

Generative Learning Algorithms

Teeradaj Racharak (ເອັກຊ້)
r.teeradaj@gmail.com



Generative vs. Discriminative

So far, the methods we have tried attempt to learn $p(y|x)$ directly.
i.e. Given an input x , we are **trying to map directly** to the output y .

Any algorithm that does this is called a **discriminative** learning algorithm.

Another class of algorithms instead tried to **model from** $p(x|y)$ and $p(y)$.
Such methods are called **generative** learning algorithms.

Using Bayes' rule

If we can build a model of $p(x | y)$ and $p(y)$, then Bayes' rule tells us that:

$$p(y | x) = \frac{p(x | y) \cdot p(y)}{p(x)}$$

Why don't we care about $p(x)$? Let's see ...

Using Bayes' rule

If we can build a model of $p(x|y)$ and $p(y)$, then Bayes' rule tells us that:

$$p(y|x) = \frac{p(x|y) \cdot p(y)}{p(x)}$$

Why don't we care about $p(x)$? Let's see ...

First, generative models are most often used for classification, not regression !

In classification, y is discrete, so we have:

$$p(x) = \sum_i p(x|y=y_i) \cdot p(y=y_i)$$

If y is continuous, we just use an integral instead of sum.

This shows that if we can model $p(x|y)$ and $p(y)$, we can obtain $p(y|x)$ without explicitly calculating $p(x)$.

Using Bayes' rule

We could calculate $p(x)$ directly.

But usually, we want to know which y maximizes $p(y|x)$.

เราต้องการทำ maximization problem

In this case, all we need to do is find:

$$\begin{aligned}y^* &= \arg \max_y p(y|x) \\&= \arg \max_y \frac{p(x|y) \cdot p(y)}{p(x)} \\&= \arg \max_y p(x|y) \cdot p(y)\end{aligned}$$

So, to perform classification in the generative approach,
all we need is models for $p(x|y)$ and $p(y)$.

class ไหนให้ตัวเลขสูงสุด ตอบ class นั้น

Gaussian Discriminant Analysis (GDA) Model

Gaussian Discriminant Analysis

Let's consider the case where $p(\mathbf{x} | y)$ is a **multivariate Gaussian**.

A possible example would be $\mathcal{Y} = \{\text{male}, \text{female}\}$ and $\mathcal{X} = \mathbb{R}^2$, where the features are a person's height and weight.

By the way, what is multivariate Gaussian ?

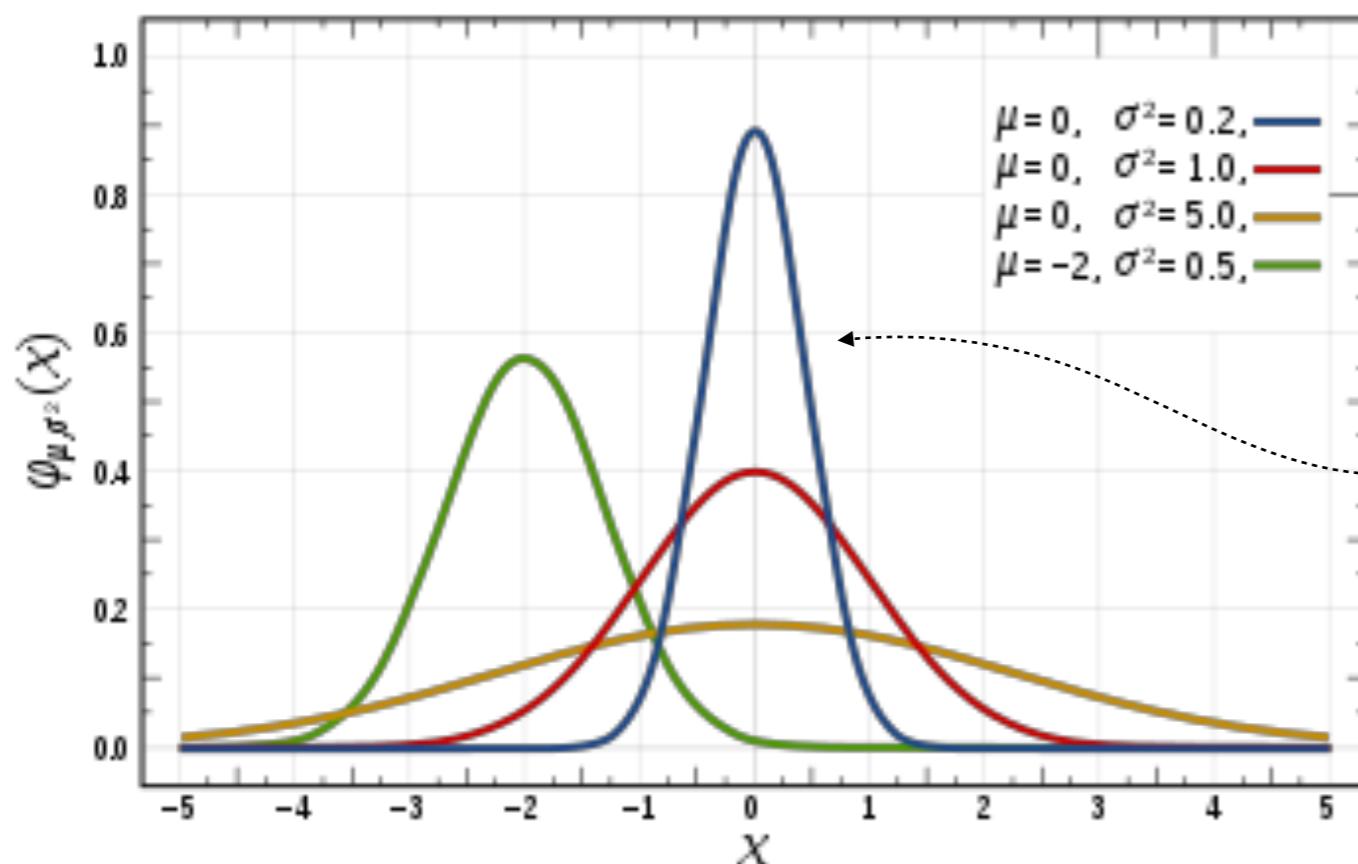
Let's revisit Gaussian Distribution

Gaussian (Normal) Distribution

Suppose $x \in \mathbb{R}$. If x is distributed Gaussian with mean μ and variance σ^2 .

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

‘distributed as’

