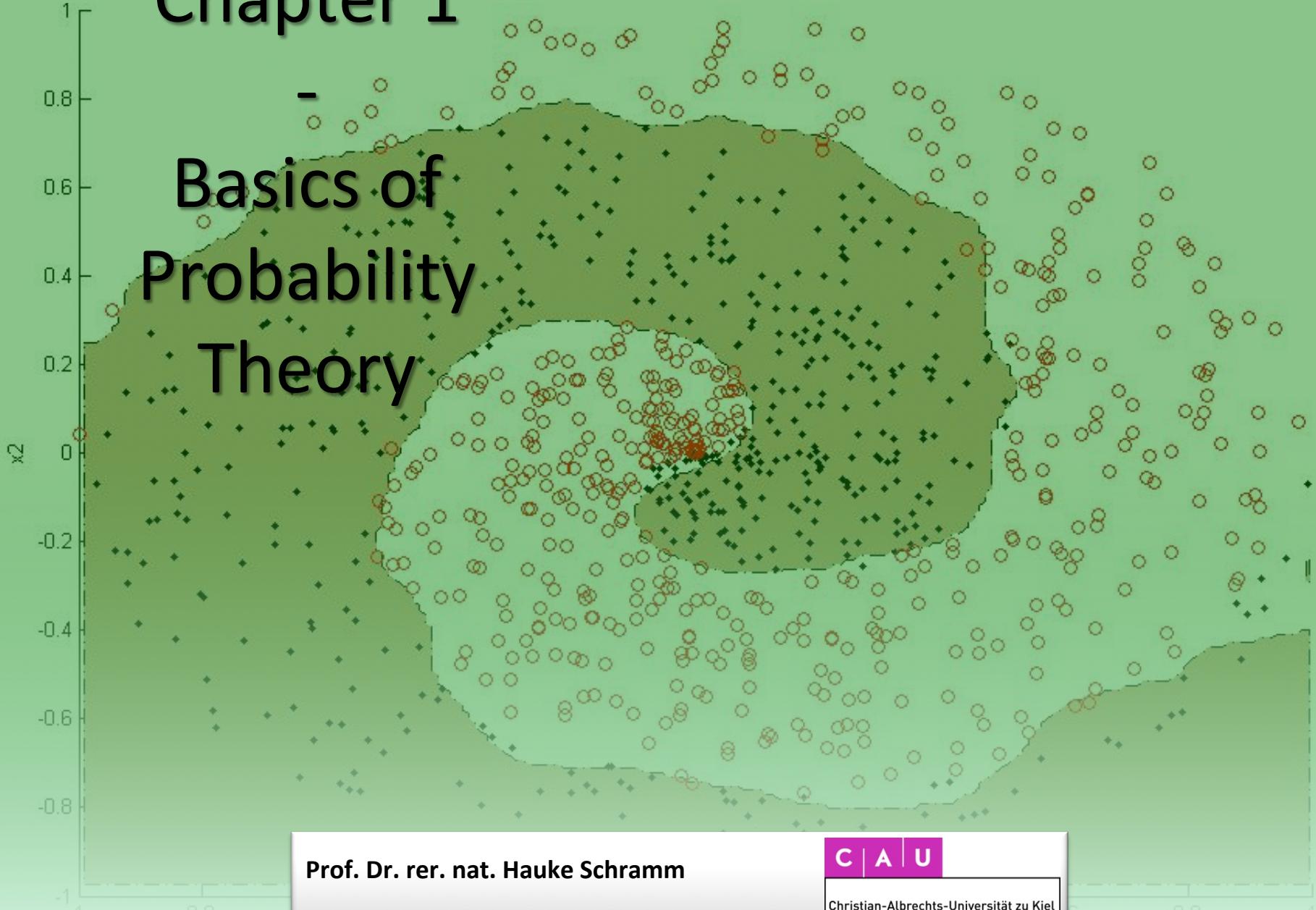


Chapter 1

Basics of Probability Theory



Prof. Dr. rer. nat. Hauke Schramm

C | A | U

Christian-Albrechts-Universität zu Kiel

Institut für Informatik

Topics of the lecture

Introduction

1. Fundamentals of probability theory

Probability distribution, expectation, covariance matrix, marginal, multivariate Gaussian, ...

2. Bayesian decision theory

Bayes theorem, prior-/a-priori-/a-posteriori-probabilities, Bayes risk, decision boundaries, ...

3. Parameter estimation

Maximum-likelihood, Bayes learning, Distances in \mathbb{R}^D

4. Non-parametric techniques

Density estimation, Parzen windows, nearest neighbor classification

Outline of Chapter 1

1. Basics of probability theory
 1. Discrete random variables
 2. Expected values
 3. Pairs of discrete random variables
 4. Marginal distributions
 5. Statistical independence
 6. Expected values of functions of two variables
 7. Covariance
 8. Conditional probability
 9. Vector random variables
 10. Continuous random variables

Outline of Chapter 1

1. Basics of probability theory (continuation)
12. Normal distributions
13. Mahalanobis distance
14. Multivariate normal densities
15. Mixture densities
16. Bayes rule

1. Basics of probability theory

1.1 Discrete random variables

The discrete random variable x can assume a finite number m of

different values in the set:

$$X = \{v_1, v_2, \dots, v_m\}$$


also called sample space

Probability that x takes the value v_i is denoted by:

$$p_i = Pr(x = v_i), \quad i = 1, \dots, m$$

Probabilities p_i must satisfy the following conditions:

$$p_i \geq 0$$

$$\sum_{i=1}^m p_i = 1$$

1.1 Discrete random variables

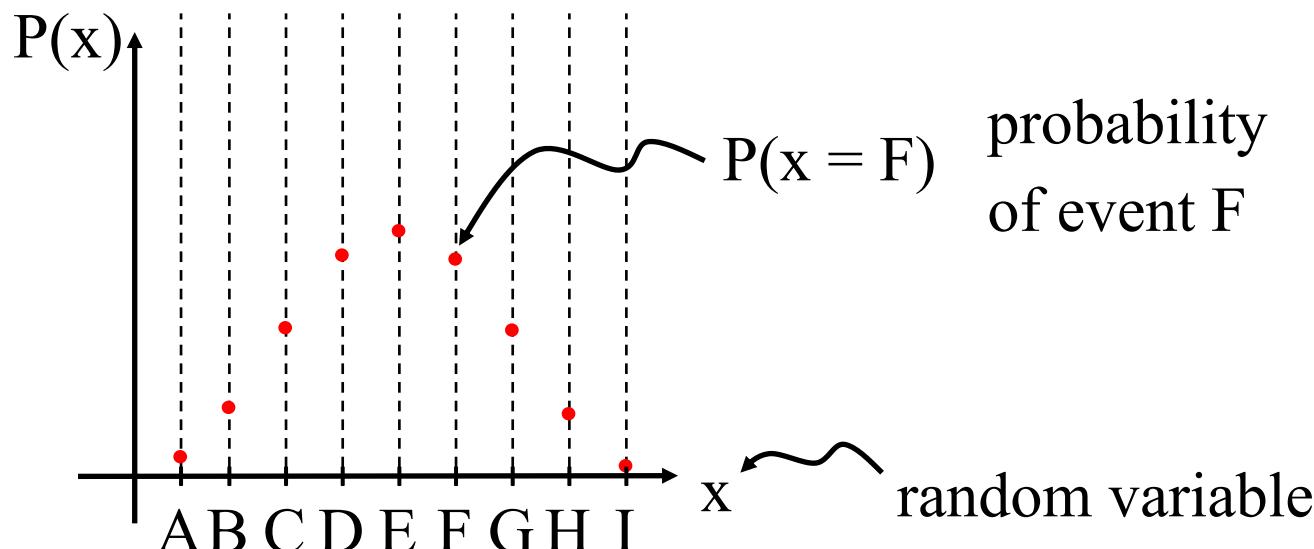
For convenience set of probabilities $\{p_1, p_2, \dots, p_m\}$ is expressed as

probability mass function $P(x)$ which must satisfy the conditions:

$$P(x) \geq 0$$

$$\sum_{x \in X} P(x) = 1$$

$$\sum_{x \notin X} P(x) = 0$$



1.1 Discrete random variables

The set of all realizations that have a strictly positive probability of being observed is called the **support** S_x .

Example:

If the discrete random variable x has the probability mass function

$$P(x) = \begin{cases} 0.5 & \text{if } x = 0 \\ 0.5 & \text{if } x = 1 \\ 0 & \text{otherwise} \end{cases}$$

the support is: $S_x = \{ 0, 1 \}$

1.2 Expected values

The expected value (mean) of a random variable x is defined as:

$$E[x] = \mu = \sum_{x \in X} x P(x) = \sum_{i=1}^m v_i p_i$$

Interpretation: μ is **arithmetic average** of values in a large random sample

1.2 Expected values

The expected value (mean) of a random variable x is defined as:

$$E[x] = \mu = \sum_{x \in X} x P(x) = \sum_{i=1}^m v_i p_i$$

Example:

Exam results:

1	2	3	4	5
0	1	3	0	0

Usual computation of the mean result of an exam?

Correspondence to the equation above?

