

COMP90051

Workshop Week 08

About the Workshops

- 7 sessions in total
 - Tue 12:00-13:00 AH211
 - Tue 12:00-13:00 AH108 *
 - Tue 13:00-14:00 AH210
 - Tue 16:15-17:15 AH109
 - Tue 17:15-18:15 AH236 *
 - Tue 18:15-19:15 AH236 *
 - Fri 14:15-15:15 AH211

About the Workshops

- Homepage

- <https://trevorcohn.github.io/comp90051-2017/workshops>

- Solutions will be released on next Friday (a week later).

Syllabus

1	Introduction; Probability theory	Probabilistic models; Parameter fitting	
2	Linear regression; Intro to regularization	Logistic regression; Basis expansion	
3	Optimization; Regularization	Perceptron	
4	Backpropagation	CNNs; Auto-encoders	
5	Hard-margin SVMs	Soft-margin SVMs	
6	Kernel methods	Ensemble Learning	
7	Clustering	EM algorithm	←
8	Dimensionality reduction; Principal component analysis	Multidimensional scaling; Spectral clustering	
9	Bayesian fundamentals	Bayesian inference with conjugate priors	
10	PGMs, fundamentals	Conditional independence	
11	PGMs, inference	Belief propagation	
12	Statistical inference; Apps	Subject review	

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Multivariate Gaussian distribution
 - ❑ Estimate parameters (for 1-d, 2-d Gaussian)
 - ❑ Probabilistic graphical models (PGM)
 - ❑ Parameters and variables
 - ❑ An example: Gaussian mixture models (GMM)
 - ❑ Generative process, joint distribution factorization, plate notation
 - ❑ Expectation maximization (EM) algorithm for GMM
- ❑ IPython notebook task: GMM

Outline

- Review the lecture, background knowledge, etc.
 - **Multivariate Gaussian distribution**
 - Estimate parameters (for 1-d, 2-d Gaussian)
- Probabilistic graphical models (PGM)
 - Parameters and variables
 - An example: Gaussian mixture models (GMM)
 - Generative process, joint distribution factorization, plate notation
- Expectation maximization (EM) algorithm for GMM

- IPython notebook task: GMM

Mean, variance, and covariance

$$\square X_1 = [1, 2, 3, 4, 5], \quad X_2 = [1, 3, 4, 5, 7]$$

$$\square \mu_1 = 3, \quad \mu_2 = 4$$

$$\square \text{Var}(X_1) = \frac{1}{5} [(-2)^2 + (-1)^2 + 0^2 + 1^2 + 2^2] = 2 = \sigma_1^2$$

$$\square \text{Var}(X_2) = \frac{1}{5} [(-3)^2 + (-1)^2 + 0^2 + 1^2 + 3^2] = 4 = \sigma_2^2$$

$$\begin{aligned} \square \text{Cov}(X_1, X_2) &= \frac{1}{5} [(-2)(-3) + (-1)(-1) + 0 \cdot 0 + 1 \cdot 1 + 2 \cdot 3] \\ &= 2.8 = \rho \sigma_1 \sigma_2 \end{aligned}$$

$$\square \text{*Standard deviation } \sigma_i = \sqrt{\text{Var}(X_i)}$$

$$\square \text{*Correlation coefficient } \rho = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1)}\sqrt{\text{Var}(X_2)}} = \frac{\text{Cov}(X_1, X_2)}{\sigma_1 \sigma_2}$$

Multivariate Gaussian distribution

□ Univariate

$$p(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

□ where σ is the standard deviation, σ^2 is the variance

□ Multivariate (k -d)

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})^T}$$

□ where $\boldsymbol{\Sigma}$ is the covariance matrix

Bivariate ($k = 2$)

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

□ 2-d point: $\mathbf{x} = [x_1 \quad x_2]$

□ Parameters:

□ Mean: $\boldsymbol{\mu} = [\mu_1 \quad \mu_2]$

□ Covariance matrix: $\boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$

□ σ_1 and σ_2 are the standard deviations

□ ρ is the correlation coefficient

The covariance matrix Σ and σ_1, σ_2, ρ

$$\square \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

\square For example:

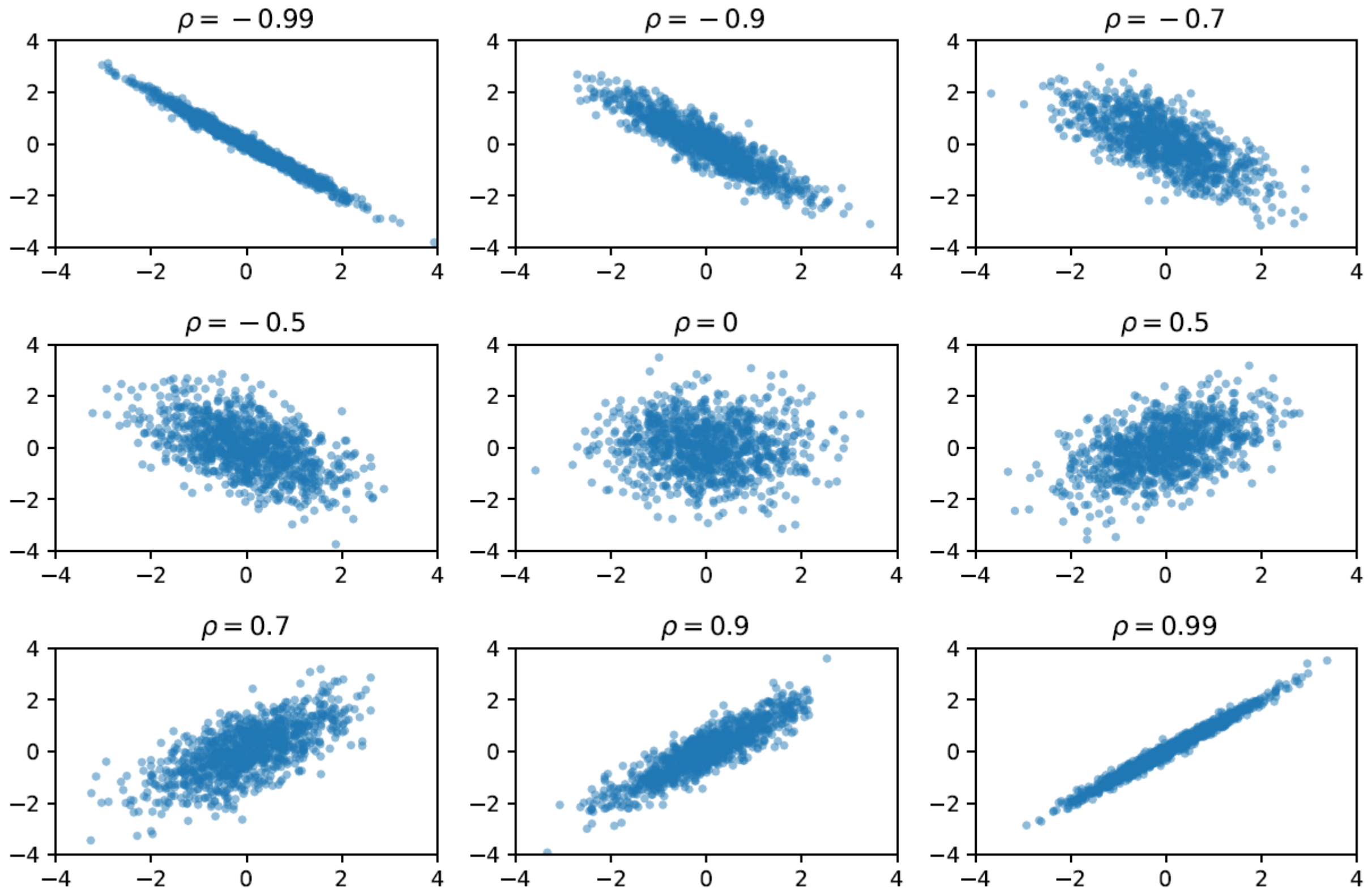
$$\square \Sigma = \begin{bmatrix} 9 & 5 \\ 5 & 4 \end{bmatrix}$$

$$\square \Sigma_{11} = 9, \Sigma_{12} = \Sigma_{21} = 5, \Sigma_{22} = 4$$

$$\square \sigma_1 = \sqrt{\Sigma_{11}} = 3, \sigma_2 = \sqrt{\Sigma_{22}} = 2$$

$$\square \Sigma_{12} = \Sigma_{21} = \rho\sigma_1\sigma_2 \rightarrow \rho = \frac{5}{3 \times 2} = \frac{5}{6}$$