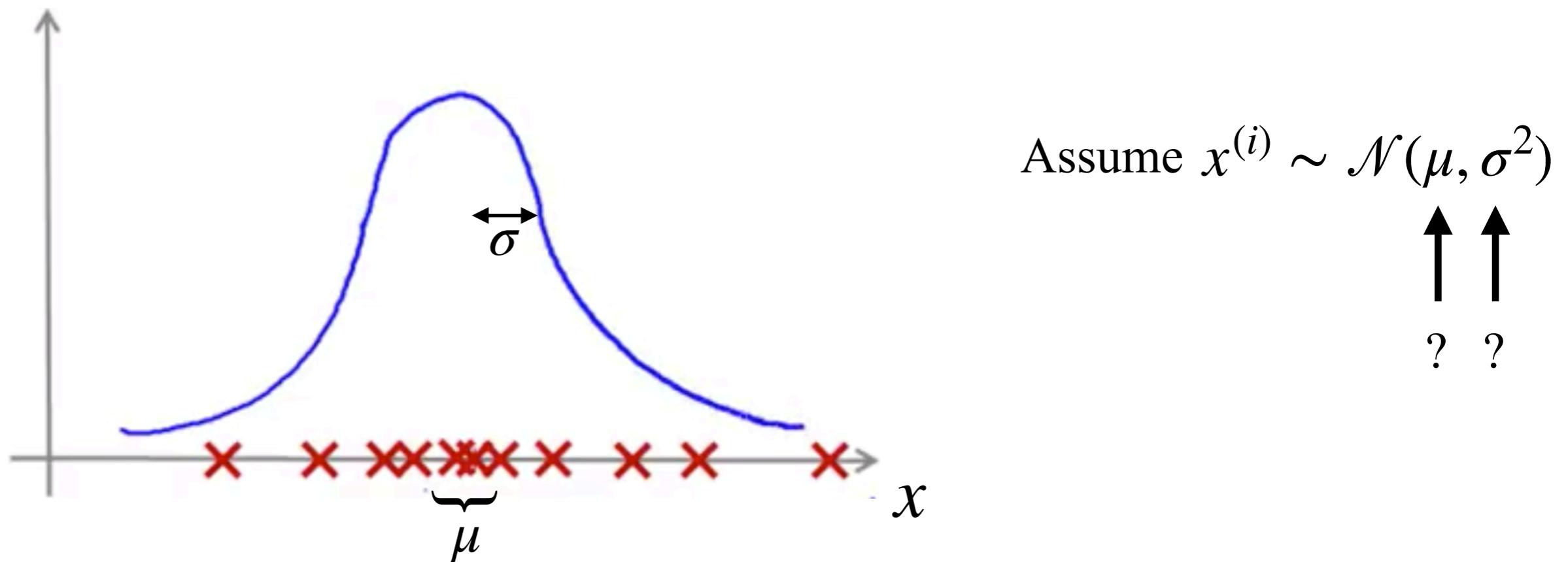


This specifies the probability of x taking on different values

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

Parameter Estimation Problem

Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ where $x^{(i)} \in \mathbb{R}$



$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$$

Question

The formula for the Gaussian density is:

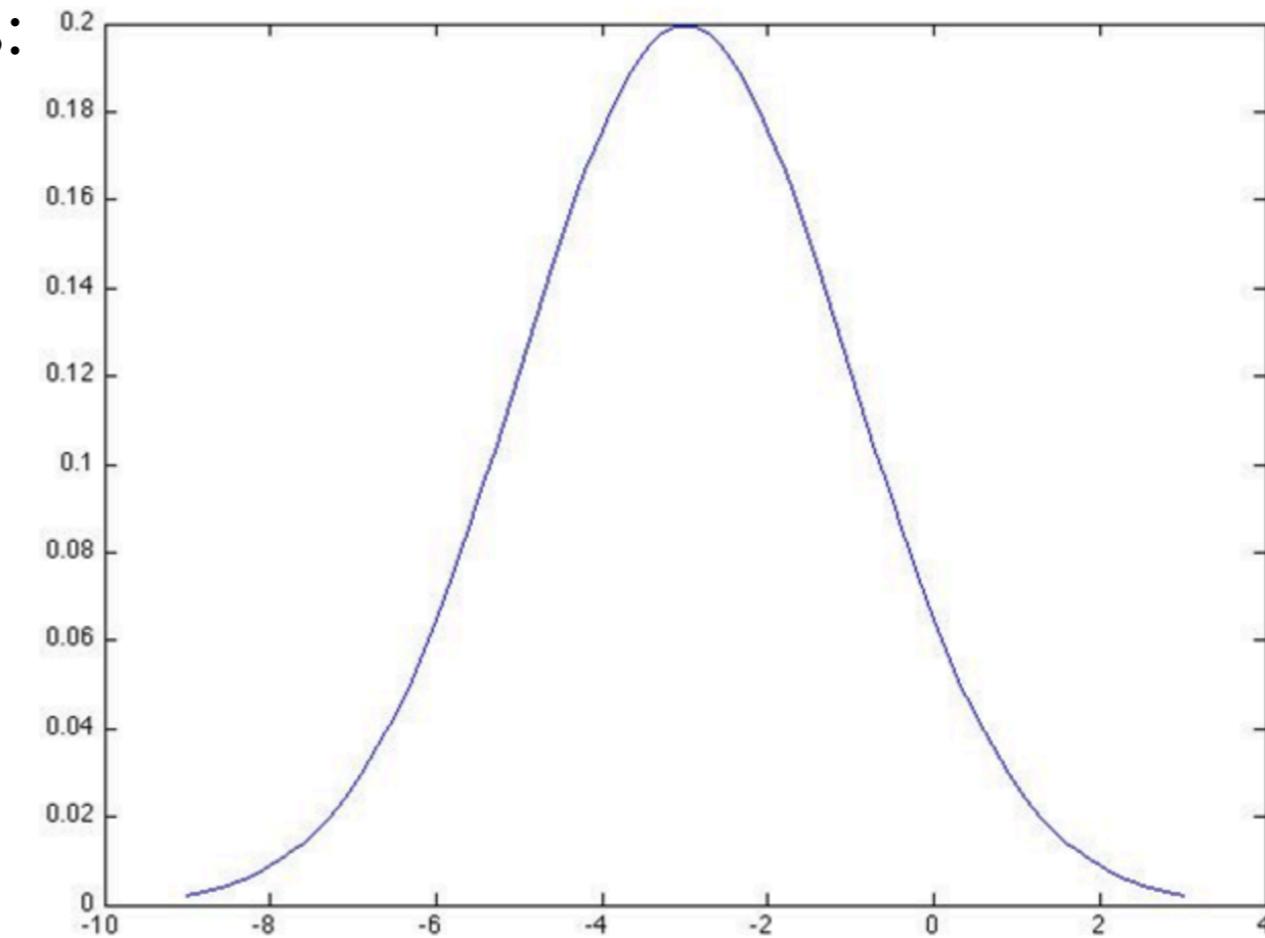
$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Which of the following is the formula for the density to the right?

(i) $p(x) = \frac{1}{\sqrt{2\pi} \times 2} \exp\left(-\frac{(x-3)^2}{2 \times 4}\right)$

(ii) $p(x) = \frac{1}{\sqrt{2\pi} \times 4} \exp\left(-\frac{(x-3)^2}{2 \times 2}\right)$

(iii) $p(x) = \frac{1}{\sqrt{2\pi} \times 2} \exp\left(-\frac{(x+3)^2}{2 \times 4}\right)$



(iv) $p(x) = \frac{1}{\sqrt{2\pi} \times 4} \exp\left(-\frac{(x+3)^2}{2 \times 2}\right)$

Multivariate Gaussian Distribution

Multivariate Gaussian Distribution

The n -dimensional multivariate Gaussian distribution has:

- a mean vector $\mu \in \mathbb{R}^n$
- a covariance matrix $\Sigma \in \mathbb{R}^{n \times n}$

We require that $\Sigma \geq 0$ is symmetric and positive semidefinite.

The distribution is written $\mathcal{N}(\mu, \Sigma)$ and the density is:

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

where $|\cdot|$ is the determinant.