1.2 Expected values

The expected value (mean) of a random variable x is defined as:

$$E[x] = \mu = \sum_{x \in X} x P(x) = \sum_{i=1}^m v_i p_i$$

Example:

Exam results:

1	2	3	4	5
0	1	3	0	0

$$\mu = (2 + 3 + 3 + 3) / 4 = 1 \cdot 0 + 2 \cdot \frac{1}{4} + 3 \cdot \frac{3}{4} + 4 \cdot 0 + 5 \cdot 0$$

$$\mu = 1 \cdot p(1) + 2 \cdot p(2) + 3 \cdot p(3) + 4 \cdot p(4) + 5 \cdot p(5)$$

1.2 Expected values

The expected value (mean) of a random variable x is defined as:

$$E[x] = \mu = \sum_{x \in X} x P(x) = \sum_{i=1}^m v_i p_i$$

Generalization to any **function** $f(x)$:

$$E[f(x)] = \sum_{x \in X} f(x) P(x) \quad \text{expected value of function } f$$

1.2 Expected values

Variance and standard deviation:

$$Var[x] = \sigma^2 = E[(x - \mu)^2] = \sum_{x \in X} (x - \mu)^2 P(x)$$

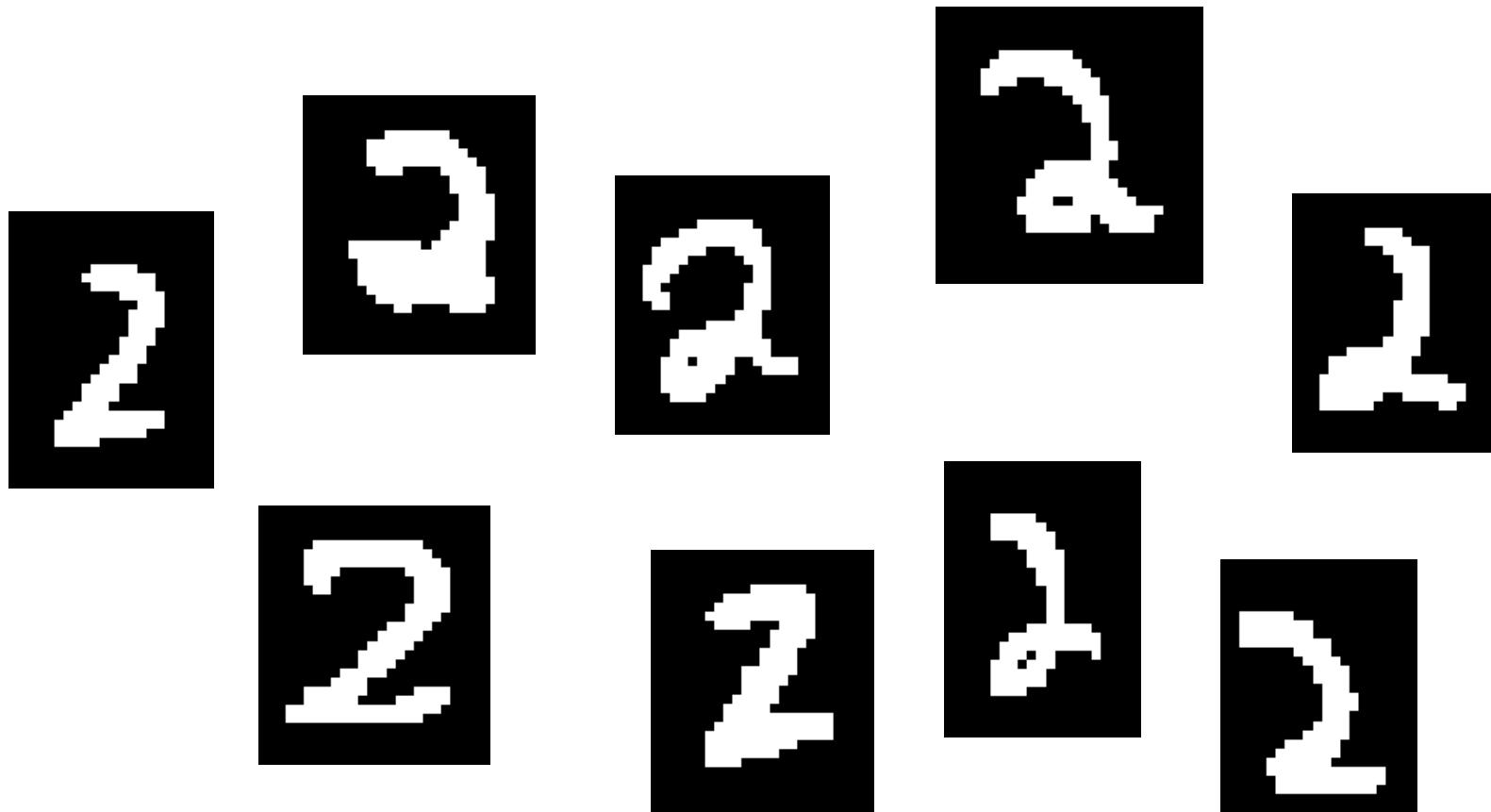
↑
Variance ↑
Standard
deviation σ

Simple but effective measure of how far values of x are likely to depart from the mean value.

1.2 Expected values

Example for practical computation of mean and variance:

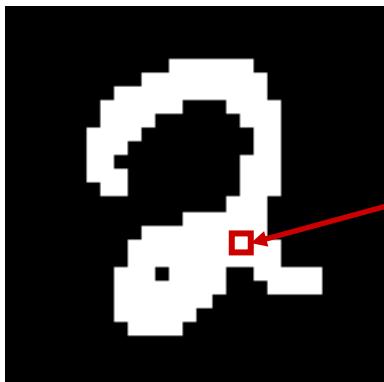
Consider a data set of 1000 handwritten digits of '2'



1.2 Expected values

Example for practical computation of mean and variance:

- Each digit image has $28 \times 28 = 784$ pixels
- Each pixel can be considered as a single random variable



Each pixel has gray values
between 0 and 255

1.2 Expected values

Goal:

- Compute the mean for each pixel and show result as new image
- Do the same for variance

1.2 Expected values

Goal:

- Compute the mean for each pixel and show result as new image
- Do the same for variance

Example for pixel 140:

$$E[x_{140}] = \sum_{x=0}^{255} x \cdot P(x) = 0 \cdot P(0) + 1 \cdot P(1) + 2 \cdot P(2) + \dots$$

Gray value 1
Probability of gray value 1 in pixel 140

1.2 Expected values

Goal:

- Compute the mean for each pixel and show result as new image
- Do the same for variance

Example for pixel 140:

$$E[x_{140}] = \sum_{x=0}^{255} x \cdot P(x) = 0 \cdot P(0) + 1 \cdot P(1) + 2 \cdot P(2) + \dots$$

Gray value 1

Probability of gray value 1 in pixel 140

Practically:

For each pixel p:

Step 1: Sum gray values of pixel p in all 1000 images

Step 2: Divide by 1000

1.2 Expected values

Use the resulting mean values as new gray values:

mean matrix

1.2 Expected values

Use the resulting mean values as new gray values:



illustration as
new image

1.2 Expected values

Variance image:



1.3 Pairs of discrete random variables

Let x and y be random variables which can take on values in:

$$X = \{v_1, v_2, \dots, v_m\} \quad Y = \{w_1, w_2, \dots, w_n\}$$

Each pair of events (v_i, w_j) has a *joint probability*

$$p_{ij} = \Pr(x = v_i, y = w_j)$$

with

$$p_{ij} \geq 0$$

$$\sum_{i=1}^m \sum_{j=1}^n p_{ij} = 1$$

1.3 Pairs of discrete random variables

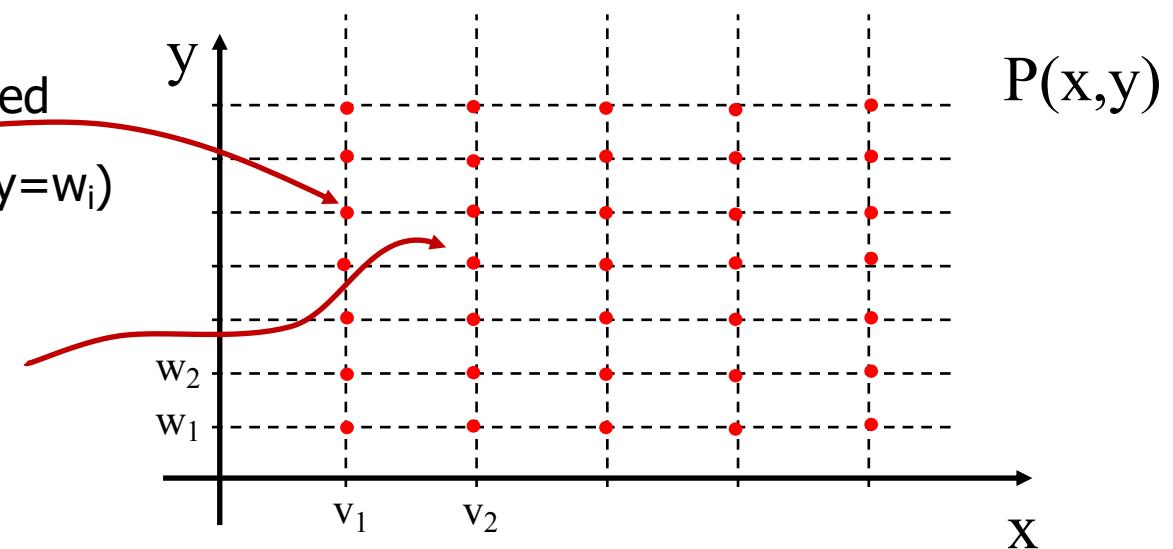
We can define a **joint probability mass function** $P(x, y)$ for which

$$P(x, y) \geq 0$$

$$\sum_{x \in X} \sum_{y \in Y} P(x, y) = 1$$

Each point (v_i, w_i) is assigned
a probability value $P(x=v_i, y=w_i)$

Note the discrete nature of
the random variables.