2-d Gaussian with $\sigma_1, \sigma_2 = 1$ & different ρ



Generate 2-d Gaussian points

```
fig = plt.figure(figsize=(9, 6))

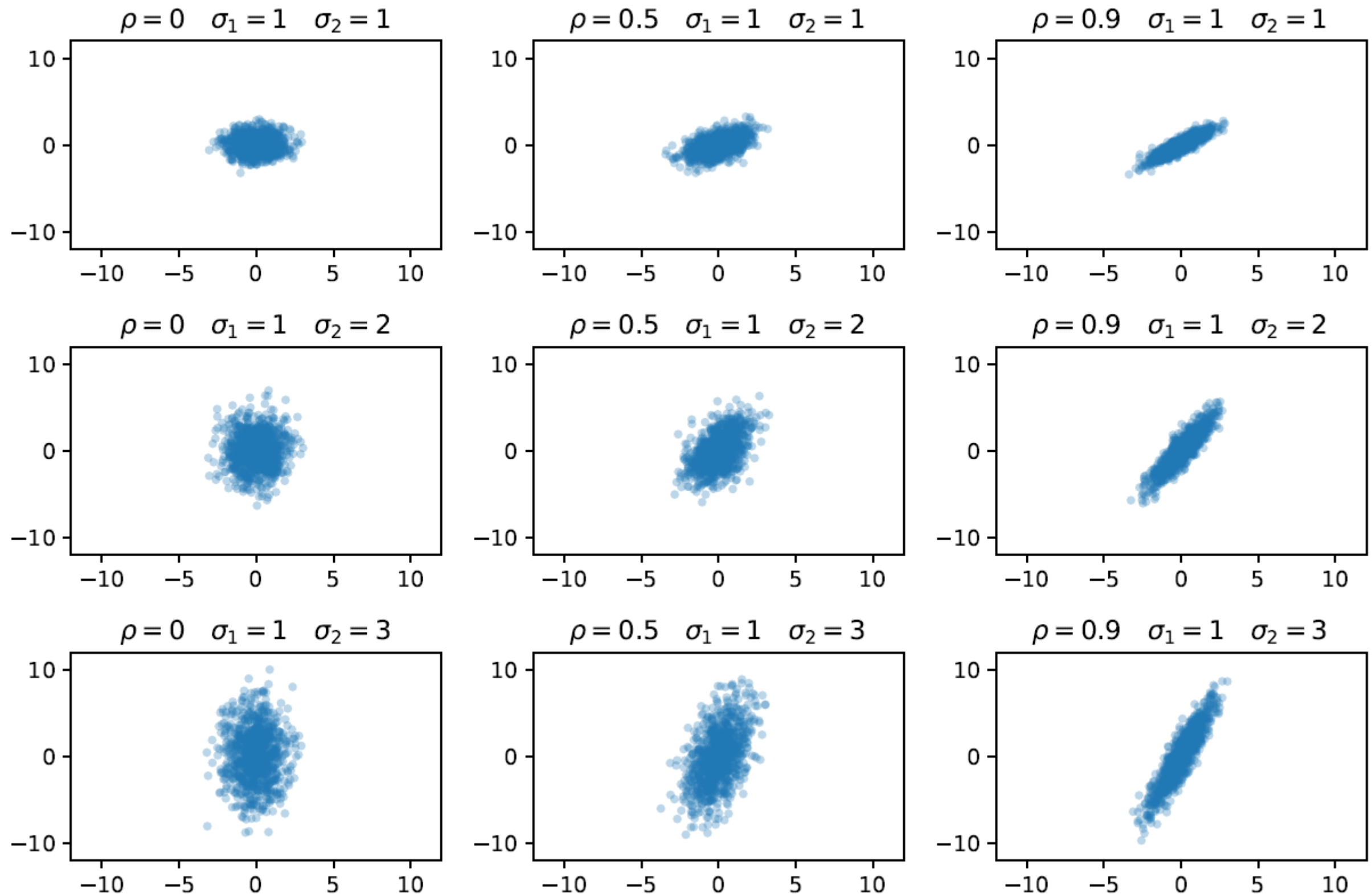
rhos = [-0.99, -0.9, -0.7, -0.5, 0, 0.5, 0.7, 0.9, 0.99]
for i, rho in enumerate(rhos):
    mean = [0, 0]
    cov = [[1, rho], [rho, 1]]

    X1, X2 = np.random.multivariate_normal(mean, cov, 1000).T

    ax = fig.add_subplot(3, 3, i+1)
    ax.plot(X1, X2, '.', alpha=0.5)
    ax.set_title(r'$\rho$=%g'%rho)
    ax.set_xlim([-4, 4])
    ax.set_ylim([-4, 4])

plt.tight_layout()
plt.show()
```

2-d Gaussian with $\sigma_1 = 1$ & different ρ, σ_2



How to estimate parameters for 1-d points

```
mean, std = 3, 2
```

```
X = np.random.normal(mean, std, 20)
```

```
-----
```

```
[1.45592203  3.68408712  3.24950642  2.15079211
 5.68247789  3.44566904  6.05949948  5.15271825
 3.81870043  1.35680408  2.73188372  2.89284301
 5.74163191 -1.77632601  0.61871797  4.13751951
 2.89829252  1.90861471  0.27220436  3.10042466]
```

$\mu = ?$ $\sigma = ?$

Maximum likelihood estimates (1-d points)

mean, std = 3, 2

X = np.random.normal(mean, std, 20)

```
[1.45592203  3.68408712  3.24950642  2.15079211
 5.68247789  3.44566904  6.05949948  5.15271825
 3.81870043  1.35680408  2.73188372  2.89284301
 5.74163191 -1.77632601  0.61871797  4.13751951
 2.89829252  1.90861471  0.27220436  3.10042466]
```

$$L = \prod_{i=1}^{20} N(x_i | \mu, \sigma) = \prod_{i=1}^{20} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(x_i - \mu)^2}$$

□ Find μ and σ that maximizes the likelihood L , let $\frac{\partial L}{\partial \mu} = 0$ and $\frac{\partial L}{\partial \sigma} = 0$

$$\mu = \frac{1}{20} \sum_{i=1}^{20} x_i \quad \sigma^2 = \frac{1}{20} \sum_{i=1}^{20} (x_i - \mu)^2 = \frac{1}{20} (\mathbf{x} - \mu)^T (\mathbf{x} - \mu)$$