Multivariate Gaussian Distribution

If $X \sim \mathcal{N}(\boldsymbol{\mu}; \Sigma)$, then we can write:

$$E[X] = \int_x x p(x; \boldsymbol{\mu}, \Sigma) dx = \boldsymbol{\mu}$$

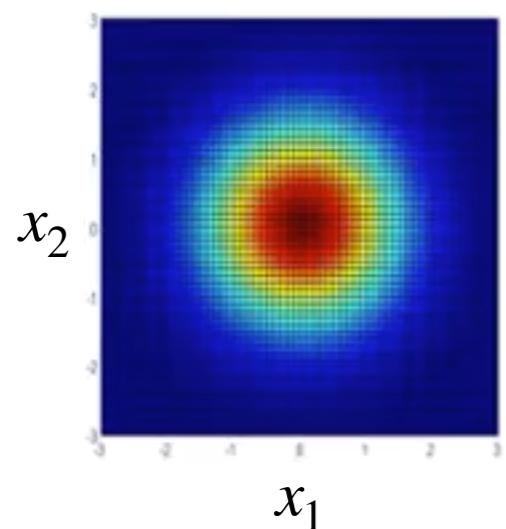
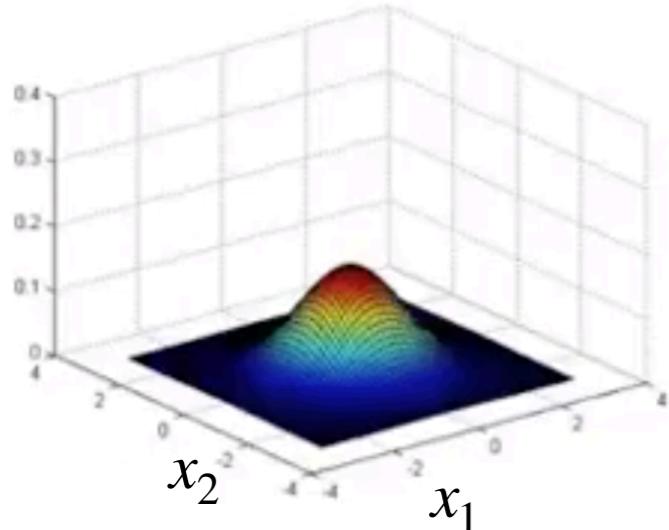
and

$$\mathbf{Cov}(X) = E[(X - E[X])(X - E[X])^T] = \Sigma$$

We will practice about the multivariate Gaussian in the lab to get an intuition about interpreting the covariance matrix.

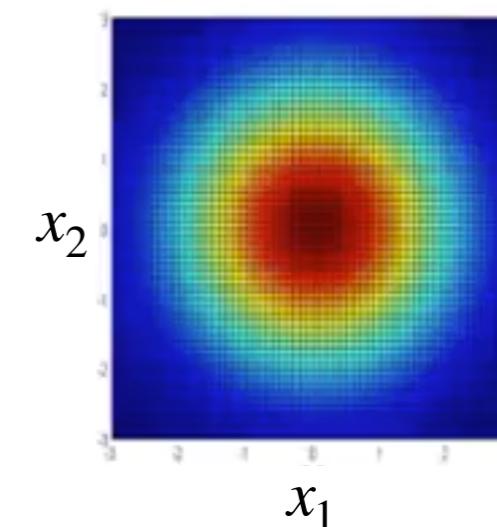
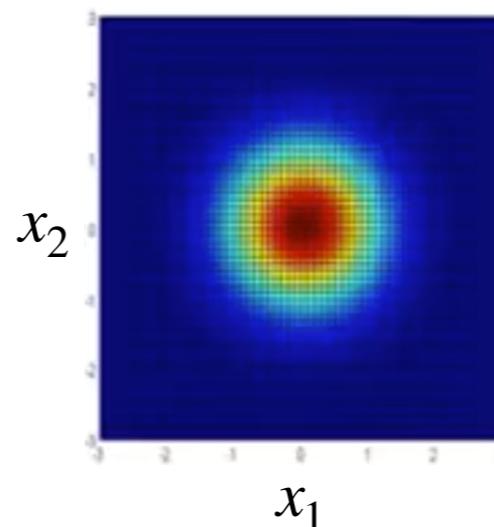
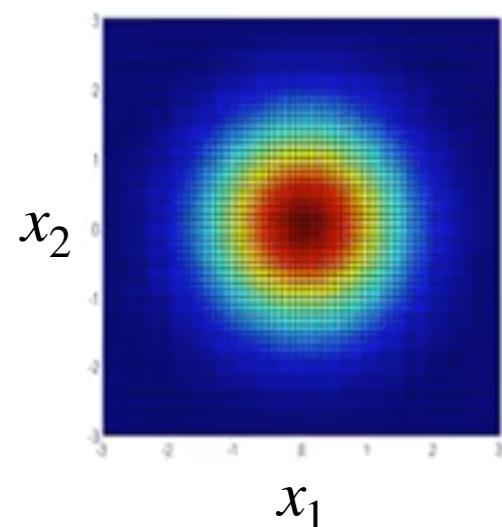
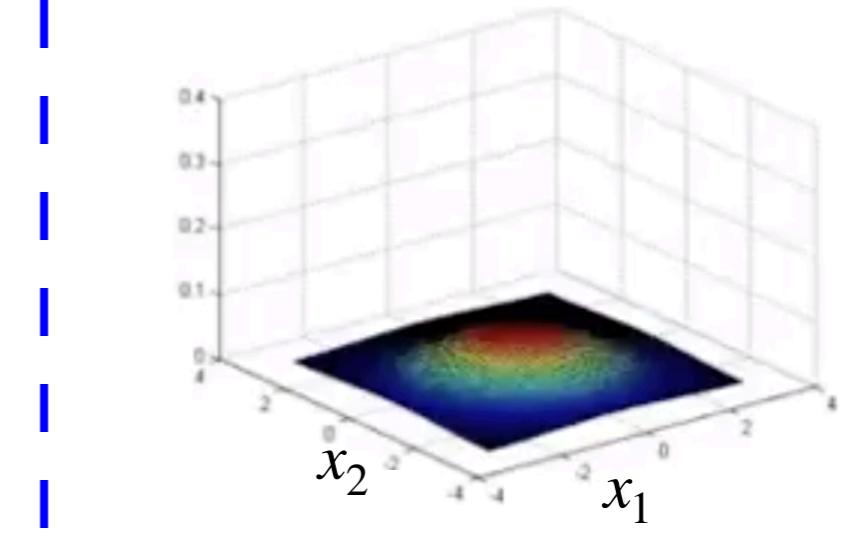
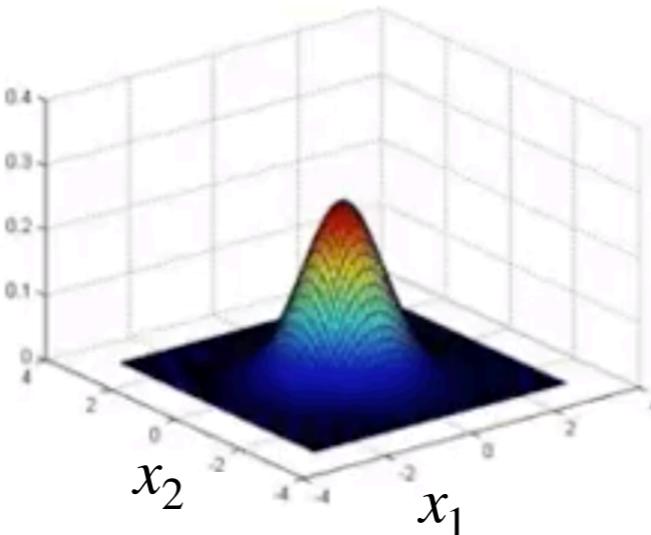
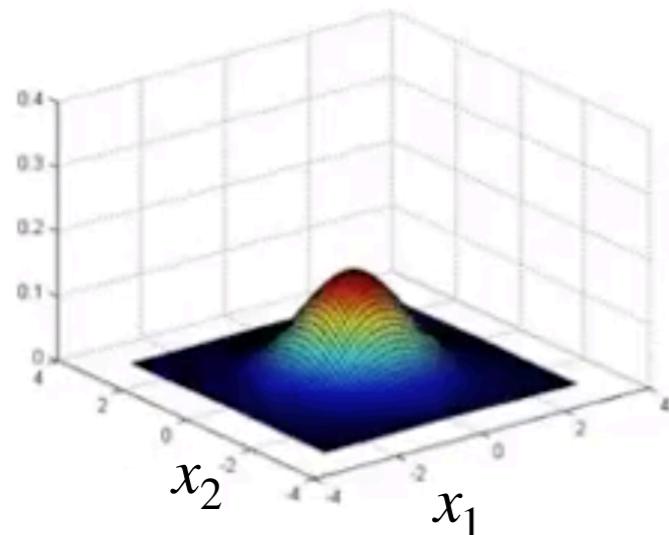
Multivariate Gaussian (Normal) Distribution Example