1.3 Pairs of discrete random variables

Example:

Let x and y be random variables with two possible values (events):

$$X = \{v_1, v_2\} \quad \text{and} \quad Y = \{w_1, w_2\}$$

v_1 : rain v_2 : no rain

w_1 : sunshine w_2 : no sunshine

Let also be $P(x, y)$ the known joint probability mass function.

1.3 Pairs of discrete random variables

v_1 : rain

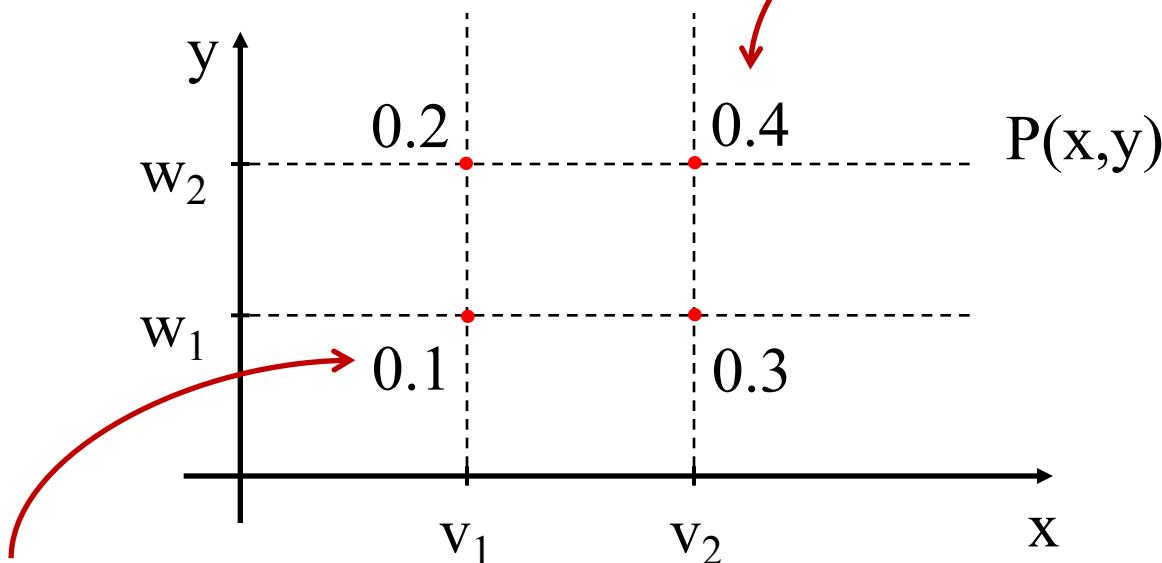
v_2 : no rain

w_1 : sunshine

w_2 : no sunshine

Illustration of $P(x, y)$

Combination of no rain
and no sunshine is likely.



Combination of rain and
sunshine is unlikely.

$$P(x, y) \geq 0$$

$$\sum_{x \in X} \sum_{y \in Y} P(x, y) = 1$$

1.4 Marginal distributions

Let $P(x, y)$ be a known joint probability mass function

The **marginal distribution** for x is given by summing over all y values:

$$P_x(x) = \sum_{y \in Y} P(x, y)$$

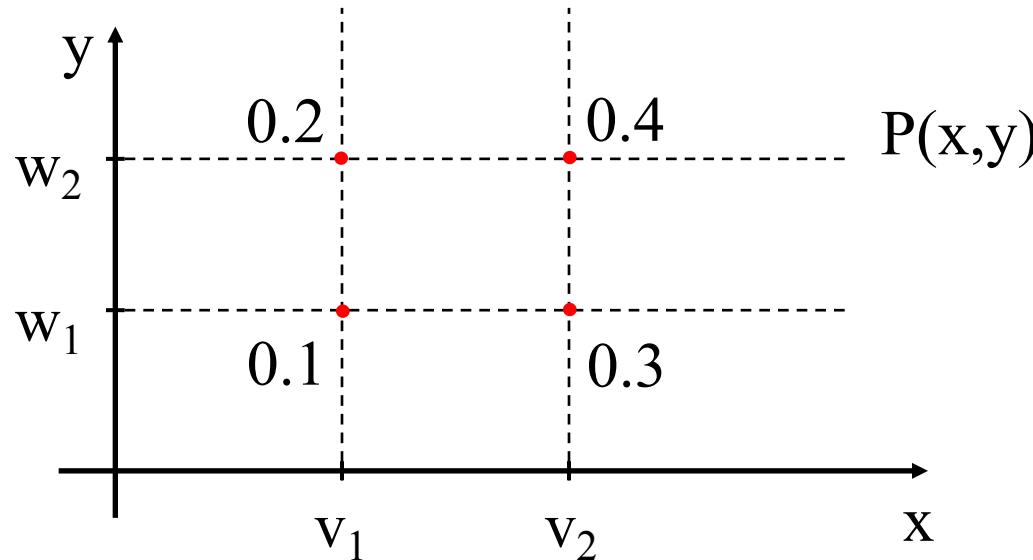
The marginal distribution for y is given by summing over all x values:

Remark: Subscript is used to emphasize that $P_x()$ has a different functional form than $P_y()$.

$$P_y(y) = \sum_{x \in X} P(x, y)$$

1.4 Marginal distributions

Example: Marginal distribution of x

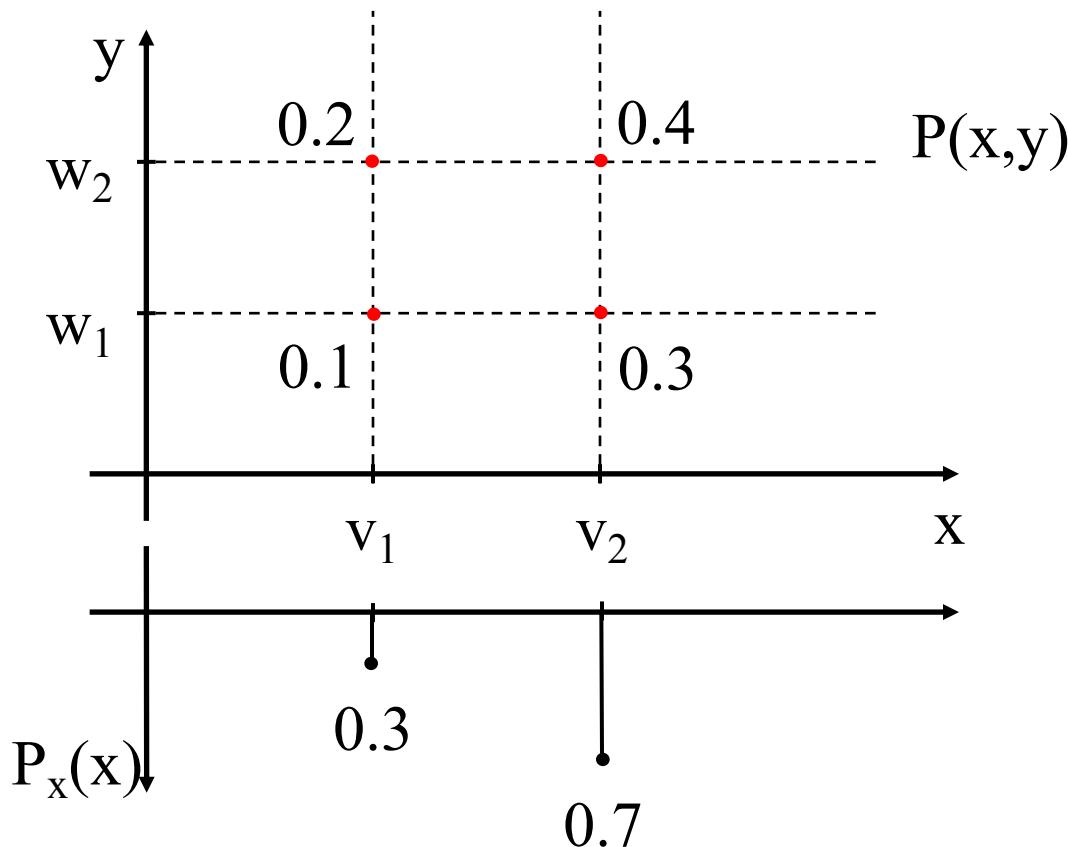


$$P_x(v_1) = \sum_{\text{all events in } Y} P(v_1, y) = P(v_1, w_1) + P(v_1, w_2) = 0.1 + 0.2$$

$$P_x(v_2) = \sum_{\text{all events in } Y} P(v_2, y) = P(v_2, w_1) + P(v_2, w_2) = 0.3 + 0.4$$

