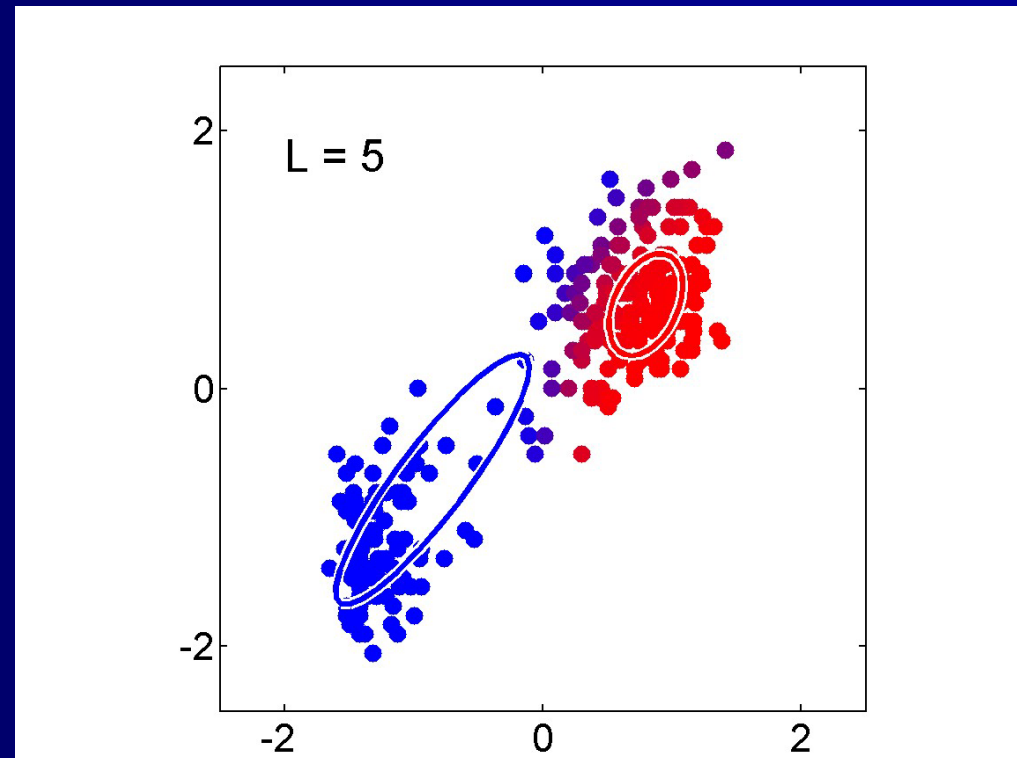
- مراحل E و M را تا همگرایی ادامه می‌دهد.