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (\text{Identity matrix})$$



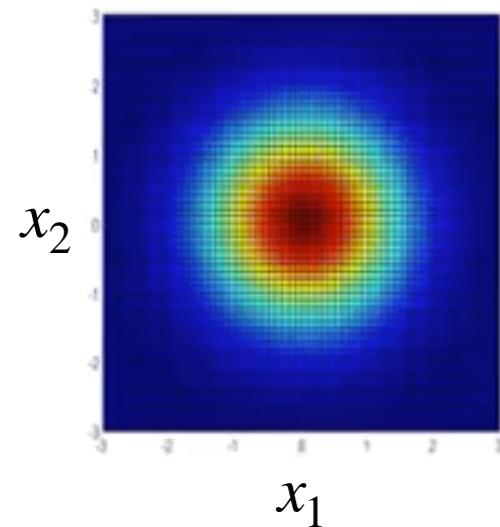
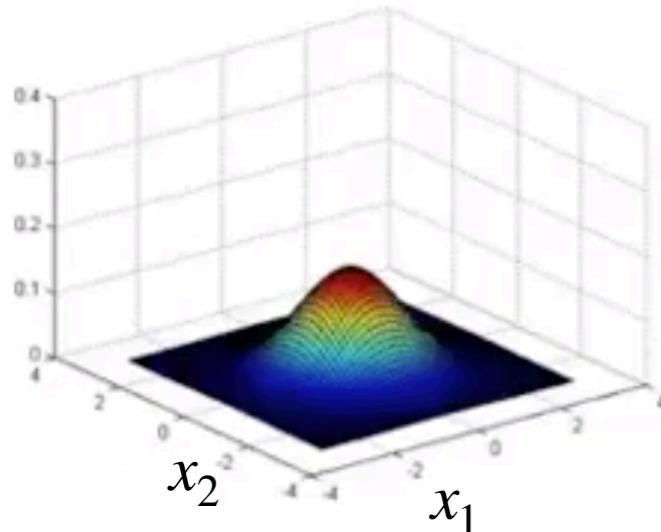
Multivariate Gaussian Distribution Example

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

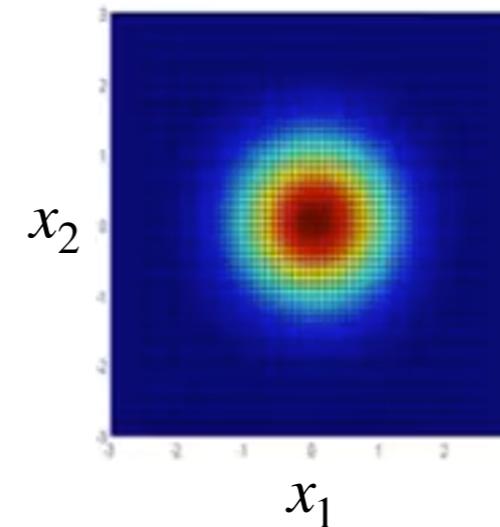
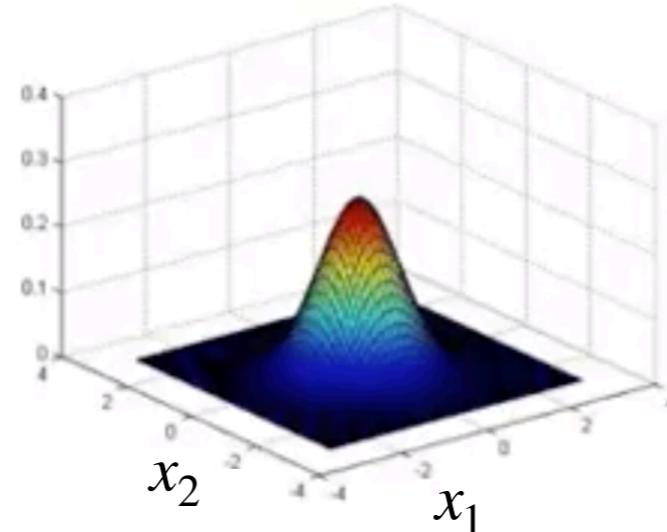


Multivariate Gaussian Distribution Example

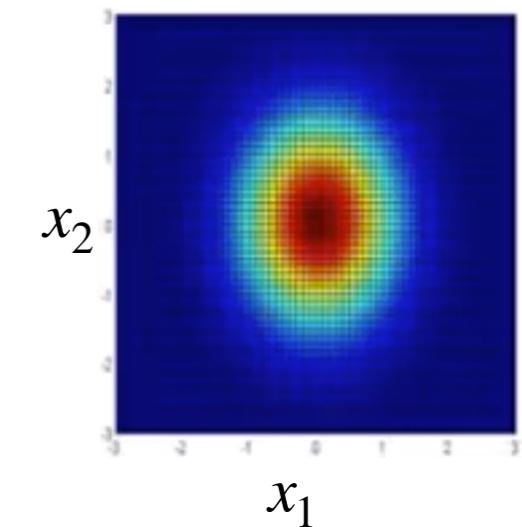
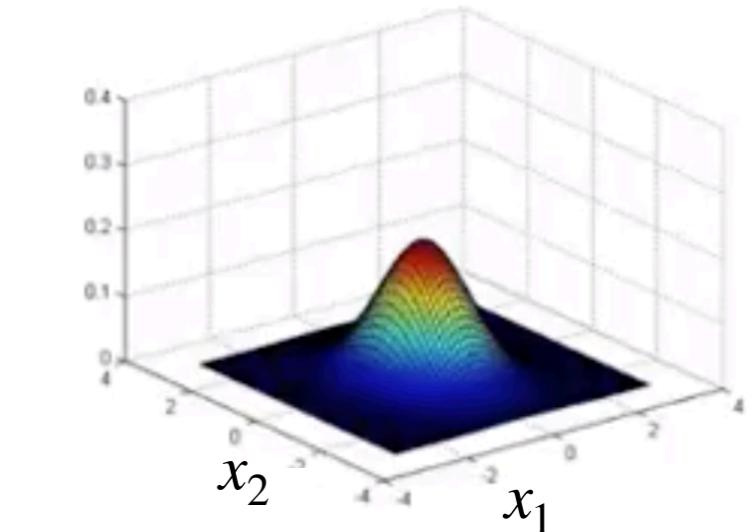
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