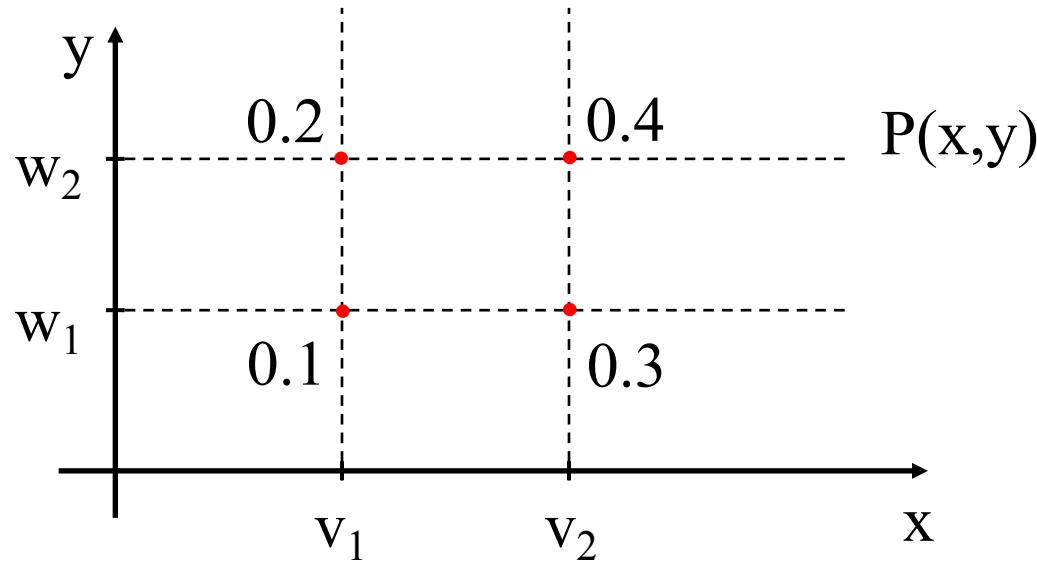
1.4 Marginal distributions

Example: Marginal distribution of x



1.4 Marginal distributions

Example: Marginal distribution of y

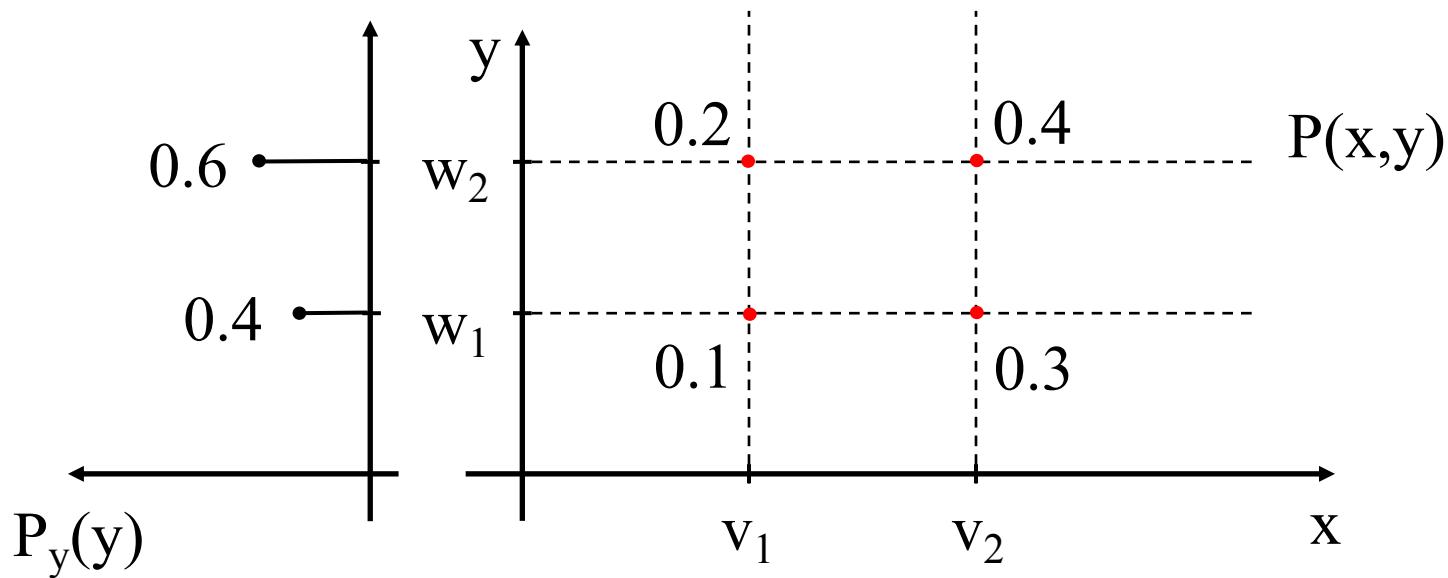


$$P_y(w_1) = \sum_{\text{all events in } X} P(x, w_1) = P(v_1, w_1) + P(v_2, w_1) = 0.1 + 0.3$$

$$P_y(w_2) = \sum_{\text{all events in } X} P(x, w_2) = P(v_1, w_2) + P(v_2, w_2) = 0.2 + 0.4$$

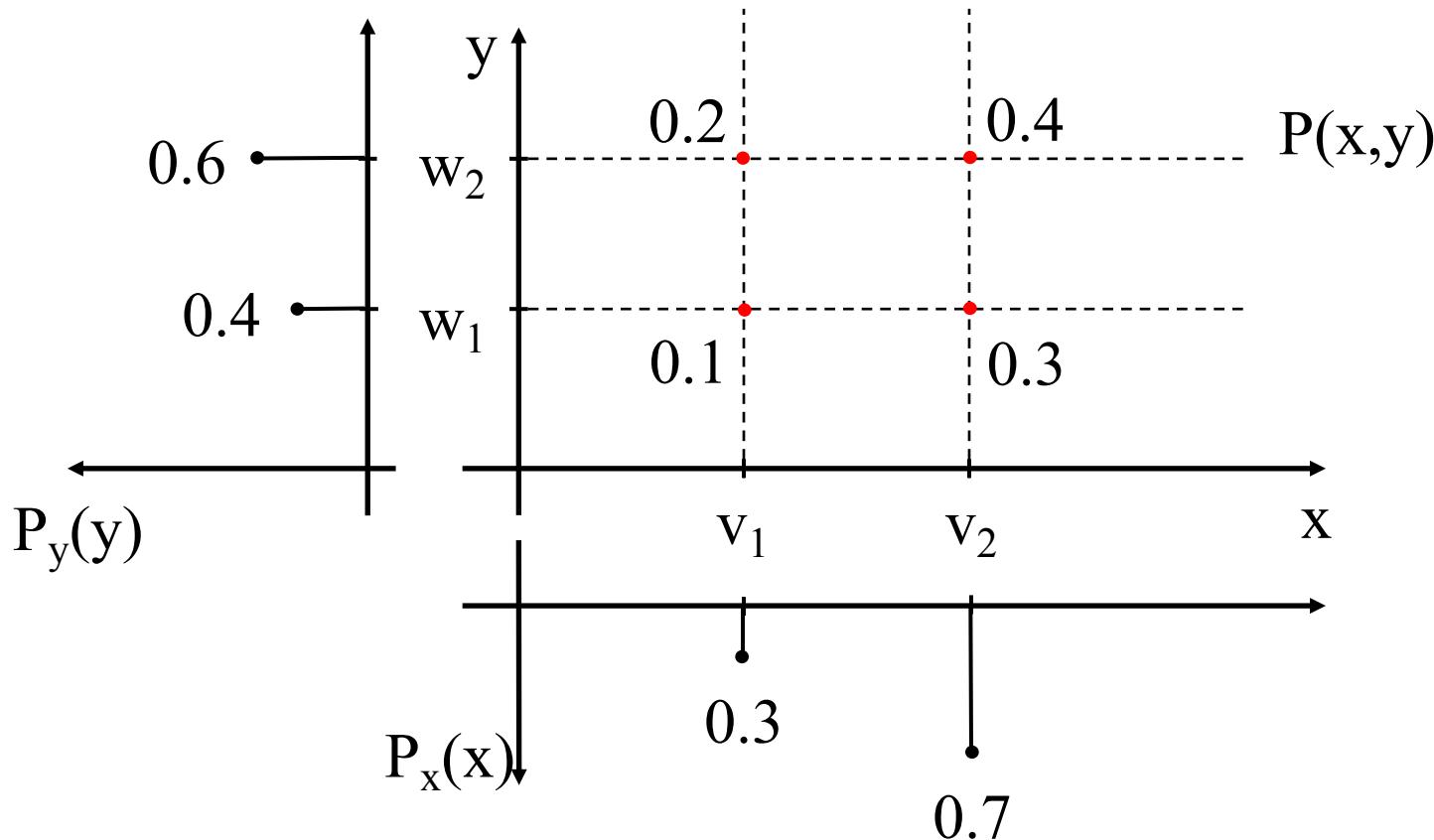
1.4 Marginal distributions

Example: Marginal distribution of y



1.4 Marginal distributions

Example: Marginal distribution of x and y



1.5 Statistical independence

Variables x and y are said to be **statistically independent** if and only if:

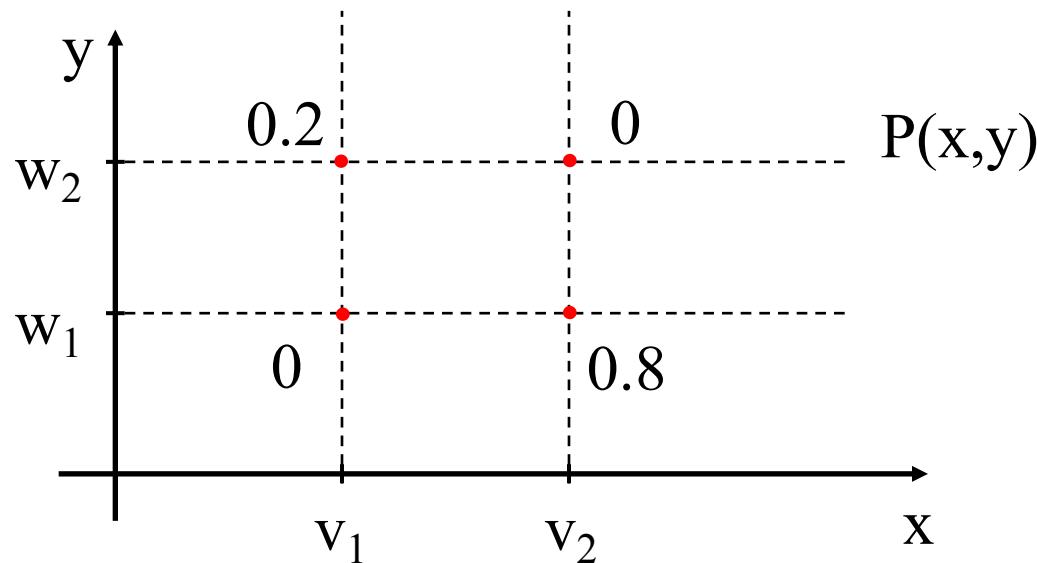
$$P(x, y) = P_x(x) P_y(y)$$

Interpretation:

Knowing the value of x does not give us additional knowledge about the possible values of y and vice versa.