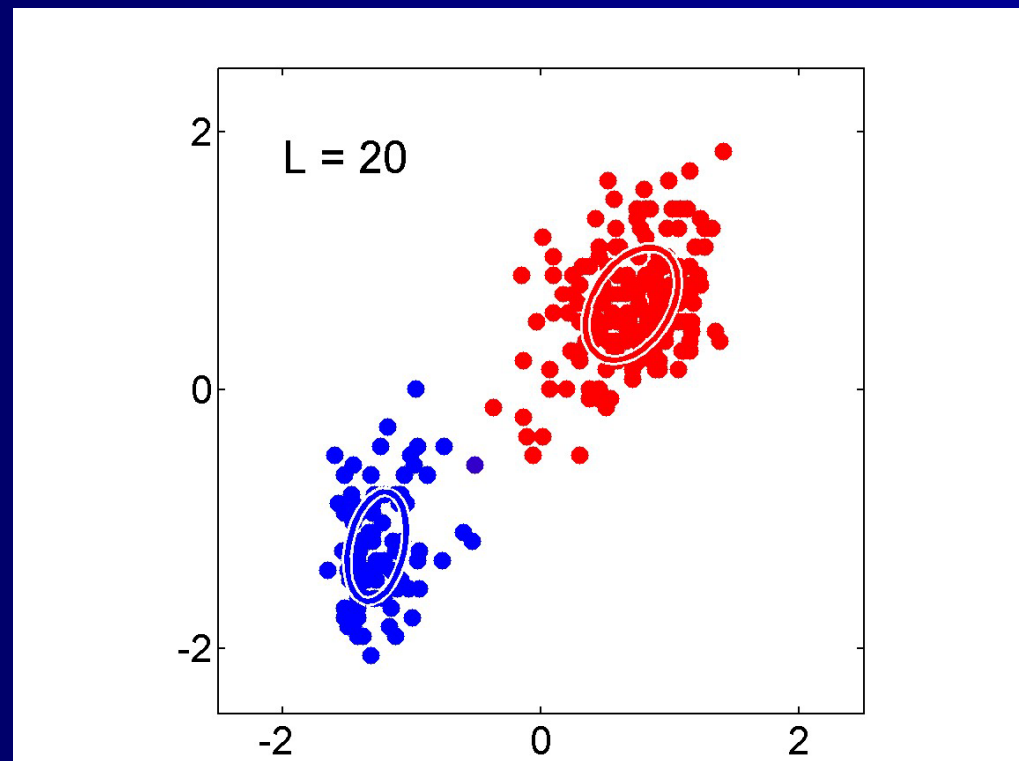












K-means بعنوان حالت خاص EM:

■ با فرض $P(w_j) = 1/K, \Sigma_j = \sigma^2 I$ داریم:

EM \rightarrow K-means

■ الگوریتم K-means:

- ۱- مقداردهی اولیه و تصادفی k مرکز خوشه.
- ۲- هر داده را بر اساس معیار مینیمم فاصله به خوشه‌ها متناظر می‌کنیم.

$$P(w_j | x_i) = \begin{cases} 1 & \text{if } \|x_i - \mu_j\| \leq \|x_i - \mu_v\| \\ 0 & \text{otherwise} \end{cases}$$

- ۳- مراکز خوشه را با استفاده از میانگین مجدداً محاسبه می‌کنیم:

$$\hat{\mu}_j^{(t+1)} = \frac{\sum_{i=1}^N P(w_j | x_i) x_i}{\sum_{i=1}^N P(w_j | x_i)}$$

- ۴- به مرحله ۲ می‌رویم تا تغییری در مراکز خوشه رخ ندهد.

خوشه‌بندی نرم و سخت:

■ اغلب GMM را بعنوان روش خوشه‌بندی نرم می‌دانند.

$$P(w_j | x_i, \Theta)$$

- احتمال پسین اینکه آمین مشاهده از آمین گوسی ایجاد بشود.

- برای هر نقطه x_i ، GMM احتمالات پسین $\{P(w_1 | x_i, \Theta), \dots, P(w_K | x_i, \Theta)\}$

را پیدا می‌کند که برچسب خوشه را برای هر x_i می‌توانیم با کمک رابطه زیر بیابیم:

$$C(x_i) = \arg \max_j P(w_j | x_i, \Theta)$$

■ K-means یک روش خوشه‌بندی سخت است.

- جواب فقط ۰ یا ۱ است.

GMM بعنوان روش خوشه‌بندی نرم:

■ K-Means

R clusters per class

Iteration:

Step ۱: Set each data to
one cluster

$\{0, 1\}$

Step ۲: Tune cluster
centers

■ GMM

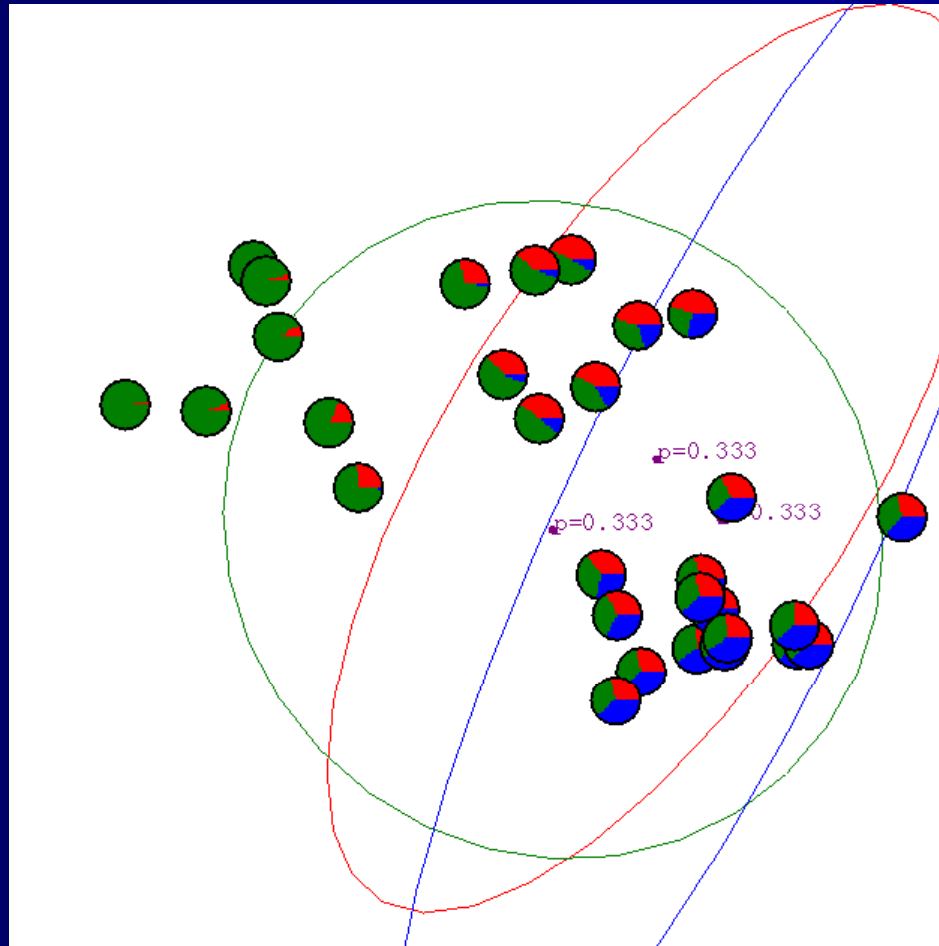
R Gaussians per class

Iteration:

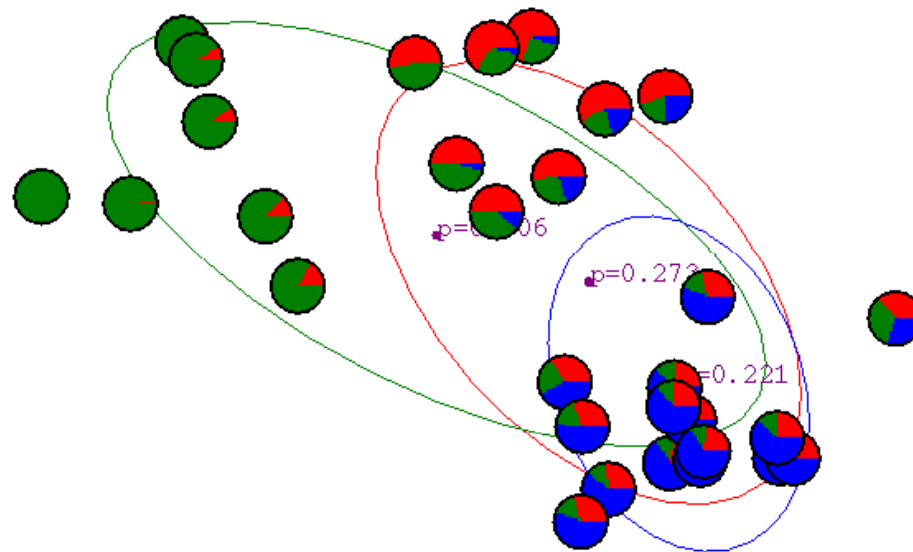
E step: Set each data's
responsibilities for all clusters
 $[0, 1]$