Maximum likelihood estimates (1-d points)

```
mean, std = 3, 2
```

```
X = np.random.normal(mean, std, 20)
```

```
-----
```

```
[1.45592203  3.68408712  3.24950642  2.15079211
 5.68247789  3.44566904  6.05949948  5.15271825
 3.81870043  1.35680408  2.73188372  2.89284301
 5.74163191 -1.77632601  0.61871797  4.13751951
 2.89829252  1.90861471  0.27220436  3.10042466]
```

```
mu      = X.sum() / 20
```

```
sigma = np.sqrt((X-mu).dot(X-mu) / 20)
```

```
print('mu      =', mu)
```

```
print('sigma =', sigma)
```

```
-----
```

```
mu      = 2.58609767501
```

```
sigma = 2.17869973073
```

$$\mu = \frac{1}{20} \sum_{i=1}^{20} x_i$$

$$\sigma^2 = \frac{1}{20} (\mathbf{x} - \mu)^T (\mathbf{x} - \mu)$$

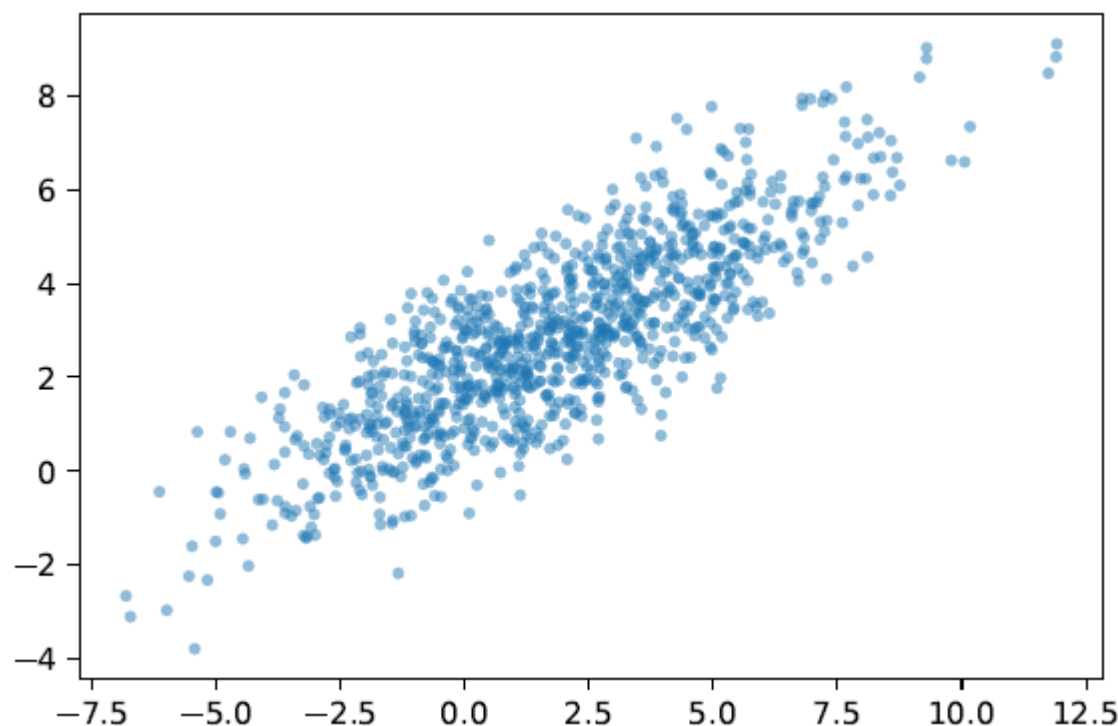
What about MLE for 2-d points?

```
mean = [2, 3]
```

```
cov = [[9, 5], [5, 4]]
```

```
X = np.random.multivariate_normal(mean, cov, 1000)
```

```
plt.plot(X[:, 1], X[:, 2], '.', alpha=0.5)
```



$$\boldsymbol{\mu} = \frac{1}{1000} \begin{bmatrix} \sum_{i=1}^{1000} x_{i,1} & \sum_{i=1}^{1000} x_{i,2} \end{bmatrix}$$
$$\boldsymbol{\Sigma} = \frac{1}{1000} (\mathbf{X} - \boldsymbol{\mu})^T (\mathbf{X} - \boldsymbol{\mu})$$

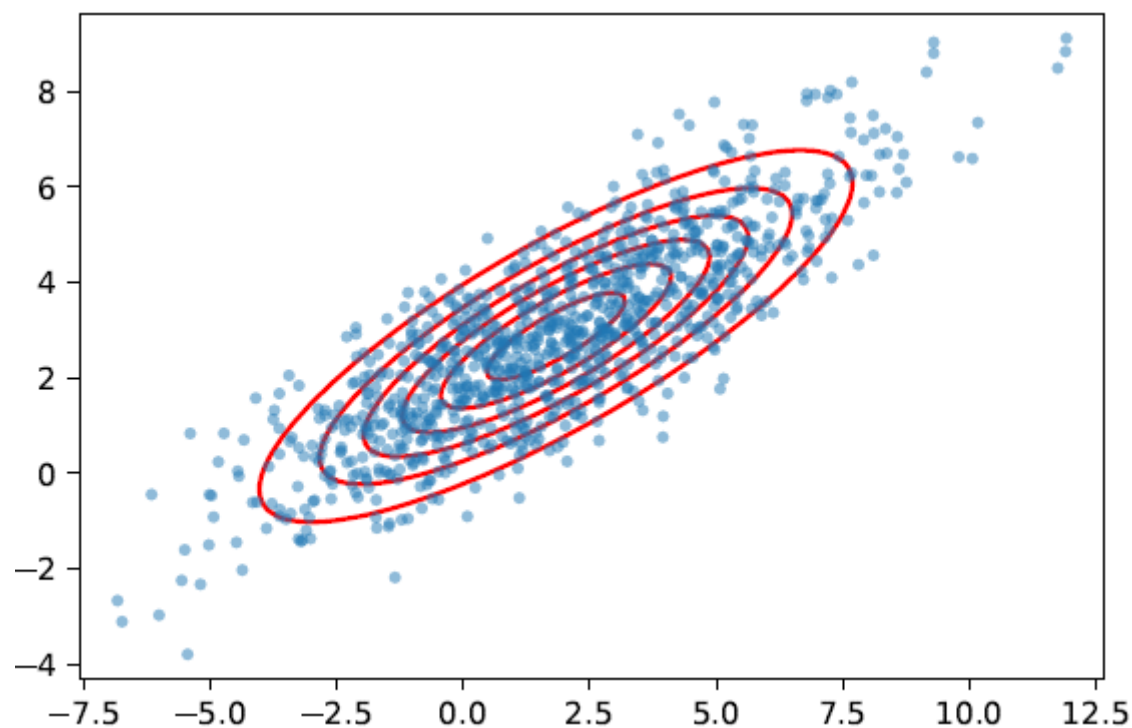
What about MLE for 2-d points?

```
mean = [2, 3]
```

```
cov = [[9, 5], [5, 4]]
```

```
X = np.random.multivariate_normal(mean, cov, 1000)
```

```
plt.plot(X[:, 1], X[:, 2], '.', alpha=0.5)
```



```
Mu = X.mean(axis=0)
```

```
Sigma = (X-Mu).T.dot(X-Mu)/1000
```

```
print('Mu:', Mu)
```

```
print('Sigma:')
```

```
print(Sigma)
```

```
-----  
Mu: [ 1.83023938  2.87127775 ]
```

```
Sigma:
```

```
[[ 9.04242839  4.95851301]
```

```
 [ 4.95851301  3.9911504 ]]
```

Outline

- ❑ Review the lecture, background knowledge, etc.
 - ❑ Multivariate Gaussian distribution
 - ❑ Estimate parameters (for 1-d, 2-d Gaussian)
 - ❑ Probabilistic graphical models (PGM)
 - ❑ Parameters and variables
 - ❑ An example: Gaussian mixture models (GMM)
 - ❑ Generative process, joint distribution factorization, plate notation
 - ❑ Expectation maximization (EM) algorithm for GMM
- ❑ IPython notebook task: GMM

Parameters and variables

- ❑ Note: there are no standard definitions for them, but you can understand them in the following way.
- ❑ Parameters
 - ❑ Known parameters
 - ❑ Unknown parameters
- ❑ (random) Variables
 - ❑ Observable variables
 - ❑ Latent variables

Parameters and variables

- ❑ Note: there are no standard definitions for them, but you can understand them in the following way.
- ❑ Input x and output y are usually considered as variables
- ❑ Variables are drawn from distributions
- ❑ Parameters are fixed, although could be unknown

Parameters and variables

□ Note: there are no standard definitions for them, but you can understand them in the following way.