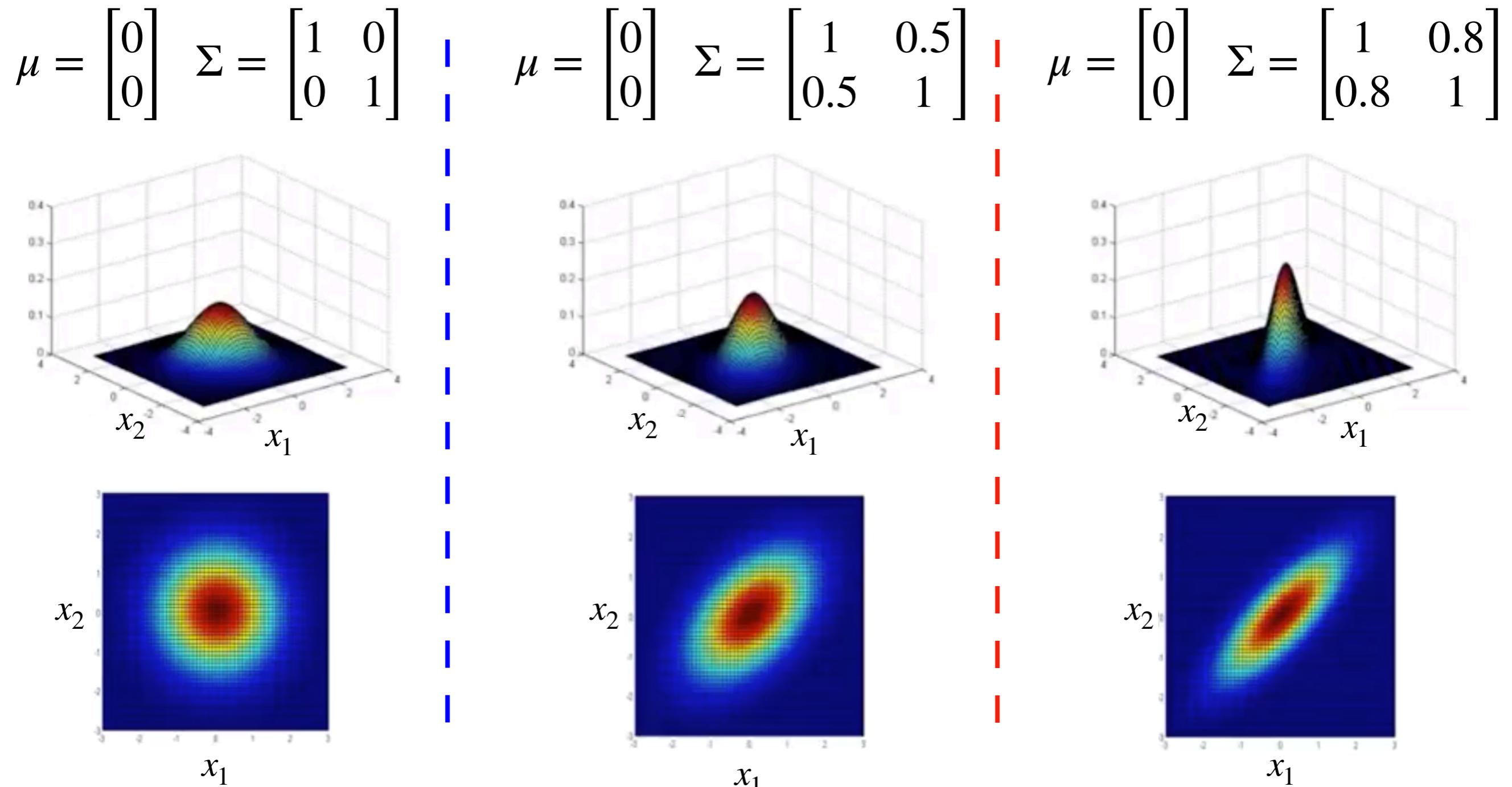
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 1 \end{bmatrix}$$

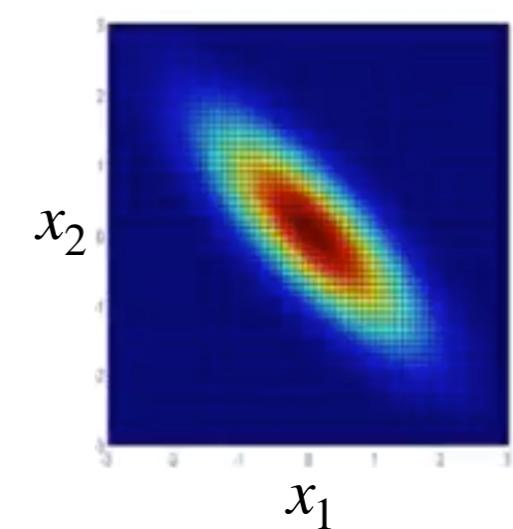
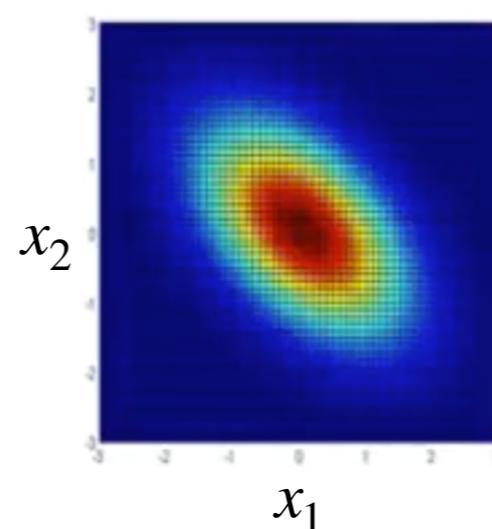
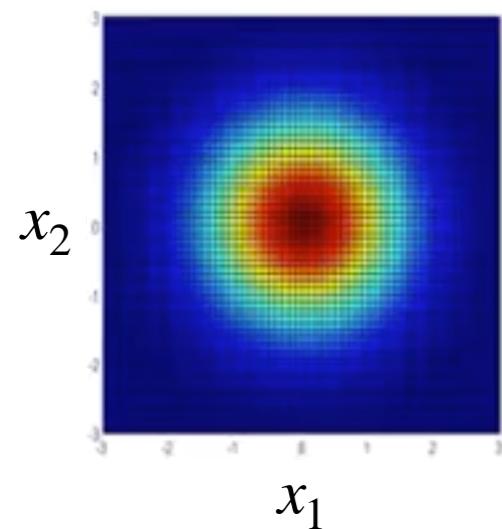
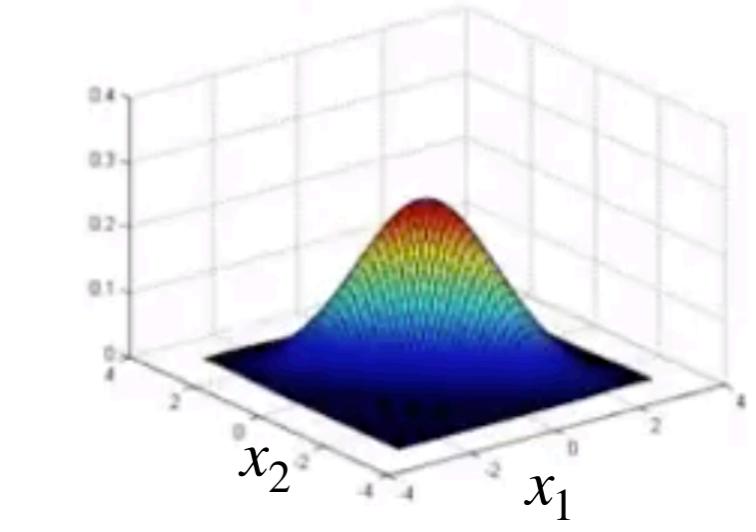
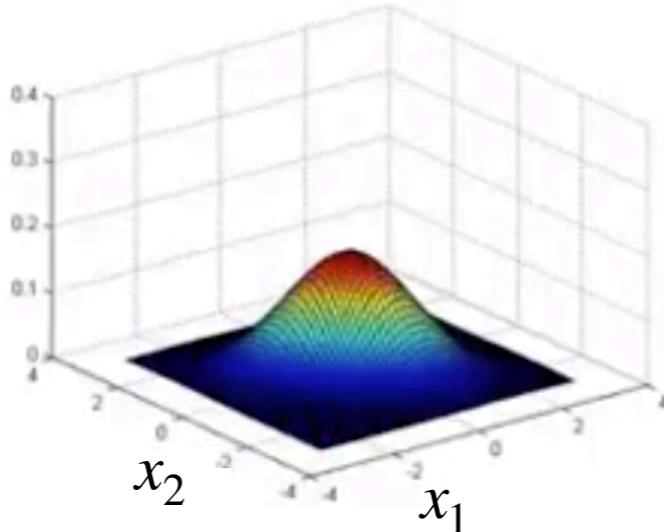
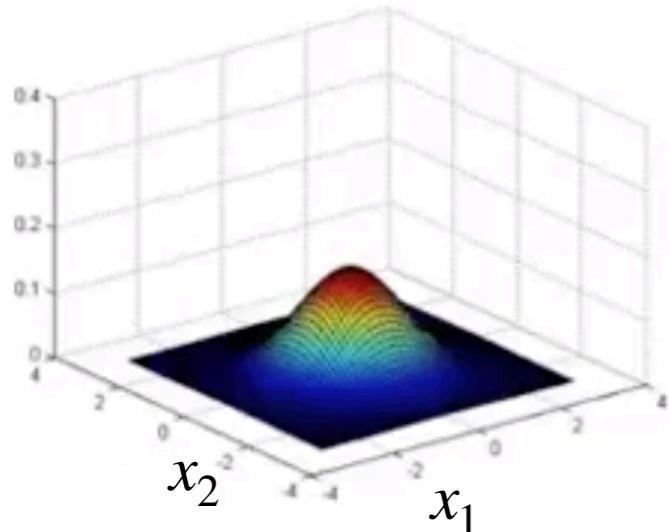