M step: Tune parameters of
Gaussians

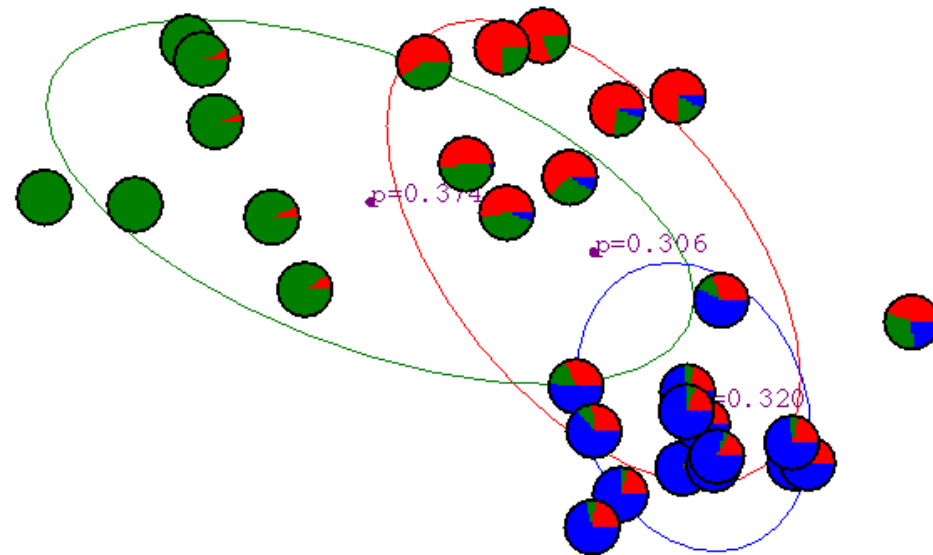
Gaussian Mixture Example: Start



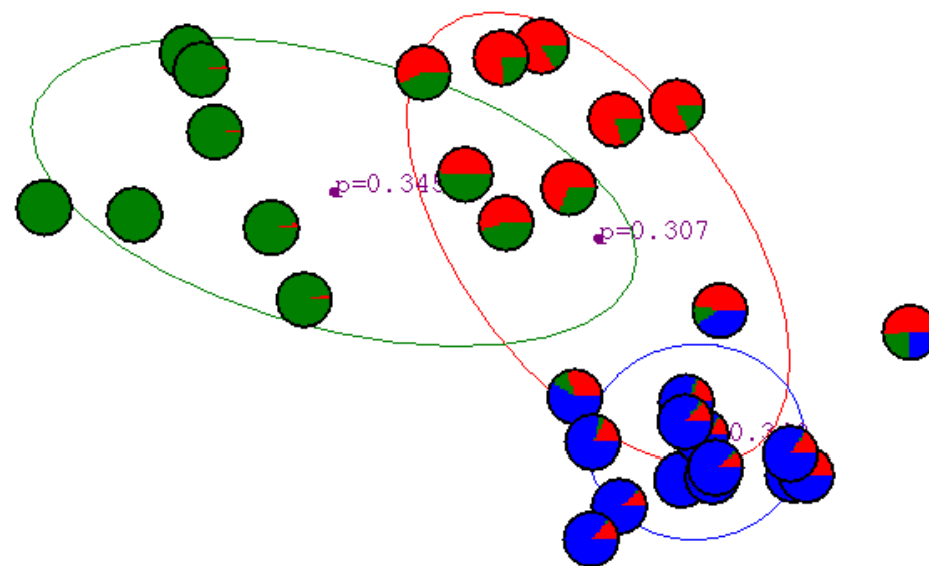
After First Iteration



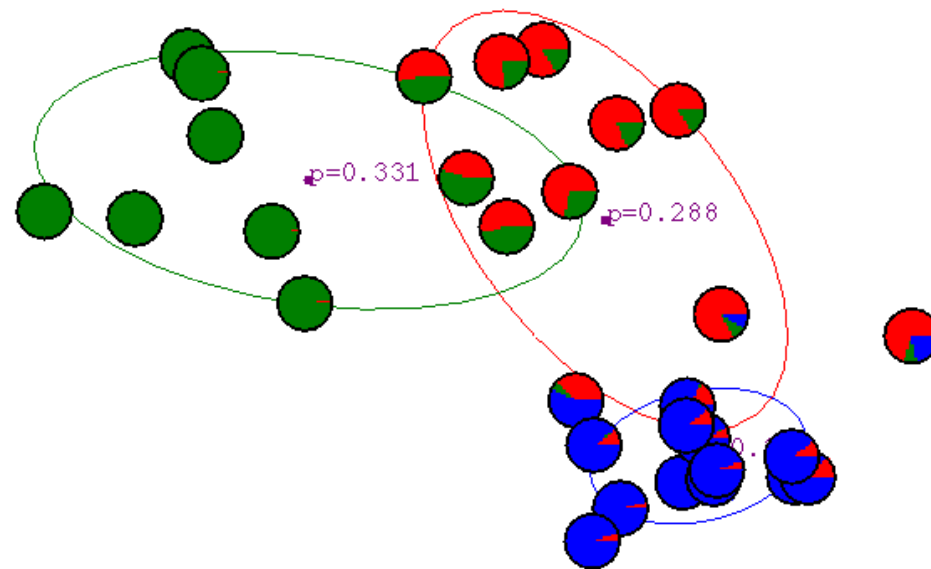