□ In logistic regression, $p(y = 1|\mathbf{x}; \mathbf{w}, b) = \sigma(\mathbf{w} \cdot \mathbf{x} + b)$

□ \mathbf{w} , b are unknown parameters

□ MLE for \mathbf{w} and b :
$$\max_{\mathbf{w}, b} p(\mathbf{y}|\mathbf{X}; \mathbf{w}, b)$$

□ But we can add a prior distribution for \mathbf{w} : $\mathbf{w} \sim N(0, \lambda^{-1}\mathbf{I})$

□ \mathbf{w} is a random variable, b is still an unknown parameter

□ MAP estimate for \mathbf{w} :
$$\max_{\mathbf{w}, b} p(\mathbf{y}|\mathbf{X}; \mathbf{w}, b)p(\mathbf{w})$$

Three ways to define the GMM

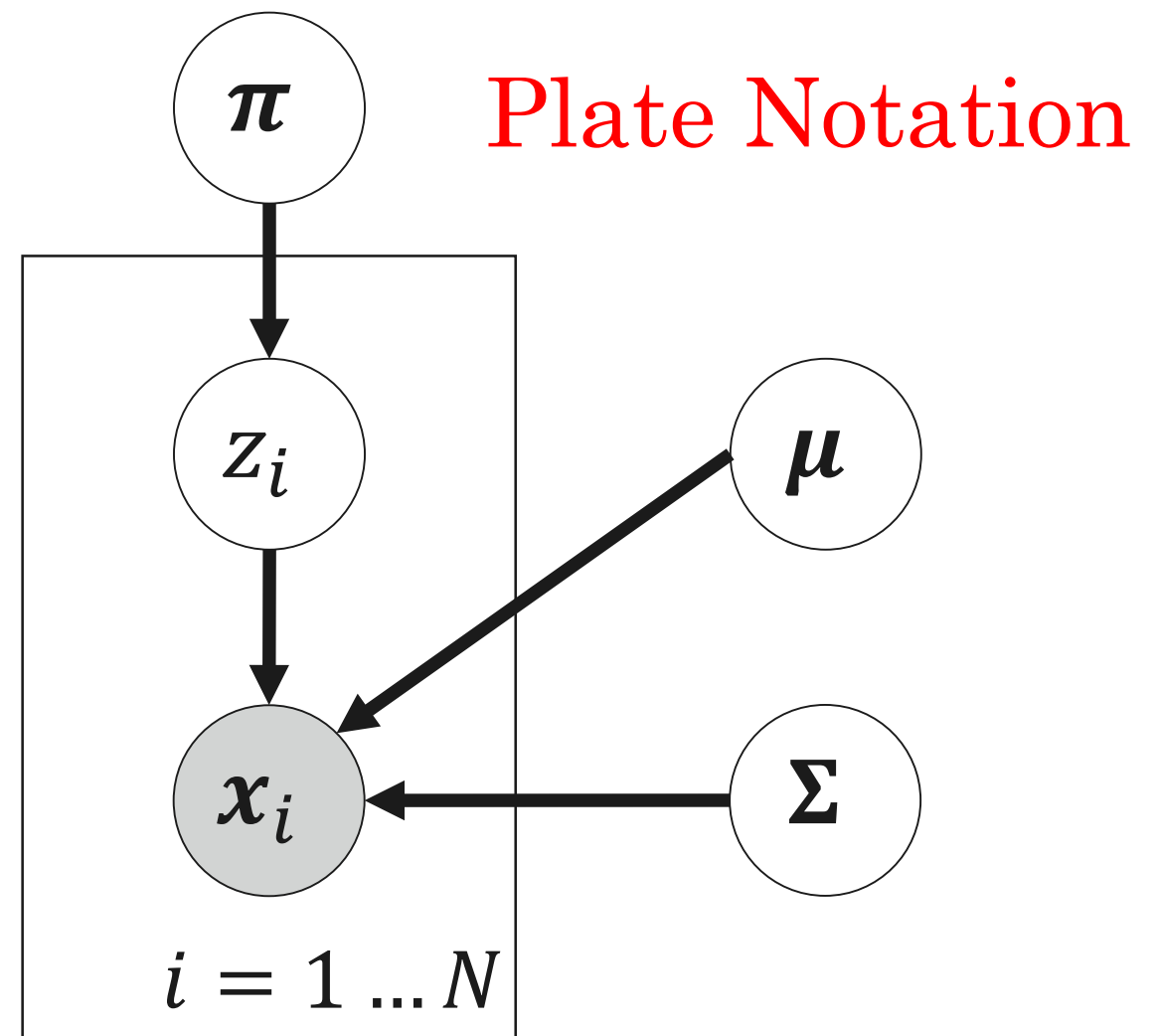
- Generative process

- Parameters: $\boldsymbol{\pi}$, $\boldsymbol{\mu}$, $\boldsymbol{\Sigma}$

- For i in $1 \dots N$:

 - $z_i \sim \text{Categorical}(\boldsymbol{\pi})$

 - $\mathbf{x}_i \sim N(\boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$



- Factorization of the joint distribution

- $$\begin{aligned} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \prod_{i=1}^N p(z_i | \boldsymbol{\pi}) p(\mathbf{x}_i | z_i, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \prod_{i=1}^N \boldsymbol{\pi}_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) \end{aligned}$$

To solve a GMM

- To maximize the log-likelihood

$$\max_{\pi, \mu, \Sigma} \log p(X | \pi, \mu, \Sigma)$$

- But how to calculate the log-likelihood?

To solve a GMM

- To maximize the log-likelihood

$$\max_{\pi, \mu, \Sigma} \log p(X|\pi, \mu, \Sigma)$$

- But how to calculate the log-likelihood?