Multivariate Gaussian Distribution Example



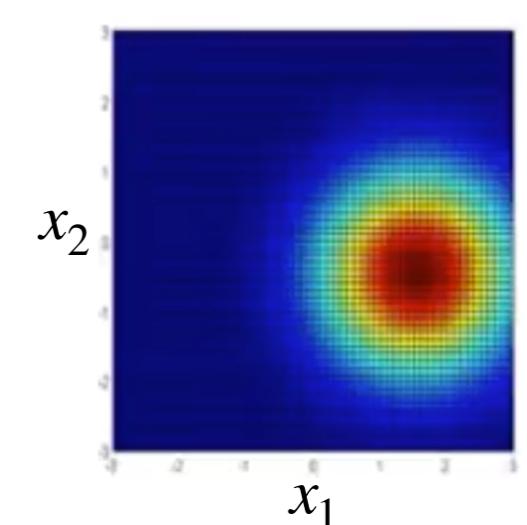
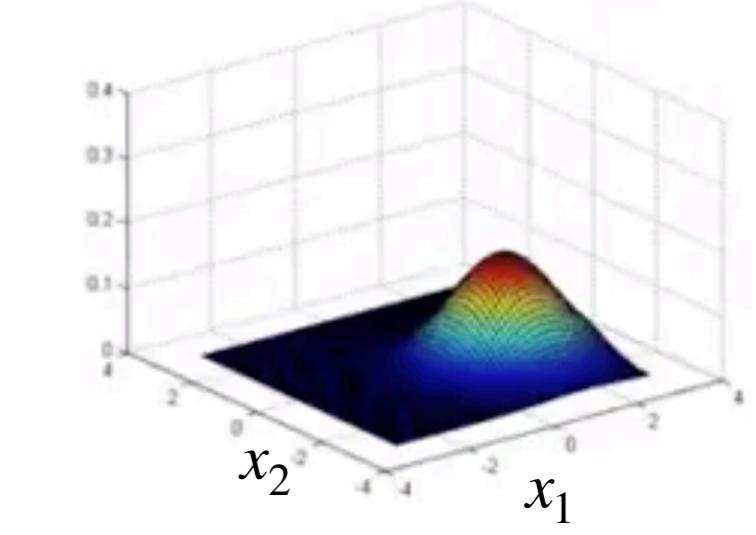
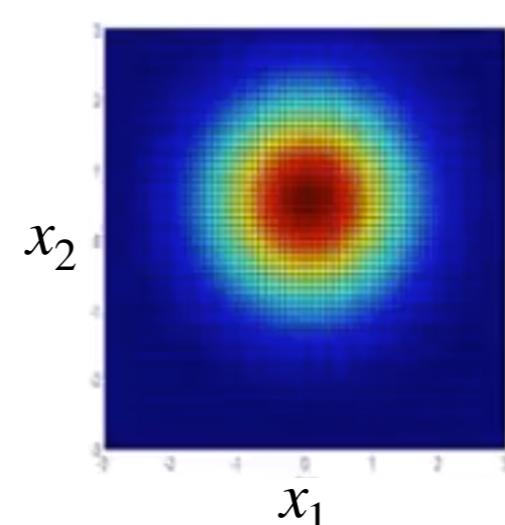
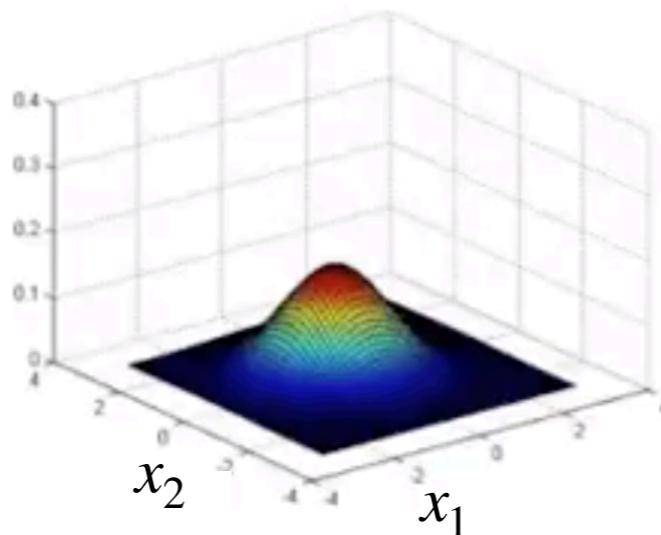
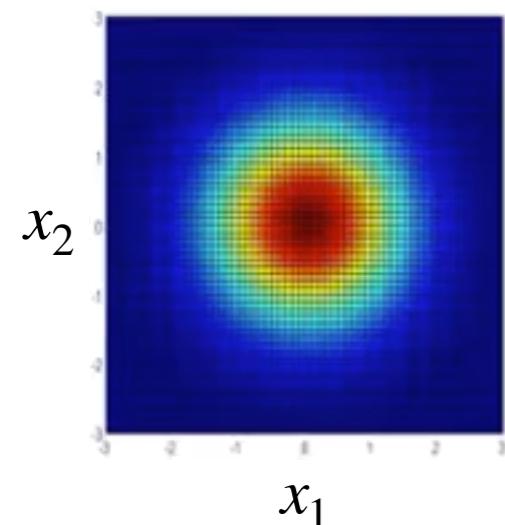
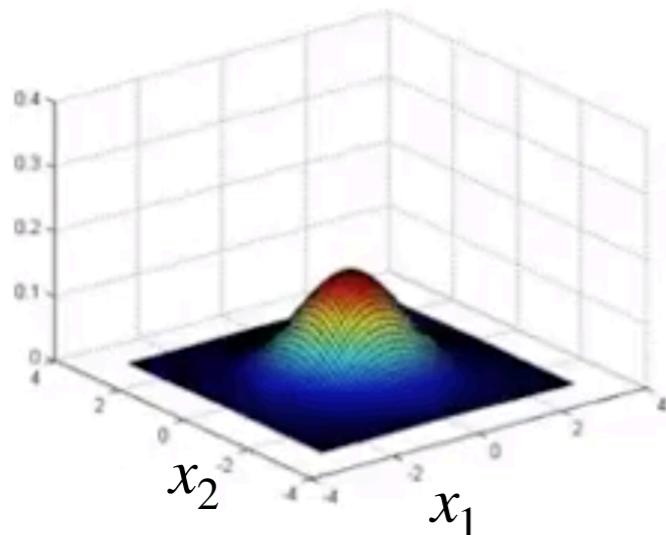
Multivariate Gaussian Distribution Example