After γ nd Iteration



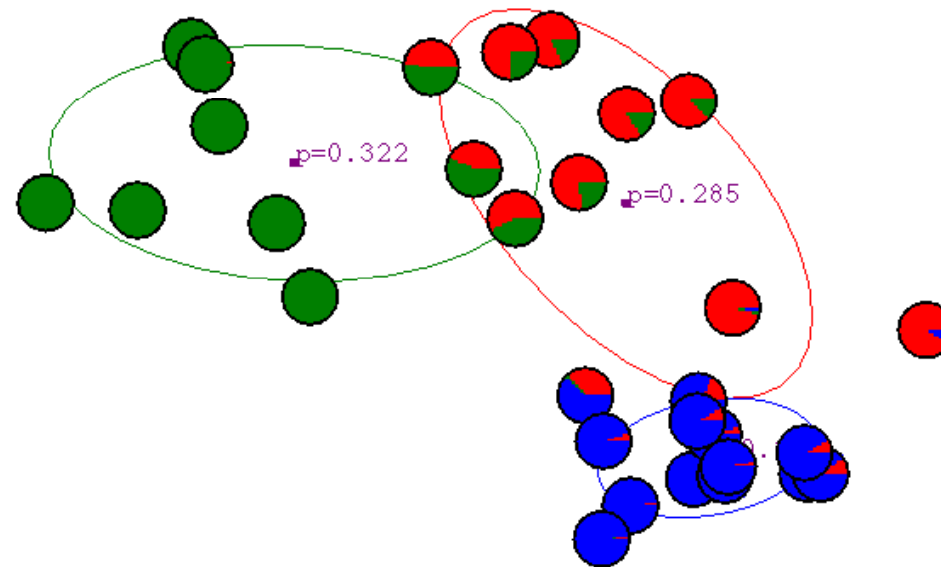
After γ rd Iteration



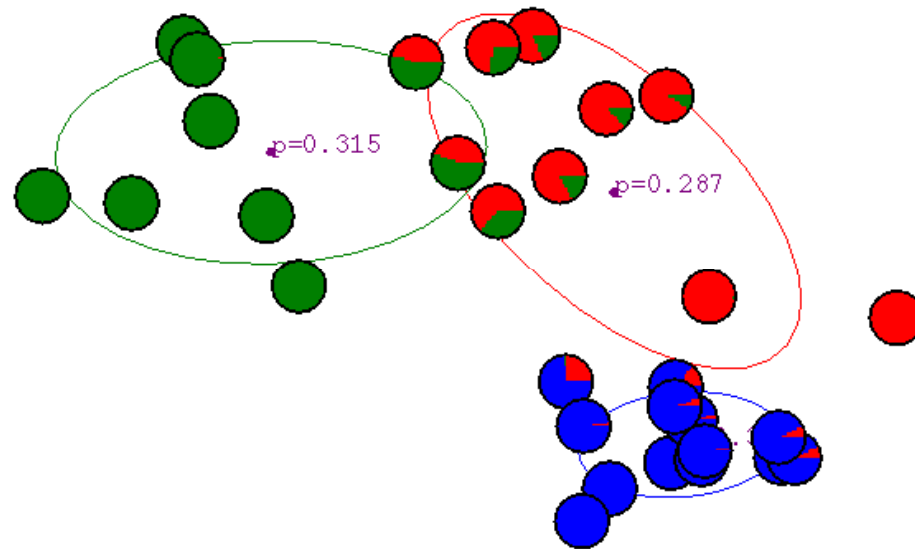
After ν th Iteration



After δ th Iteration



After t th Iteration



After γ -th Iteration