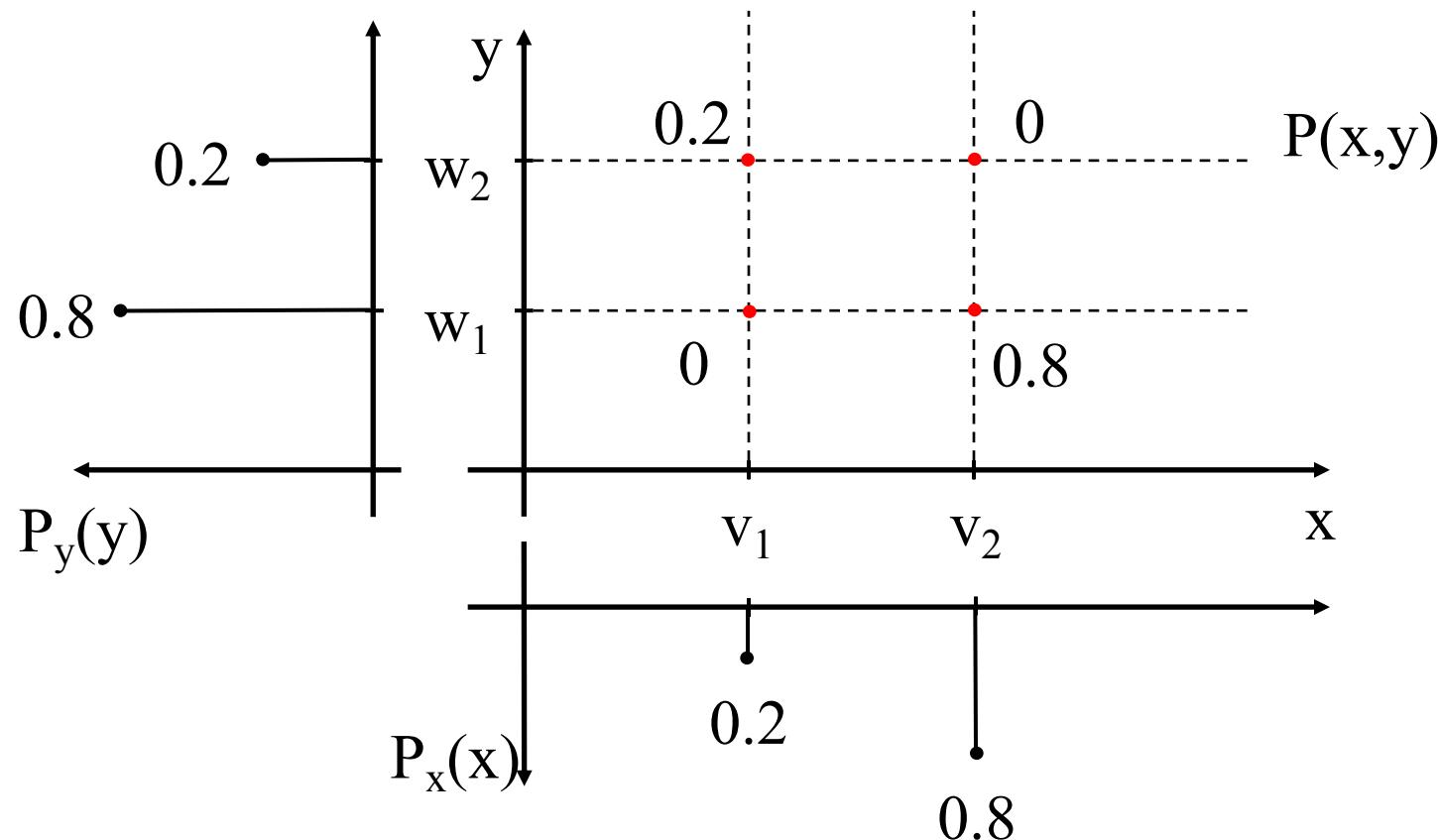
1.5 Statistical independence

Question: Are x and y statistically independent in this example?