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$



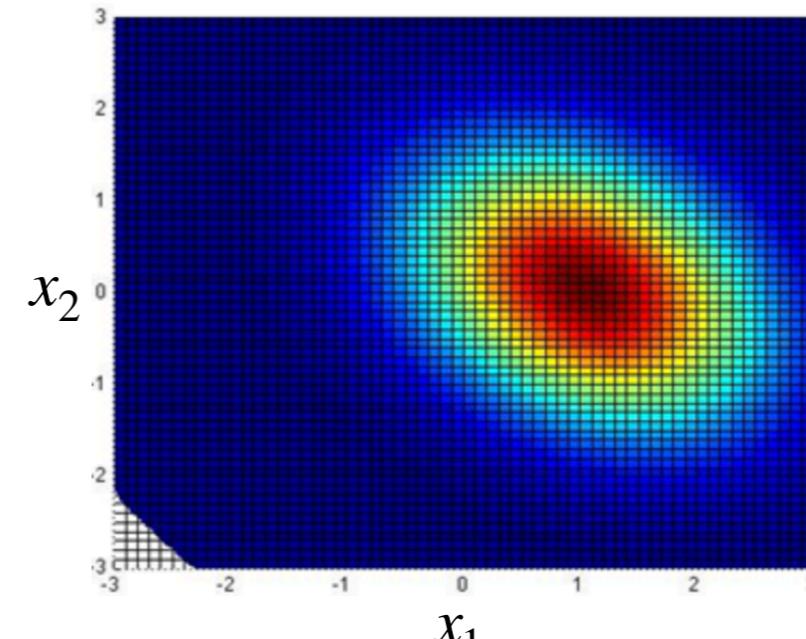
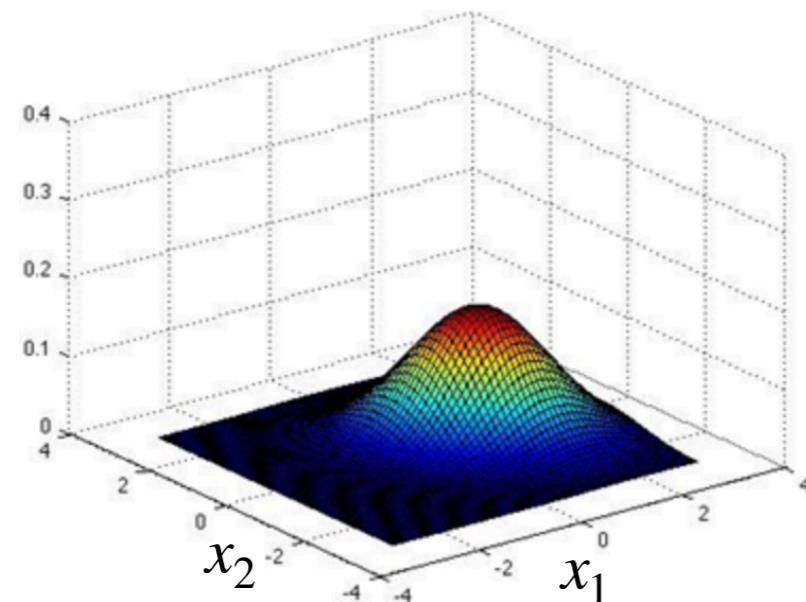
Multivariate Gaussian Distribution Example

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad | \quad \mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



Question

Consider the following multivariate Gaussian, which of the following are the values of μ and Σ for this distribution?



(i) $\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}$

(ii) $\mu = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & 0.3 \\ 0.3 & 1 \end{bmatrix}$

(iii) $\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$

(iv) $\mu = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \Sigma = \begin{bmatrix} 1 & -0.3 \\ -0.3 & 1 \end{bmatrix}$

Now, we're ready to talk
about GDA

Gaussian Discriminant Analysis

Now, the GDA model for a 2-class problem is:

$$y \sim \mathbf{Bernoulli}(\phi)$$

$$(x | y = 0) \sim \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$$

$$(x | y = 1) \sim \mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$$

We can then write the distributions as follows:

$$p(y) = \phi^y(1 - \phi)^{1-y}$$

$$p(x | y = 0) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x - \boldsymbol{\mu}_0)^T \Sigma_0^{-1} (x - \boldsymbol{\mu}_0)}$$

$$p(x | y = 1) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x - \boldsymbol{\mu}_1)^T \Sigma_1^{-1} (x - \boldsymbol{\mu}_1)}$$

Gaussian Discriminant Analysis

Now, it's time to optimize !

$$\begin{aligned} l(\phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma_0, \Sigma_1) &= \log \prod_{i=1}^m p(\mathbf{x}^{(i)}, y^{(i)}; \phi, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma_0, \Sigma_1) \\ &= \log \prod_{i=1}^m p(\mathbf{x}^{(i)} | y^{(i)}; \boldsymbol{\mu}_0, \boldsymbol{\mu}_1, \Sigma_0, \Sigma_1) p(y^{(i)}; \phi) \end{aligned}$$

Gaussian Discriminant Analysis

It turns out that if we assume $\Sigma_0 = \Sigma_1 = \Sigma$, we obtain the maximum likelihood estimates as follows:

$$\phi = \frac{1}{m} \sum_{i=1}^m 1\{y^{(i)} = 1\}$$

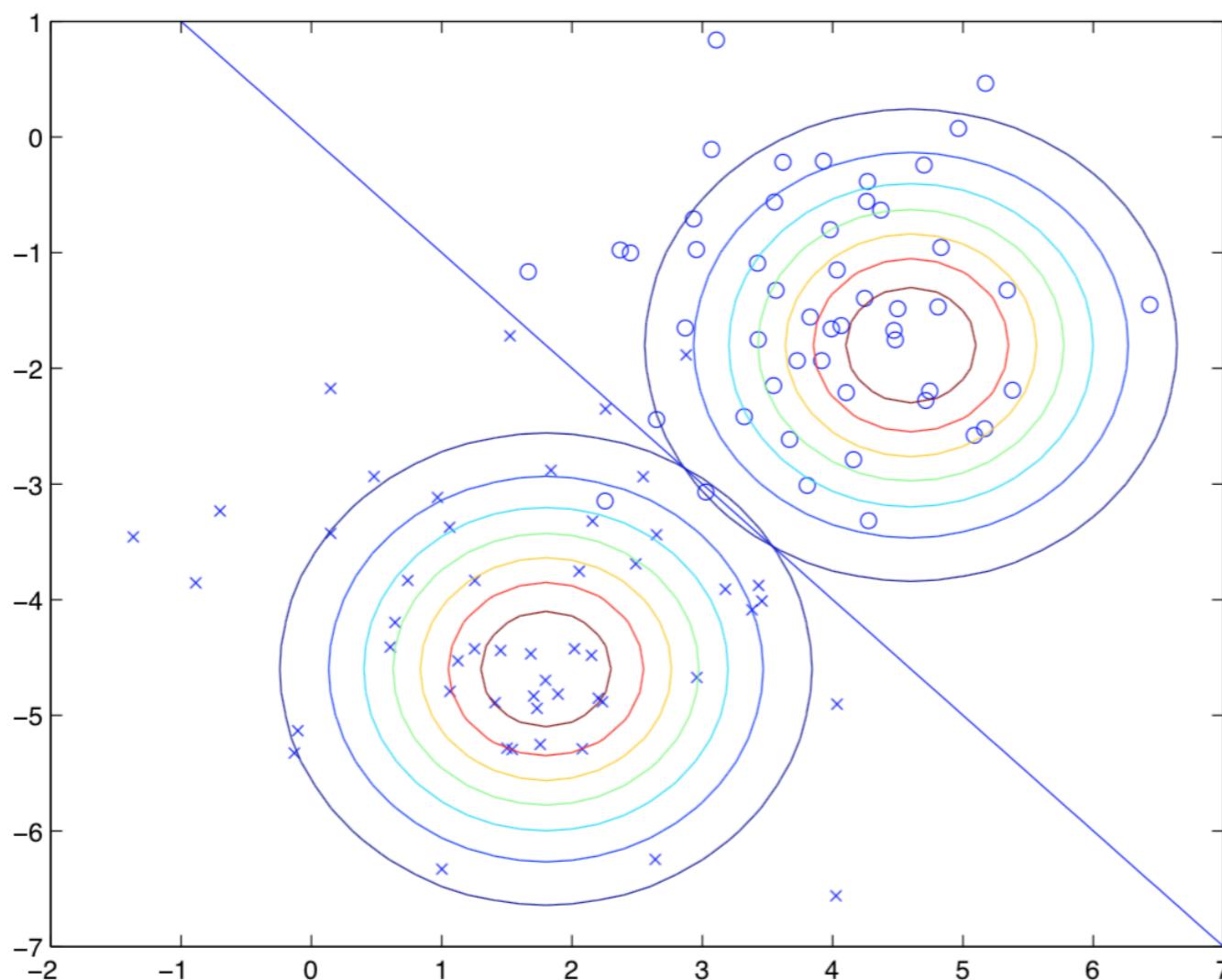
$$\boldsymbol{\mu}_0 = \frac{\sum_{i=1}^m 1\{y^{(i)} = 0\} \mathbf{x}^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}}$$