- – by marginalization

$$p(X|\pi, \mu, \Sigma) = \sum_Z p(X, Z|\pi, \mu, \Sigma)$$

□ We know how to calculate the joint distribution

$$p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^N \pi_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$$

□ So the likelihood can be calculated as

$$\begin{aligned} \log p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \log \sum_{\mathbf{Z}} p(\mathbf{X}, \mathbf{Z} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \\ &= \log \sum_{\mathbf{Z}} \prod_{i=1}^N \pi_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) = \log \prod_{i=1}^N \sum_{z_i=1}^C \pi_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) \\ &= \sum_{i=1}^N \log \sum_{z_i=1}^C \pi_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i}) \end{aligned}$$

To solve a GMM

- To maximize the log-likelihood

$$\max_{\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}} \log p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^N \log \sum_{z_i=1}^C \boldsymbol{\pi}_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$$

- Two options to solve
 - Gradient-based algorithms
 - Expectation maximization (EM) algorithm

An example: 2-D, 3 clusters, 4 points

$$\square \boldsymbol{\pi} = [\pi_1 \quad \pi_2 \quad \pi_3]$$

$$\square \boldsymbol{\mu}_1 = [\mu_{1,1} \quad \mu_{1,2}] \quad \boldsymbol{\mu}_2 = [\mu_{2,1} \quad \mu_{2,2}] \quad \boldsymbol{\mu}_3 = [\mu_{3,1} \quad \mu_{3,2}]$$

$$\square \boldsymbol{\Sigma}_1 = \begin{bmatrix} \Sigma_{1,11} & \Sigma_{1,12} \\ \Sigma_{1,21} & \Sigma_{1,22} \end{bmatrix} \quad \boldsymbol{\Sigma}_2 = \begin{bmatrix} \Sigma_{2,11} & \Sigma_{2,12} \\ \Sigma_{2,21} & \Sigma_{2,22} \end{bmatrix} \quad \boldsymbol{\Sigma}_3 = \begin{bmatrix} \Sigma_{3,11} & \Sigma_{3,12} \\ \Sigma_{3,21} & \Sigma_{3,22} \end{bmatrix}$$

$$\square \boldsymbol{x}_1 = [x_{1,1} \quad x_{1,2}] \quad \boldsymbol{x}_2 = [x_{2,1} \quad x_{2,2}]$$

$$\square \boldsymbol{x}_3 = [x_{3,1} \quad x_{3,2}] \quad \boldsymbol{x}_4 = [x_{4,1} \quad x_{4,2}]$$

$$\log p(X|\boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^N \log \sum_{z_i=1}^C \pi_{z_i} N(\mathbf{x}_i | \boldsymbol{\mu}_{z_i}, \boldsymbol{\Sigma}_{z_i})$$

$$\begin{aligned} &= \log(\pi_1 N(\mathbf{x}_1 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \pi_2 N(\mathbf{x}_1 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + \pi_3 N(\mathbf{x}_1 | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)) \\ &+ \log(\pi_1 N(\mathbf{x}_2 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \pi_2 N(\mathbf{x}_2 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + \pi_3 N(\mathbf{x}_2 | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)) \\ &+ \log(\pi_1 N(\mathbf{x}_3 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \pi_2 N(\mathbf{x}_3 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + \pi_3 N(\mathbf{x}_3 | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)) \\ &+ \log(\pi_1 N(\mathbf{x}_4 | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) + \pi_2 N(\mathbf{x}_4 | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) + \pi_3 N(\mathbf{x}_4 | \boldsymbol{\mu}_3, \boldsymbol{\Sigma}_3)) \end{aligned}$$

□ where

$$\begin{aligned} N(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) &= \frac{1}{2\pi \sqrt{|\boldsymbol{\Sigma}_k|}} e^{-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)} \\ &= \frac{1}{2\pi \sqrt{\Sigma_{k,11}\Sigma_{k,22} - \Sigma_{k,12}\Sigma_{k,21}}} e^{-\frac{[x_{i,1} - \mu_{k,1} \quad x_{i,1} - \mu_{k,2}] \begin{bmatrix} \Sigma_{k,22} & -\Sigma_{k,12} \\ -\Sigma_{k,21} & \Sigma_{k,11} \end{bmatrix} \begin{bmatrix} x_{i,1} - \mu_{k,1} \\ x_{i,1} - \mu_{k,2} \end{bmatrix}}{2(\Sigma_{k,11}\Sigma_{k,22} - \Sigma_{k,12}\Sigma_{k,21})}} \end{aligned}$$

To solve a GMM

- To maximize the log-likelihood

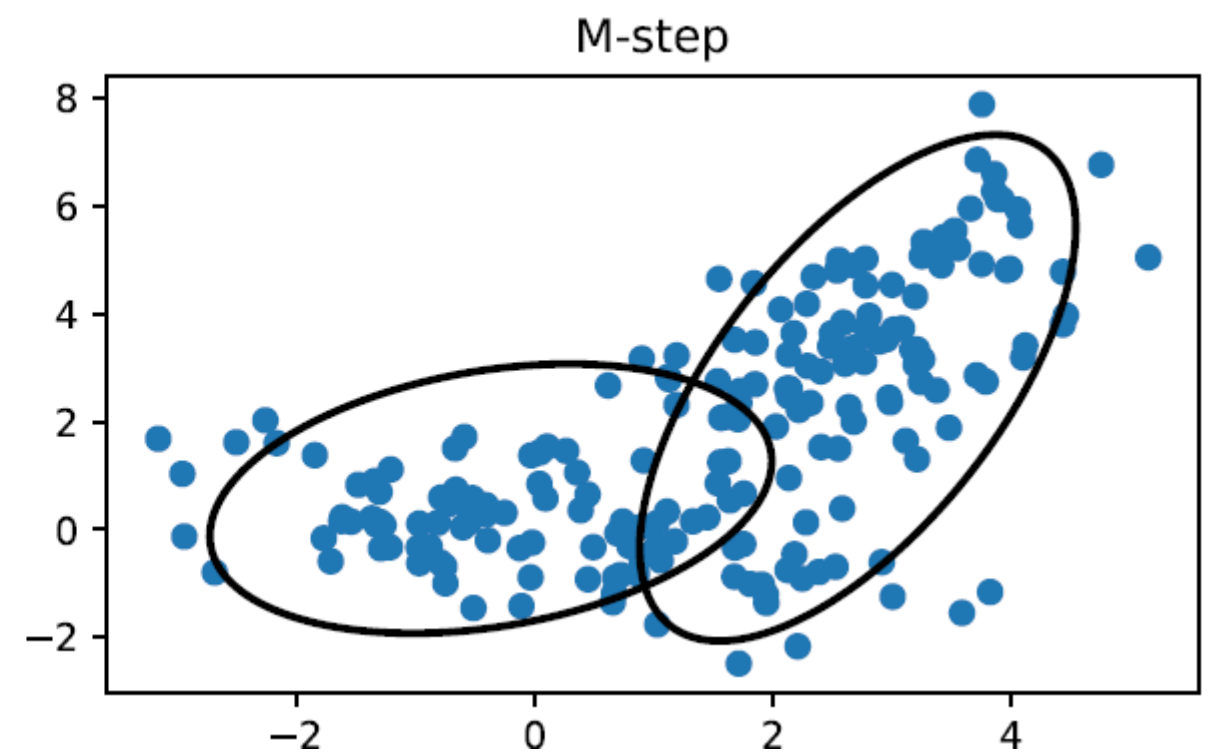
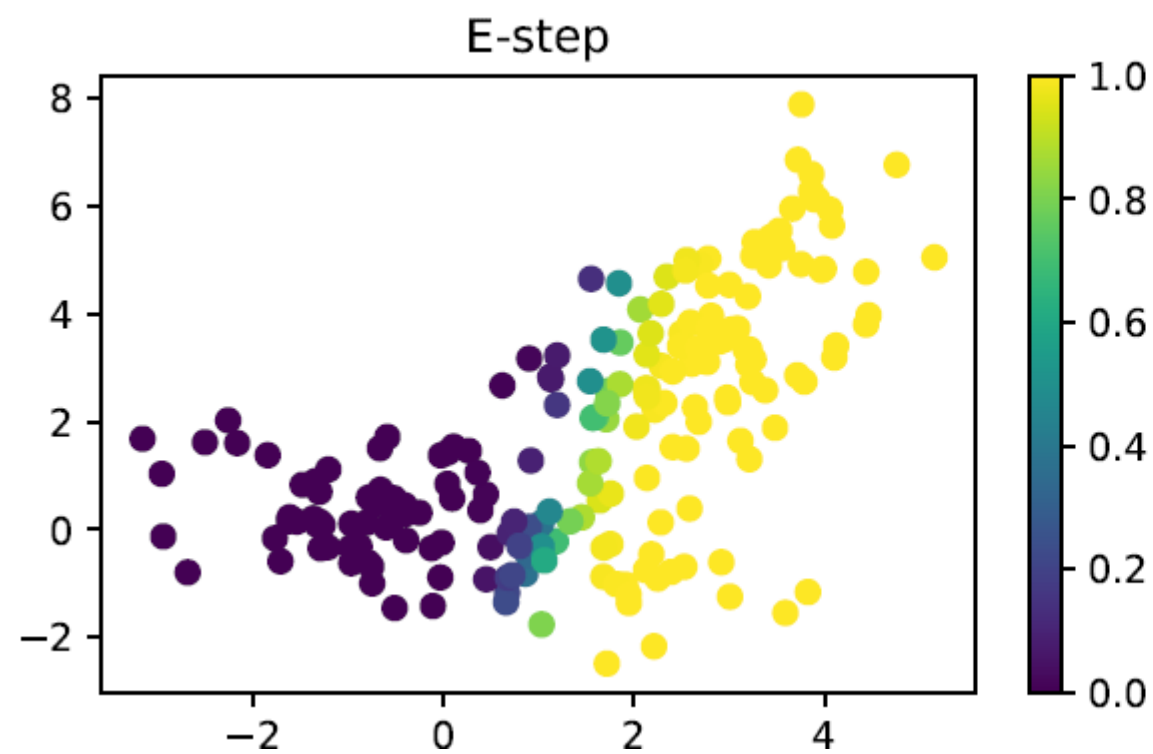
$$\max_{\pi, \mu, \Sigma} \log p(X | \pi, \mu, \Sigma) = \sum_{i=1}^N \log \sum_{z_i=1}^C \pi_{z_i} N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