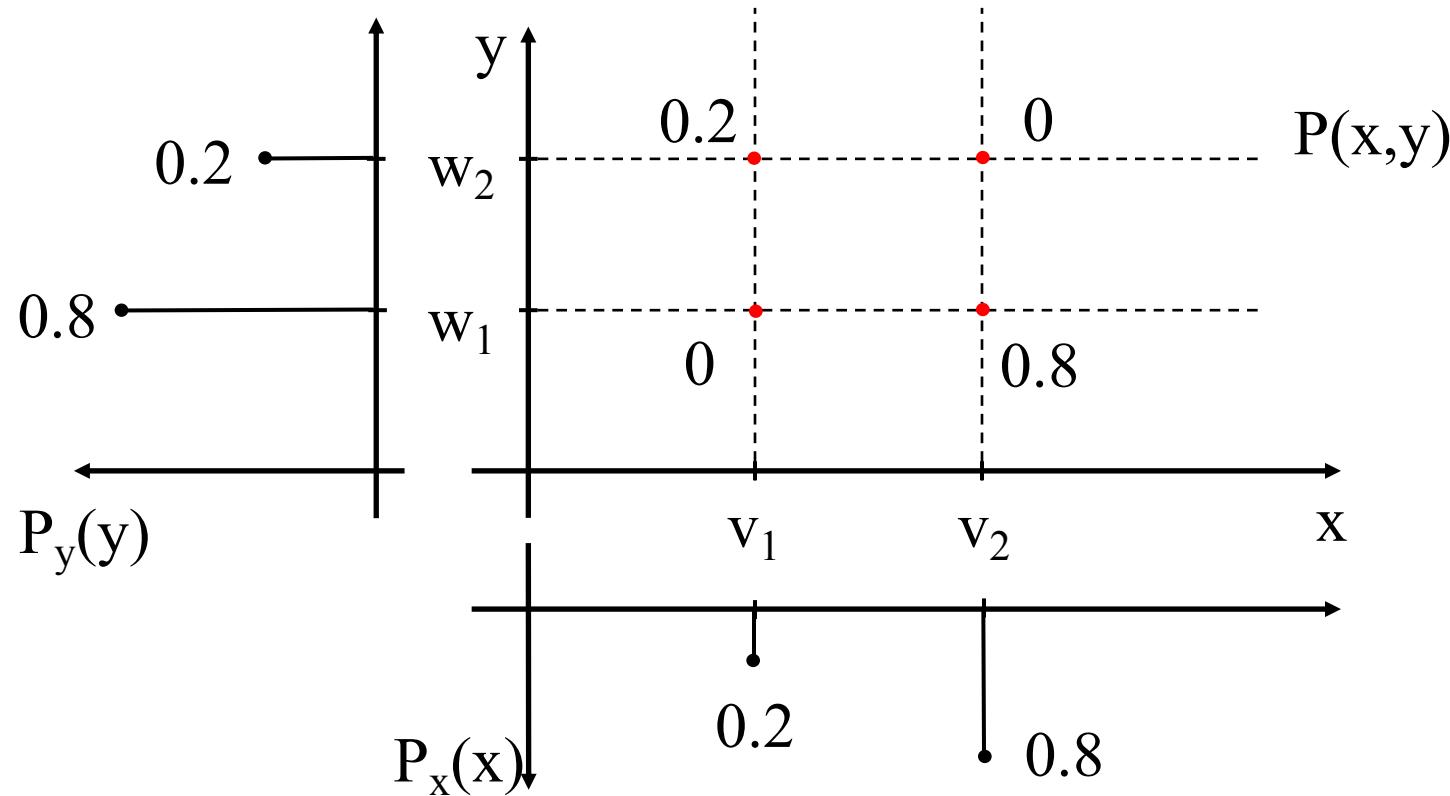
1.5 Statistical independence

Marginal distributions



1.5 Statistical independence

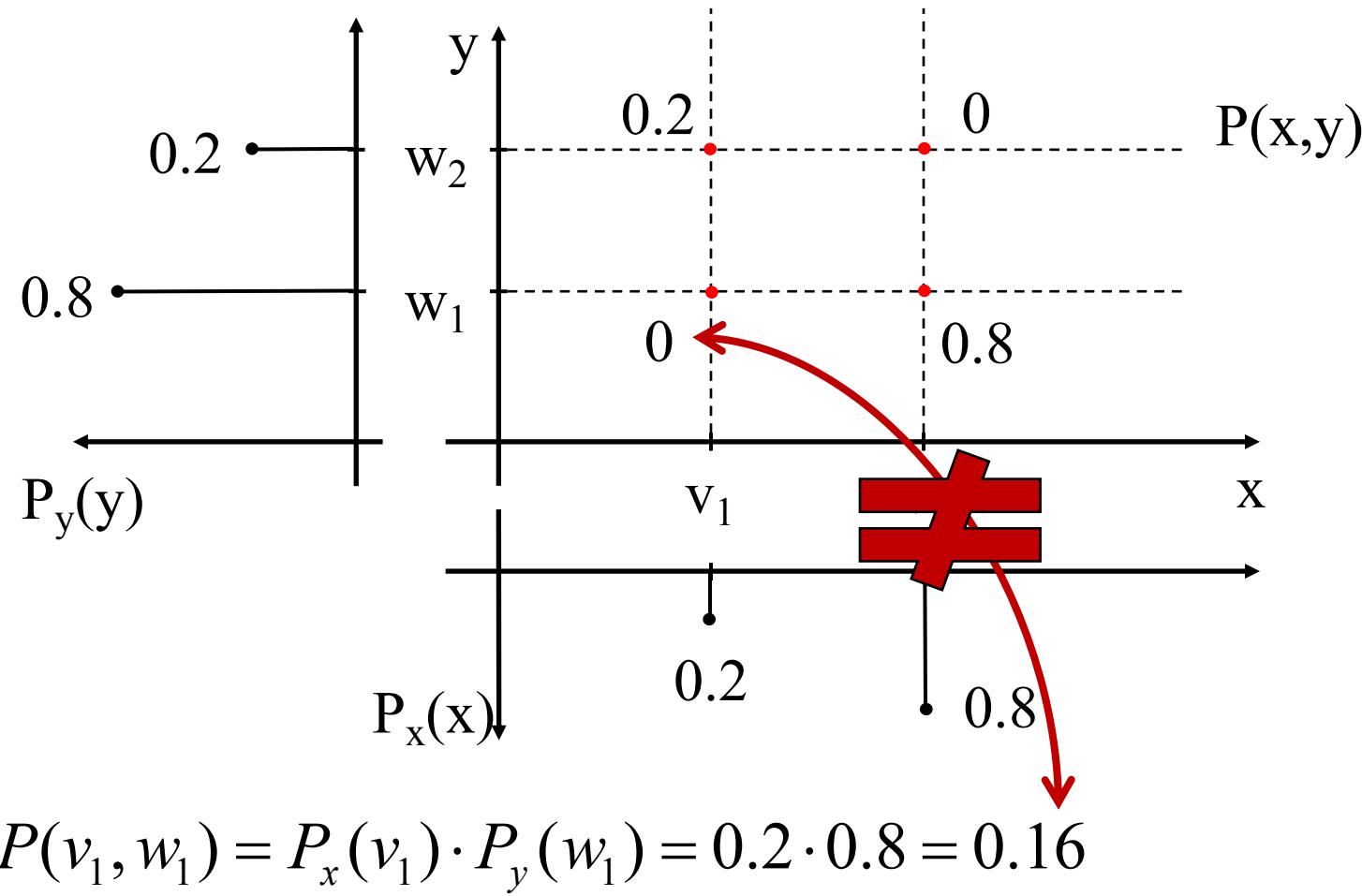
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_1, w_1) = P_x(v_1) \cdot P_y(w_1) = 0.2 \cdot 0.8 = 0.16$$

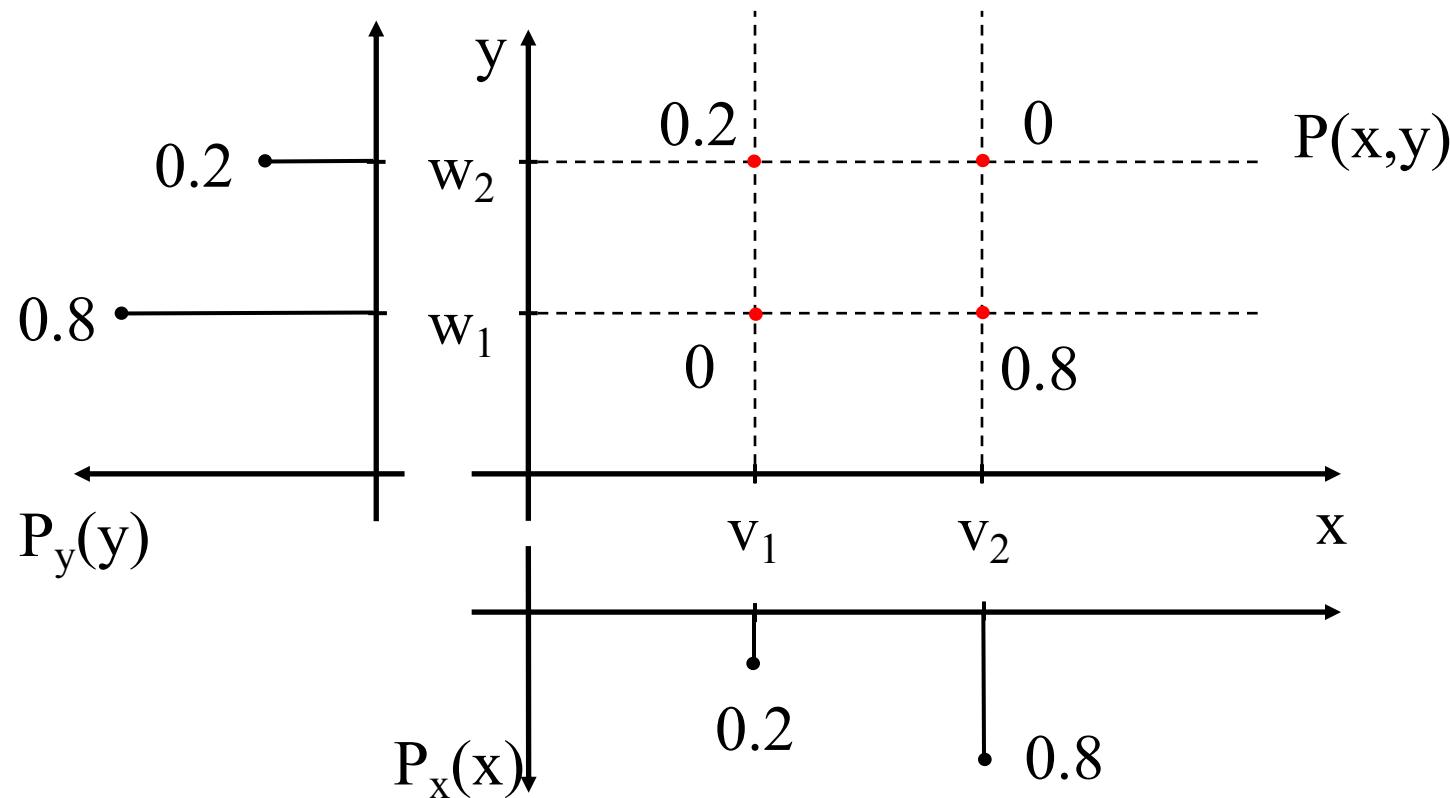
1.5 Statistical independence

In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



1.5 Statistical independence

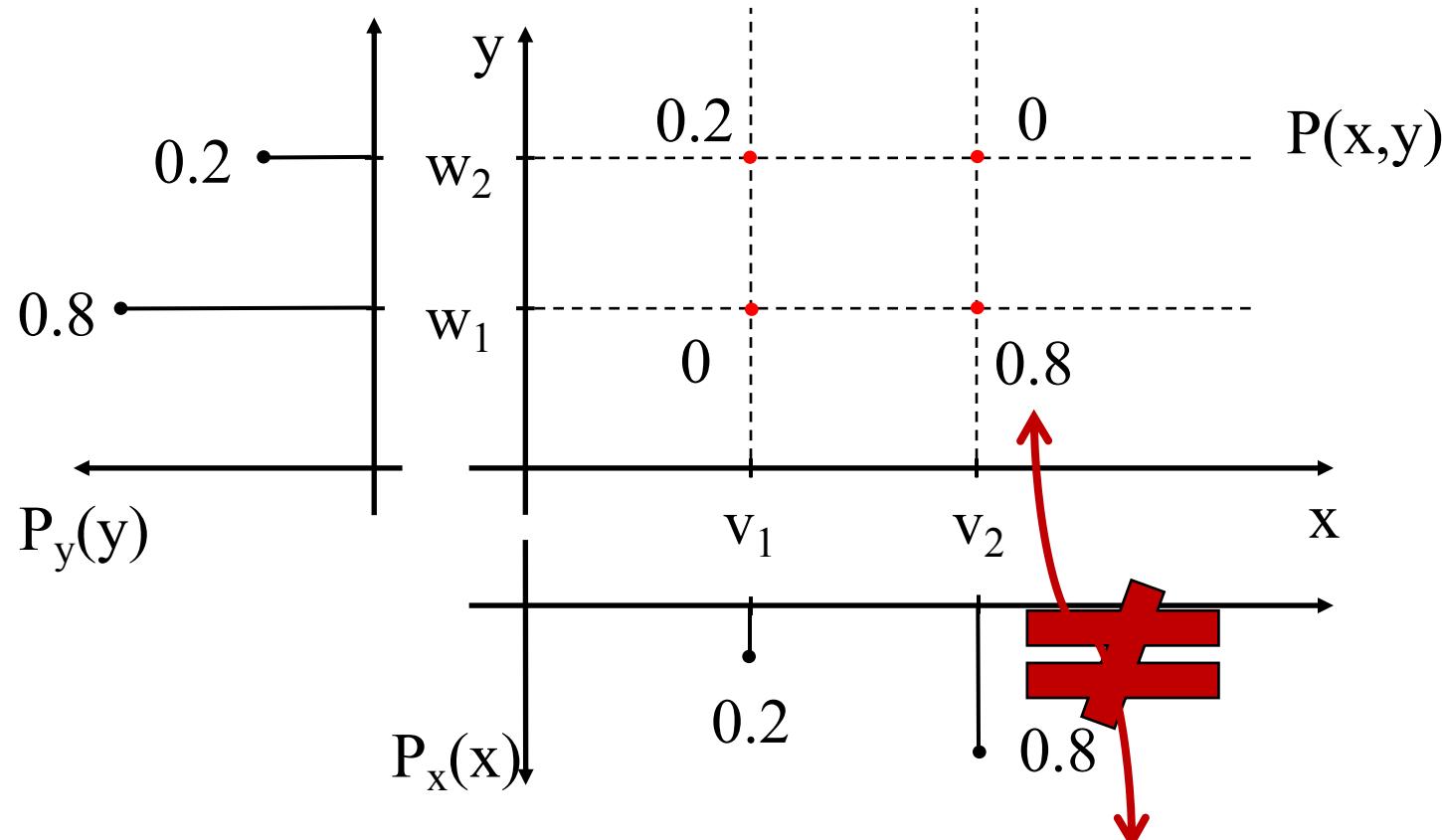
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_2, w_1) = P_x(v_2) \cdot P_y(w_1) = 0.8 \cdot 0.8 = 0.64$$