$$\boldsymbol{\mu}_1 = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\} \mathbf{x}^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})(\mathbf{x}^{(i)} - \boldsymbol{\mu}_{y^{(i)}})^T$$

Gaussian Discriminant Analysis

Pictorially, what the algorithm is doing can be seen in as follows:



- Two Gaussians have the same shape and orientation since they share Σ ; they have different μ .
- The straight line is the decision boundary at which $p(y = 1 | x) = 0.5$. On one side of the boundary, we'll predict $y = 1$; on the other side, we'll predict $y = 0$.

Gaussian Discriminant Analysis

The assumption $\Sigma_0 = \Sigma_1 = \Sigma$ turns out to be very useful
i.e. it is possible to show that:

$$p(y = 1 | \mathbf{x}; \phi, \Sigma, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

with $\boldsymbol{\theta}$ as a function of $\phi, \Sigma, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1$.

This shows that the GDA model is related to the logistic regression.

Gaussian Discriminant Analysis

The assumption $\Sigma_0 = \Sigma_1 = \Sigma$ turns out to be very useful
i.e. it is possible to show that:

$$p(y = 1 | \mathbf{x}; \phi, \Sigma, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

with $\boldsymbol{\theta}$ as a function of $\phi, \Sigma, \boldsymbol{\mu}_0, \boldsymbol{\mu}_1$.

This shows that the GDA model is related to the logistic regression.

GDA and logistic regression have the same form but give different decision boundaries. We'll later practice about it in the lab !

GDA vs. Logistic Regression

- ▶ GDA will be better (*i.e.* will require less training data to provide accurate predictions) if $p(\mathbf{x} | y)$ is in fact **multivariate Gaussian** or **almost multivariate Gaussian**.
- ▶ Logistic regression will probably be better if $p(\mathbf{x} | y)$ is **definitely non-Gaussian** or **unknown**.

GDA ใช้ข้อมูลน้อยกว่าการใช้ Logistic regression

Summary

Let's again modify our check-list reminder !

- If you have continuous \mathcal{X} and continuous \mathcal{Y} , your first go-to model should be **linear regression**. Also, consider non-linear transformation of the inputs.
- If you have continuous \mathcal{X} and discrete \mathcal{Y} but don't know much about $p(x|y)$, your first go-to model should be **logistic or softmax regression**, or may come up with a new **GLM** from scratch.
- If you have continuous \mathcal{X} and discrete \mathcal{Y} and know something about $p(x|y)$, you should model the distribution accurately, as a **Gaussian (GDA)** or build a new **generative** model from scratch.

Naive Bayes Classifier

Naive Bayes

GDA was an example of a generative classification method for problems in which $\mathcal{X} = \mathbb{R}^n$.

What if our features have discrete values? e.g.

- Car buying behavior:
 - $\mathcal{Y} := \{\text{buy}, \neg\text{buy}\}$ and $\mathcal{X} := [\text{Size}, \text{Color}, \text{Gender}, \text{Age}]$, where **Size** := {small, medium, large}, **Color** := {red, green, blue}, **Gender** := {male, female}, **Age** := {0 - 18, 19 - 39, 30 - 39, 40+}
- Email spam filter:
 - $\mathcal{Y} := \{\text{spam}, \neg\text{spam}\}$ and $\mathcal{X} := \{0,1\}^K$, where each variable x_1, \dots, x_k representing presence/absence of a particular word in a vocabulary.

Naive Bayes

Whatever \mathcal{Y} and \mathcal{X} , a generative model needs a form for $p(y | x)$.

The simplest approach to these problems is to **use the multinomial distribution over the set of possible outcomes** for x .

Exercise: How many outcomes for x are there for the car buying and spam filter examples?

The full multinomial model will thus have $2L - 1$ parameters for $p(x | y)$, where L is the number of outcomes for x , and 1 parameter for $p(y)$, assuming y binary.

Is this practical for the car buying example? For the spam example?

Naive Bayes

Naive Bayes attempts to reduce the number of parameters required using the (very strong but very useful) assumption that the features are conditionally independent given y :

$$\begin{aligned} p(x|y) &= p(x_1, x_2, \dots, x_n | y) \\ &= p(x_1 | y) \cdot p(x_2 | x_1, y) \cdot \dots \cdot p(x_n | x_1, x_2, \dots, x_{n-1}, y) \\ &\approx \prod_{i=1}^n p(x_i | y) \end{aligned}$$

This model requires a set of parameters $\phi_{ij|y=1} = p(x_i = j | y = 1)$ and a set of parameters $\phi_{ij|y=0} = p(x_i = j | y = 0)$.

Noted that if the variables x_i are binary, we only need two parameters $\phi_{i|y=1} = p(x_i = 1 | y = 1)$ and $\phi_{i|y=0} = p(x_i = 1 | y = 0)$.

Naive Bayes

For the case where x_i are binary, we get the joint likelihood:

$$\mathcal{L}(\phi_y, \phi_{j|y=0}, \phi_{j|y=1}) = \prod_{i=1}^m p(x^{(i)}, y^{(i)})$$

Maximizing this w.r.t. $\phi_y, \phi_{j|y=0}, \phi_{j|y=1}$ gives the following maximum likelihood estimates:

$$\phi_{j|y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}}$$

$$\phi_y = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}$$

Naive Bayes

If n is large, there may be some features that do not appear in all combinations with every class.

Example:

- The term ‘viagra’ might occur only in spam emails.
- In a given period of time, there might not be any customers 18 years or younger who bought a car.

Suppose we are then faced with an email x with the term ‘viagra’ or a customer x whose age is 0-18.

Do you predict $p(\text{spam} | x) = 1$ and $p(\text{spam} | x) = 0$?

Naive Bayes

Laplace smoothing avoids 0-probability issues by adding one pseudo-example to the dataset:

$$\phi_{j|y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 1\} + 2}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 0\} + 2}$$

This is for binary features (*i.e.* a multivariate Bernoulli event model).

Also, see Ng's lecture note (provided in the external materials) for more information about the multivariate event model !

Summary

Let's again modify our check-list reminder !

- If you have continuous \mathcal{X} and continuous \mathcal{Y} , your first go-to model should be **linear regression**. Also, consider non-linear transformation of the inputs.
- If you have continuous \mathcal{X} and discrete \mathcal{Y} but don't know much about $p(x|y)$, your first go-to model should be **logistic or softmax regression**, or may come up with a new **GLM** from scratch.
- If you have continuous \mathcal{X} and discrete \mathcal{Y} and know something about $p(x|y)$, you should model the distribution accurately, as a **Gaussian (GDA)** or build a new **generative** model from scratch.
- If you have discrete \mathcal{X} and \mathcal{Y} , you should probably start with **naive Bayes** and build up from there.

That's it for Generative Learning !
We'll come back to Discriminative
Learning again when we discuss
about kernel methods and SVMs.