- Two options to solve

- ~~□ Gradient-based algorithms~~

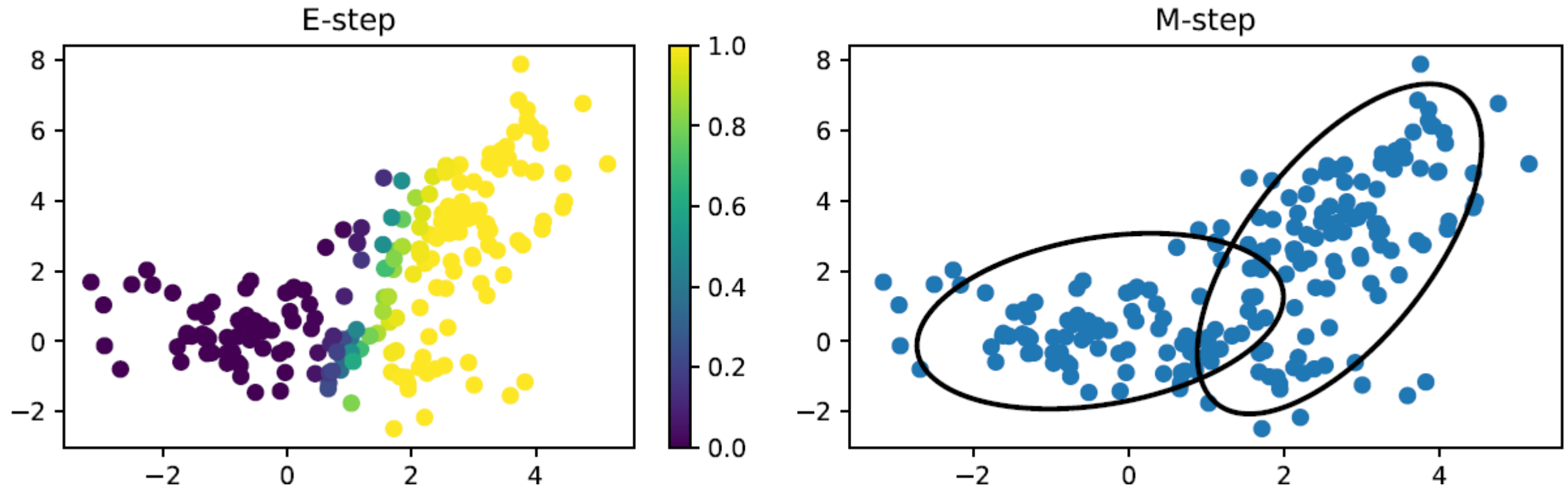
- Expectation maximization (EM) algorithm ←

EM for GMM



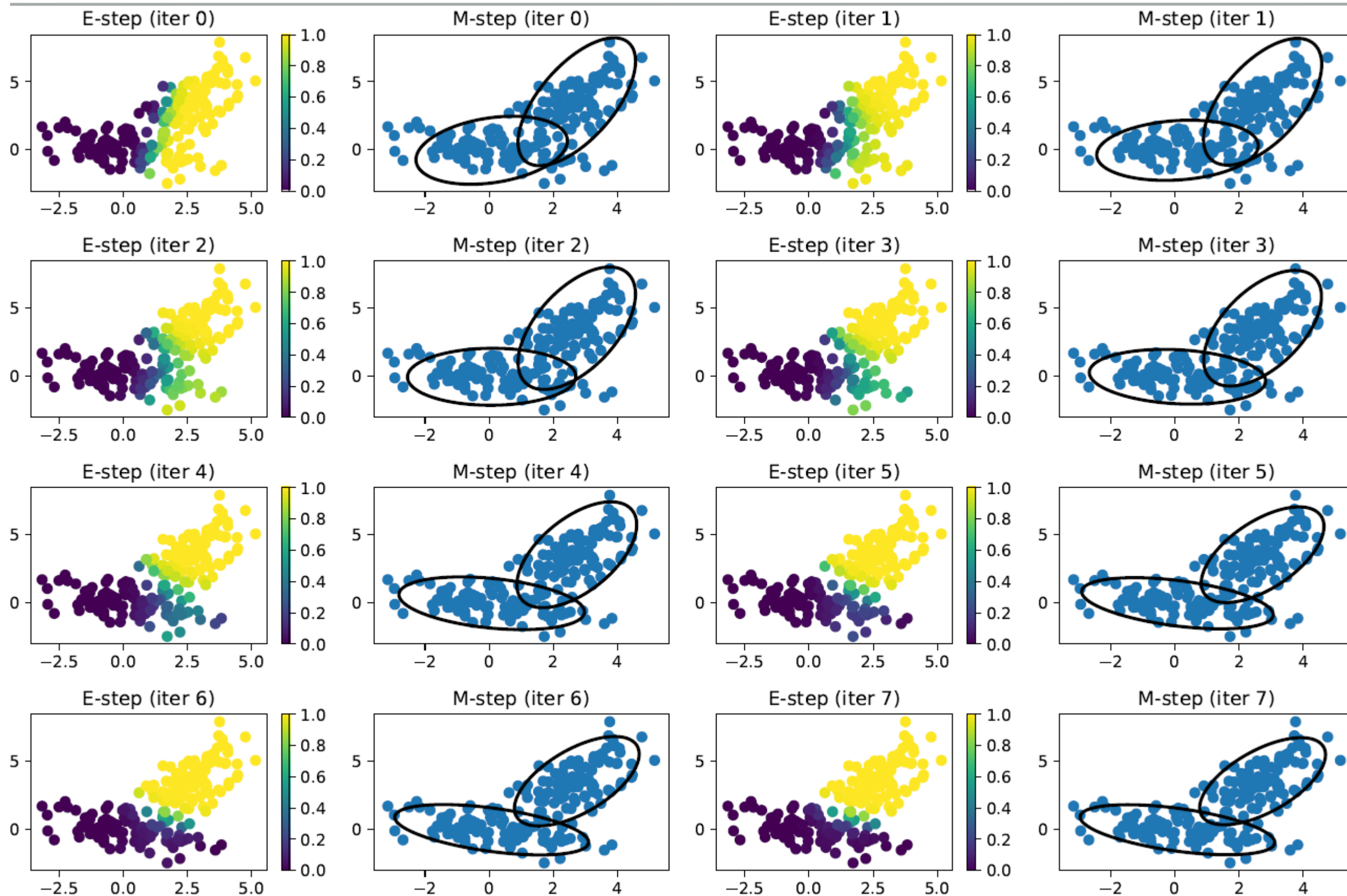
EM for GMM v.s. k-means

- EM for GMM considers the probabilities of a point belonging to different clusters
- K-means assigns every point to one cluster