1.5 Statistical independence

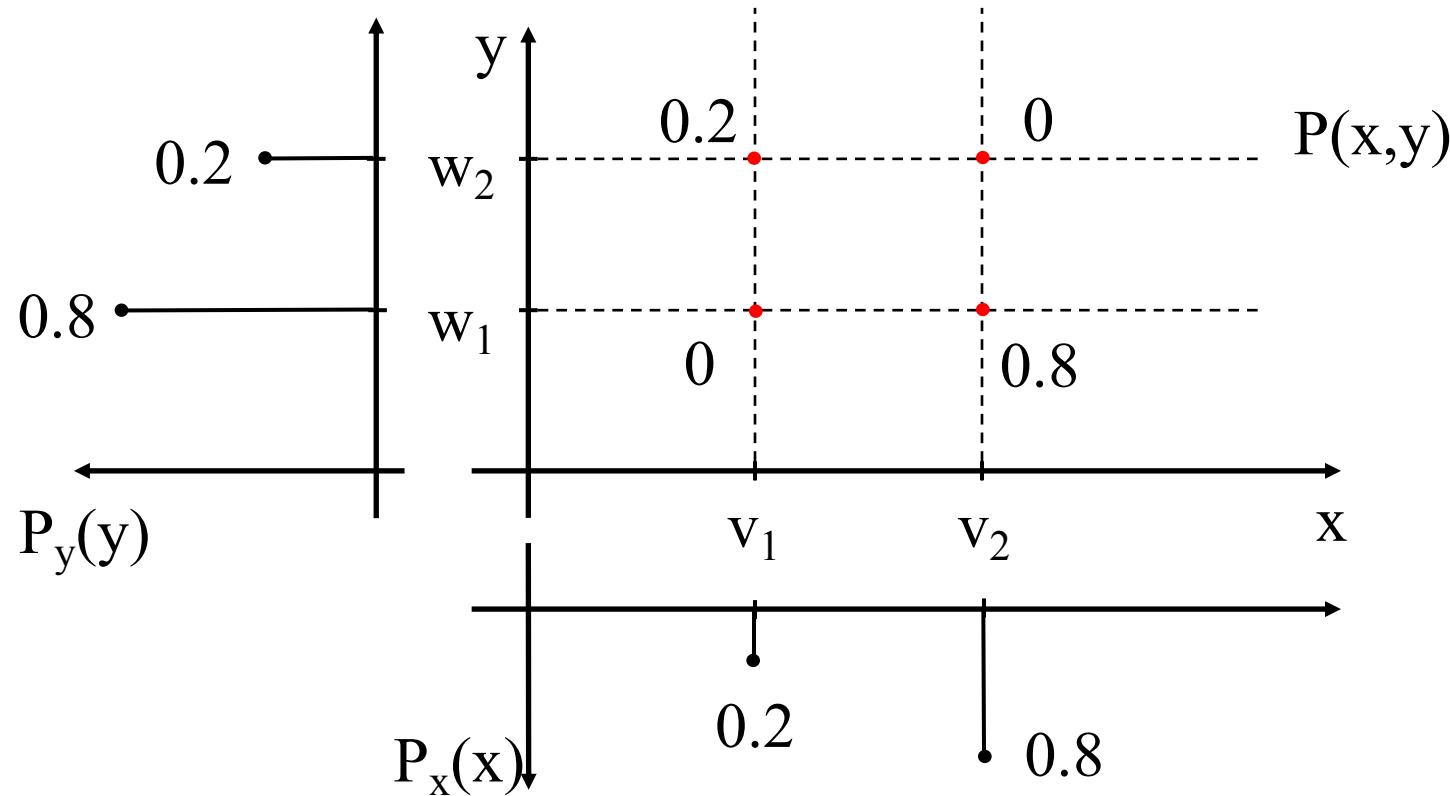
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_2, w_1) = P_x(v_2) \cdot P_y(w_1) = 0.8 \cdot 0.8 = 0.64$$

1.5 Statistical independence

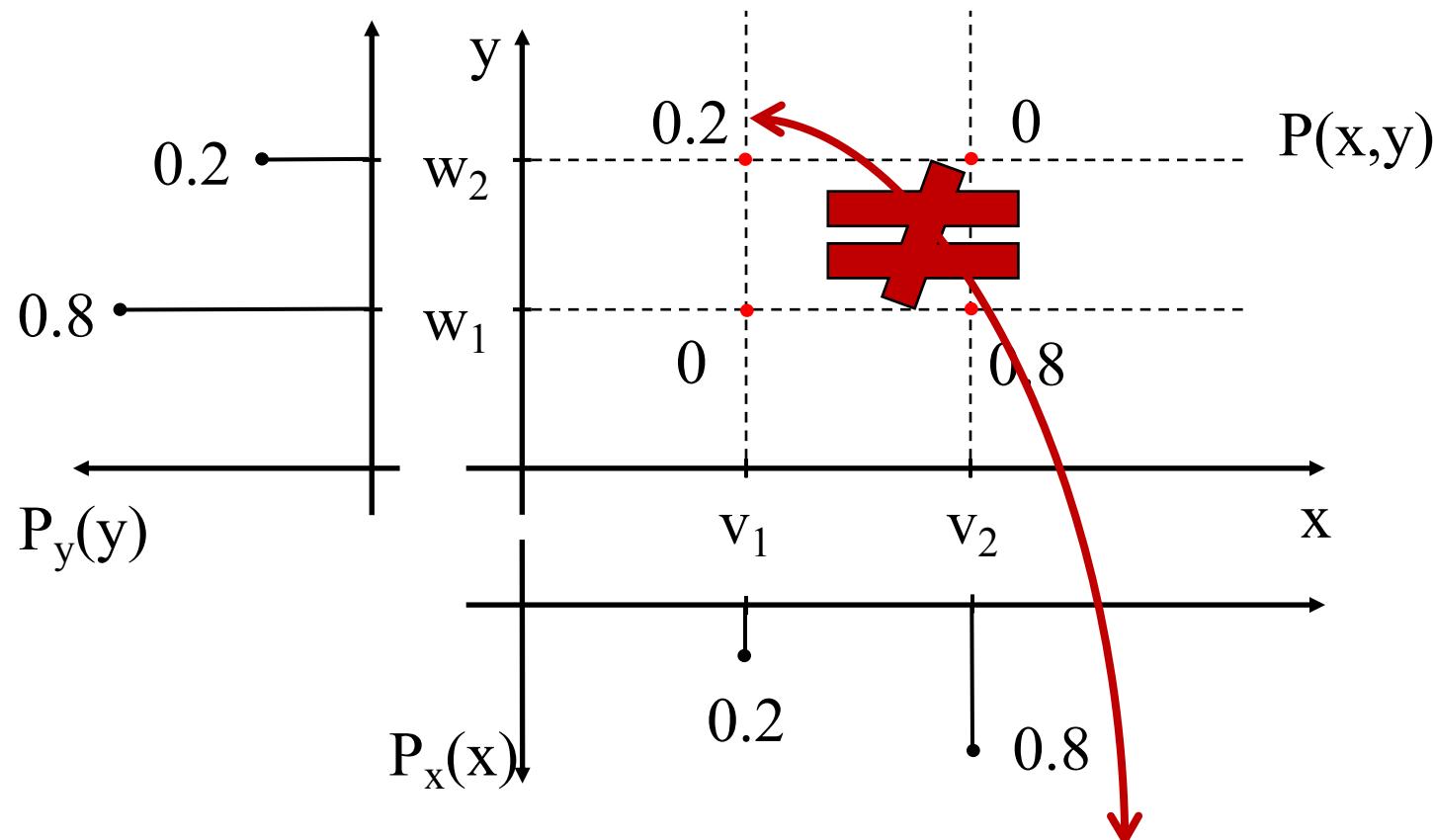
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_1, w_2) = P_x(v_1) \cdot P_y(w_2) = 0.2 \cdot 0.2 = 0.04$$