- EM for GMM estimates μ , Σ , π
- K-means only estimates μ

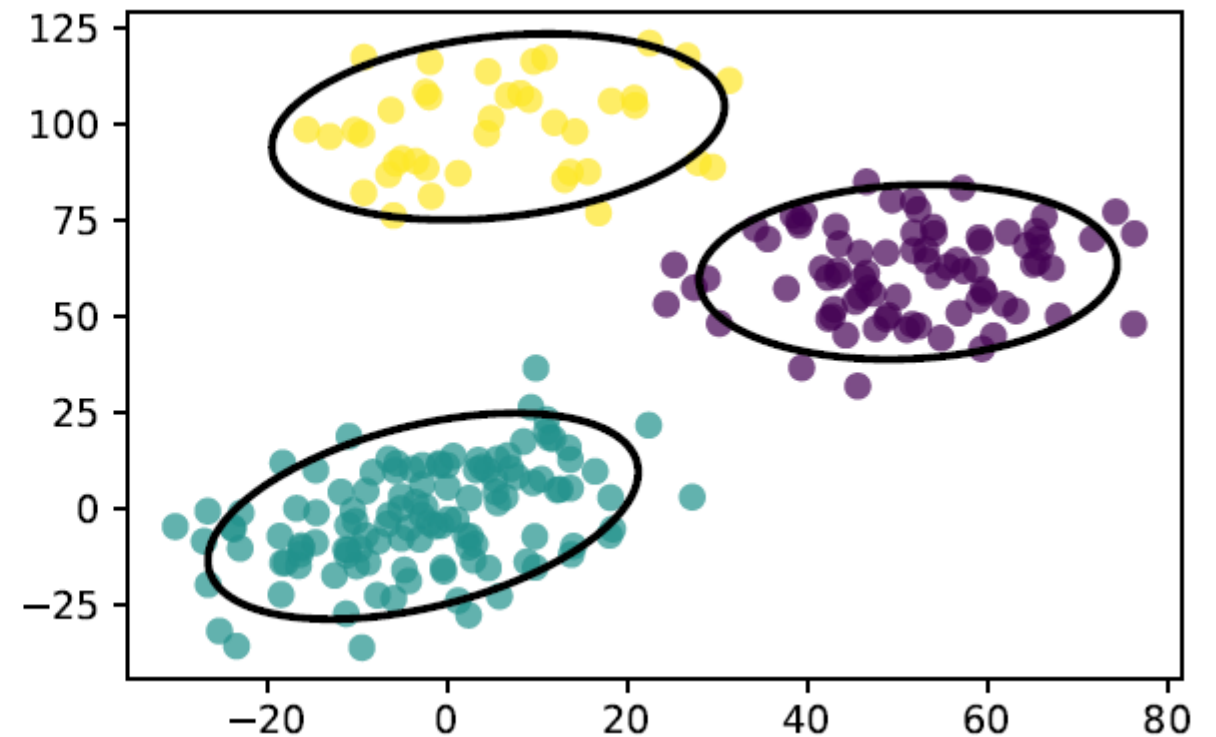
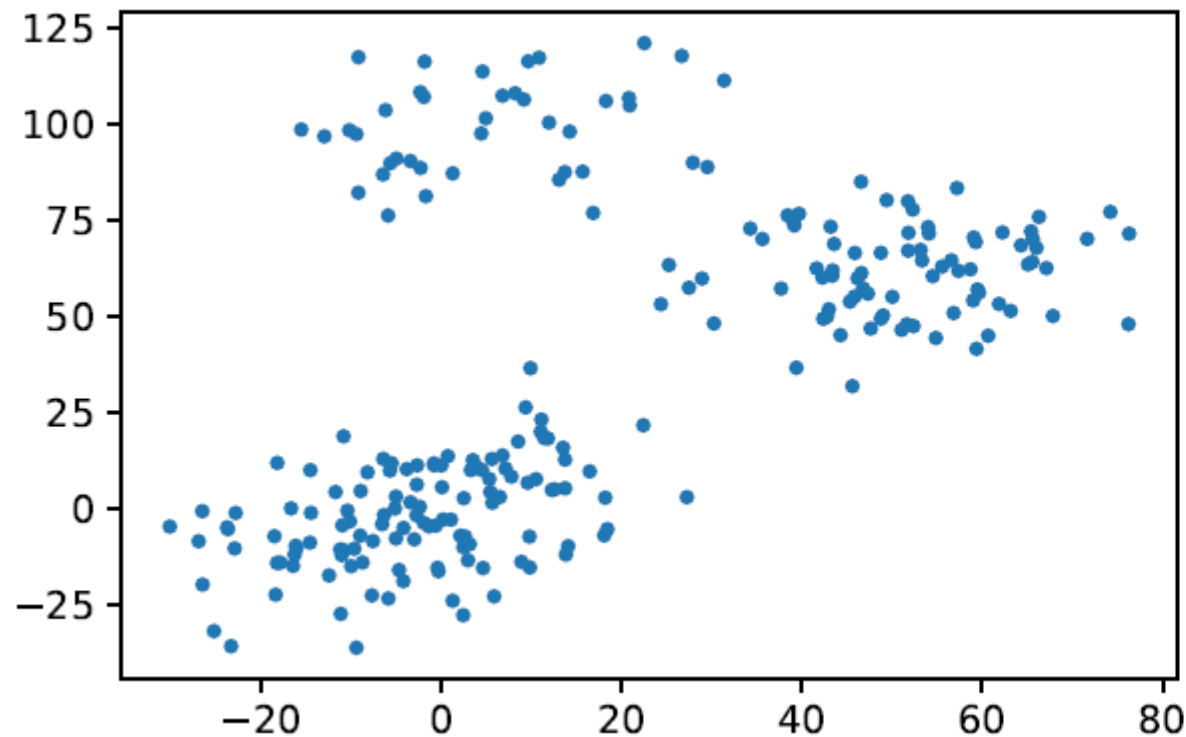
EM for GMM



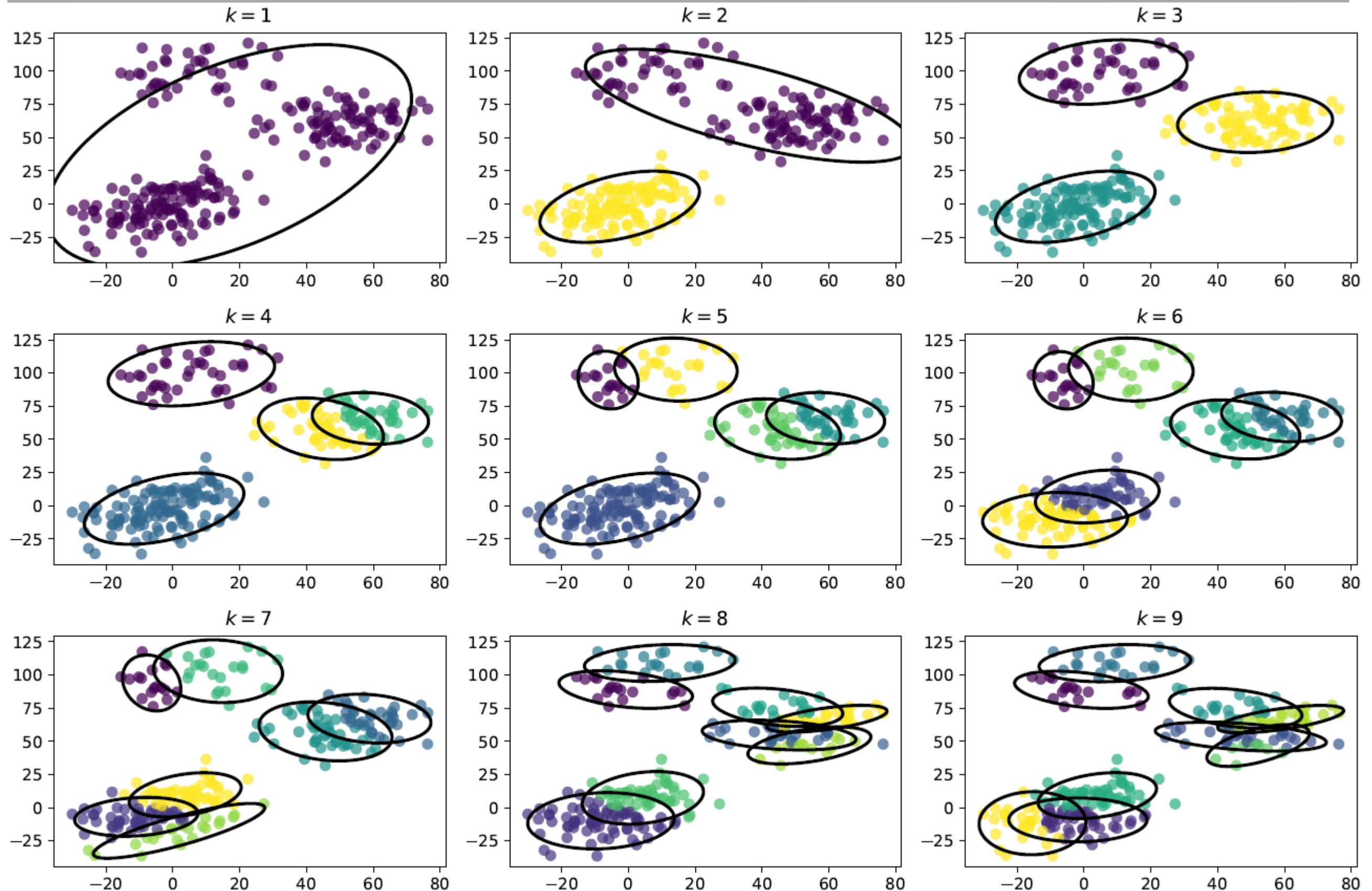
Outline

- Review the lecture, background knowledge, etc.
 - Multivariate Gaussian distribution
 - Estimate parameters (for 1-d, 2-d Gaussian)
 - Probabilistic graphical models (PGM)
 - Parameters and variables
 - An example: Gaussian mixture models (GMM)
 - Generative process, joint distribution factorization, plate notation
 - Expectation maximization (EM) algorithm for GMM
- IPython notebook task: GMM

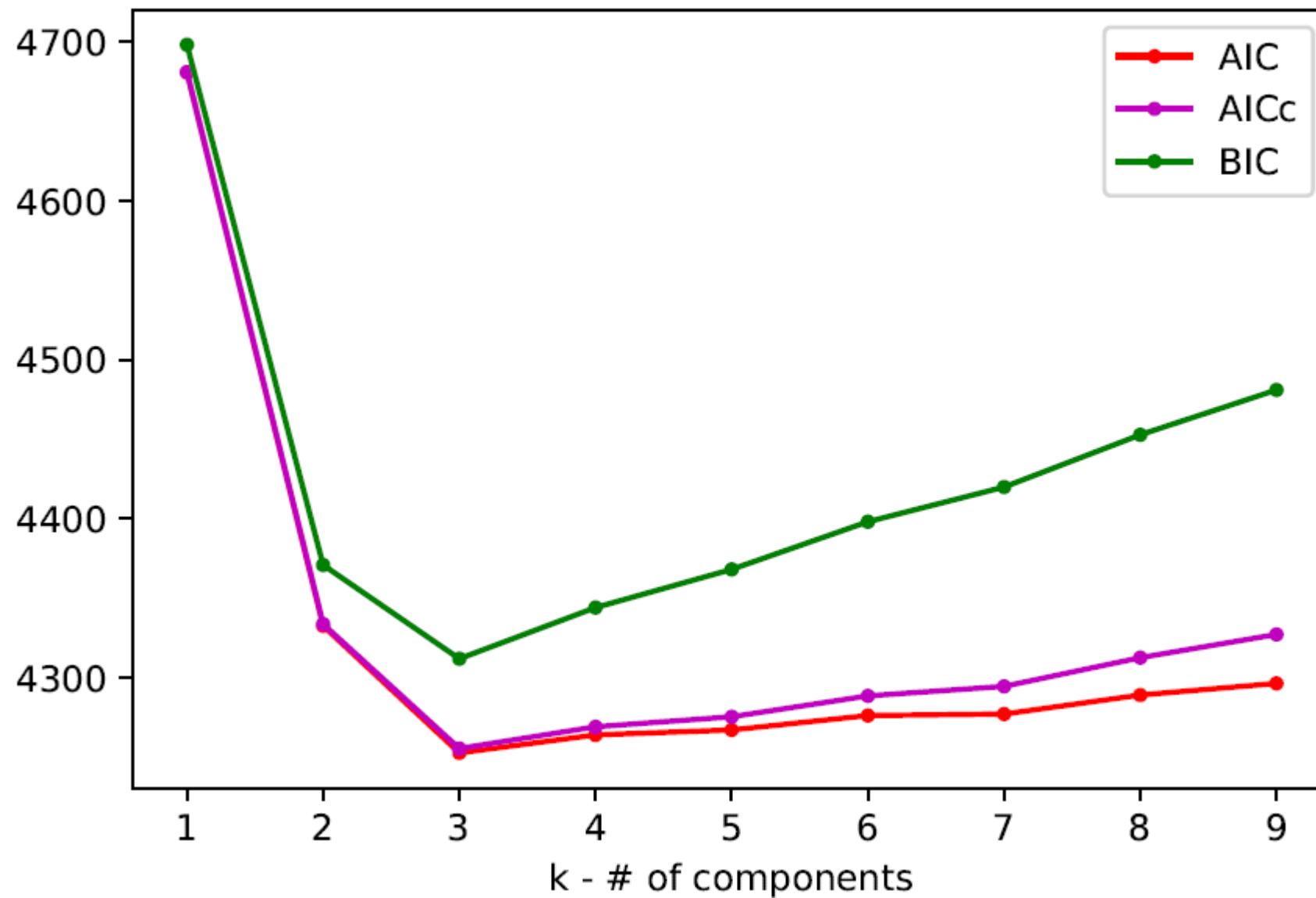
GMM – clustering (EM, $k = 3$)



GMM – clustering with different k



GMM – choose k by AIC, AICc, BIC



GMM – EM generates different solutions