1.5 Statistical independence

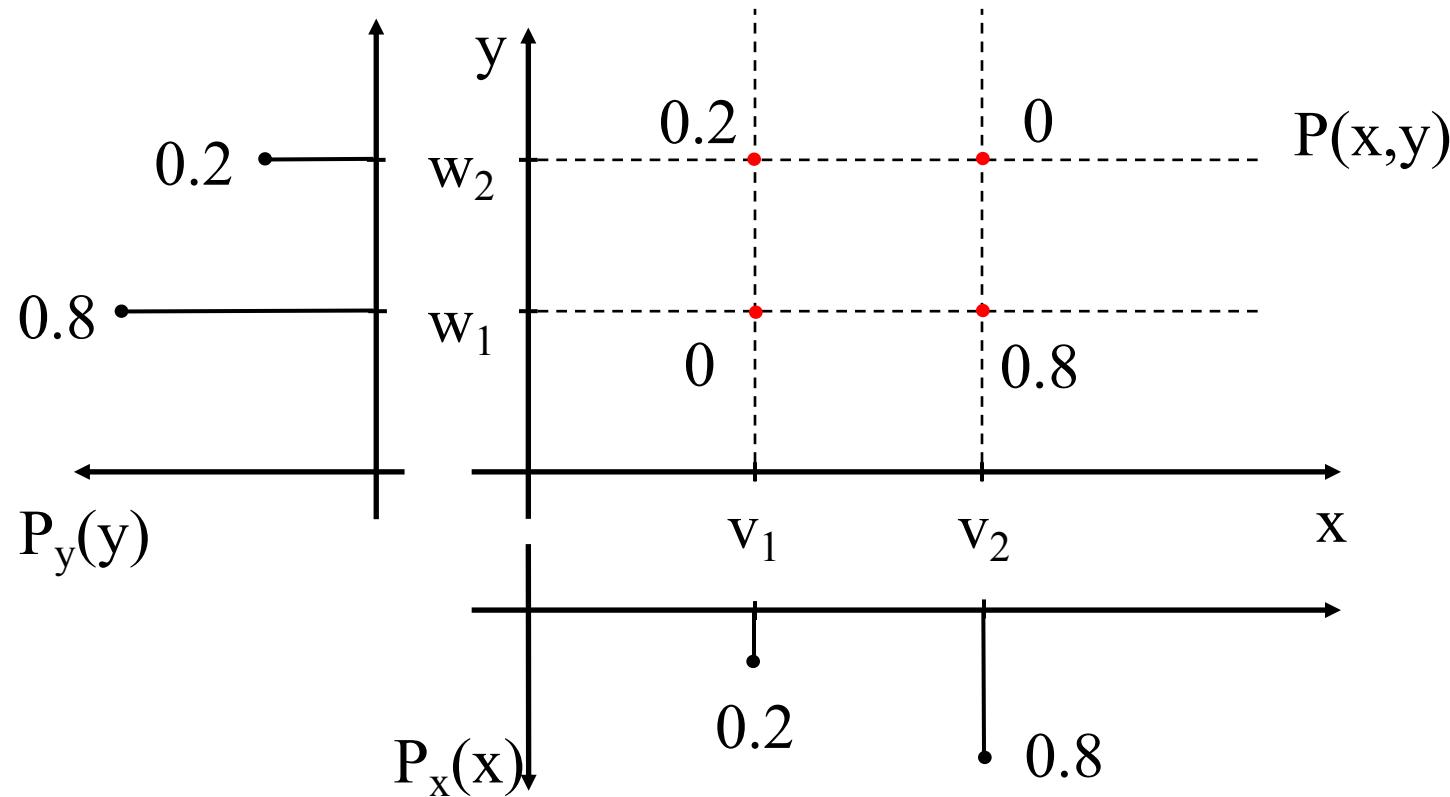
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_1, w_2) = P_x(v_1) \cdot P_y(w_2) = 0.2 \cdot 0.2 = 0.04$$

1.5 Statistical independence

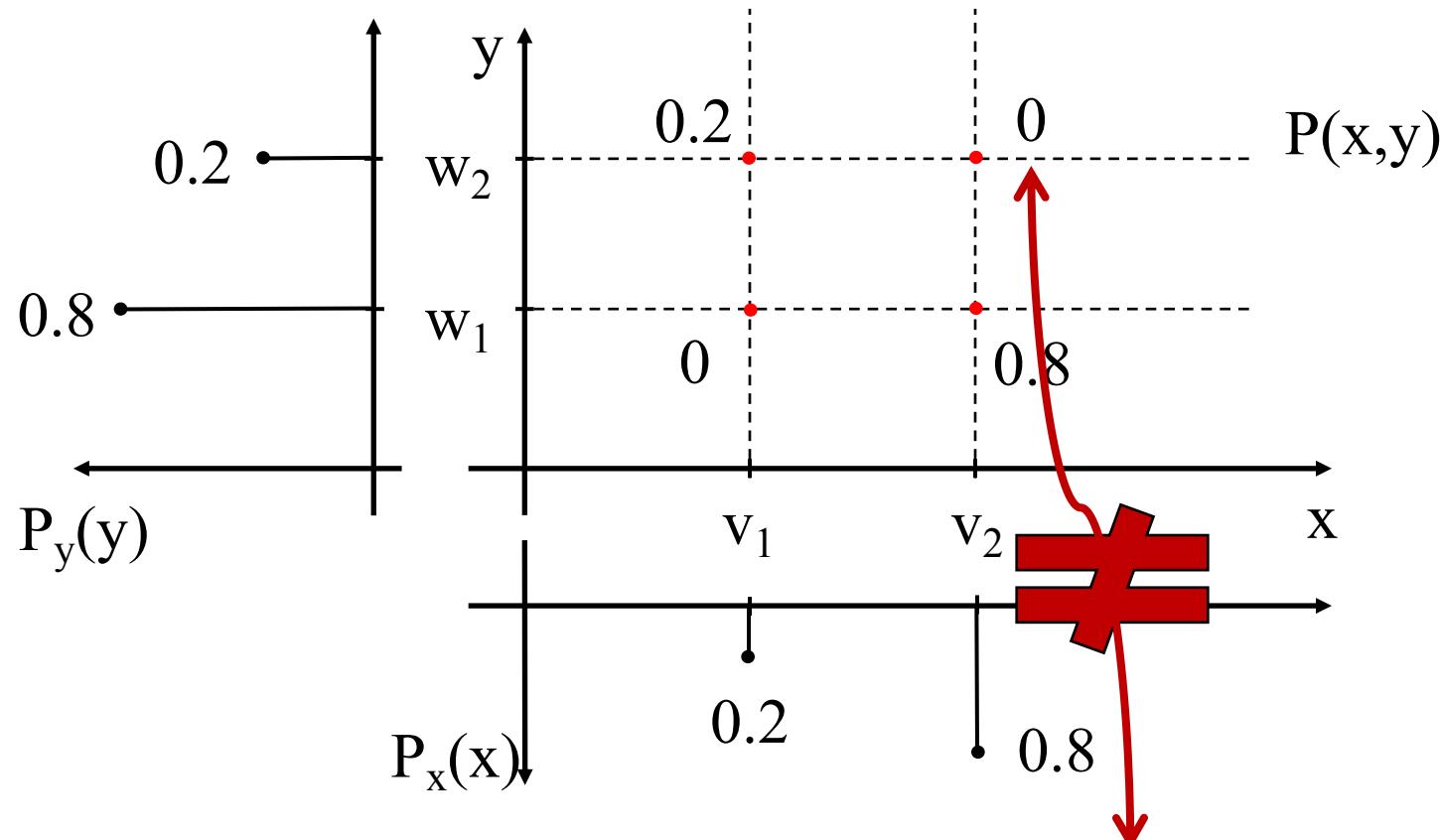
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_2, w_2) = P_x(v_2) \cdot P_y(w_2) = 0.8 \cdot 0.2 = 0.16$$

1.5 Statistical independence

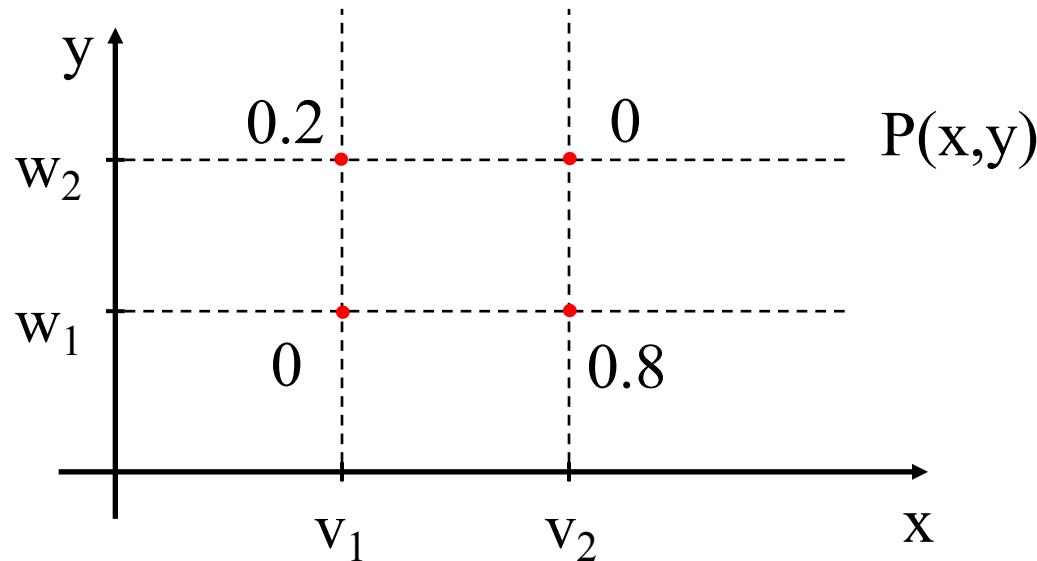
In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



$$P(v_2, w_2) = P_x(v_2) \cdot P_y(w_2) = 0.8 \cdot 0.2 = 0.16$$

1.5 Statistical independence

In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$

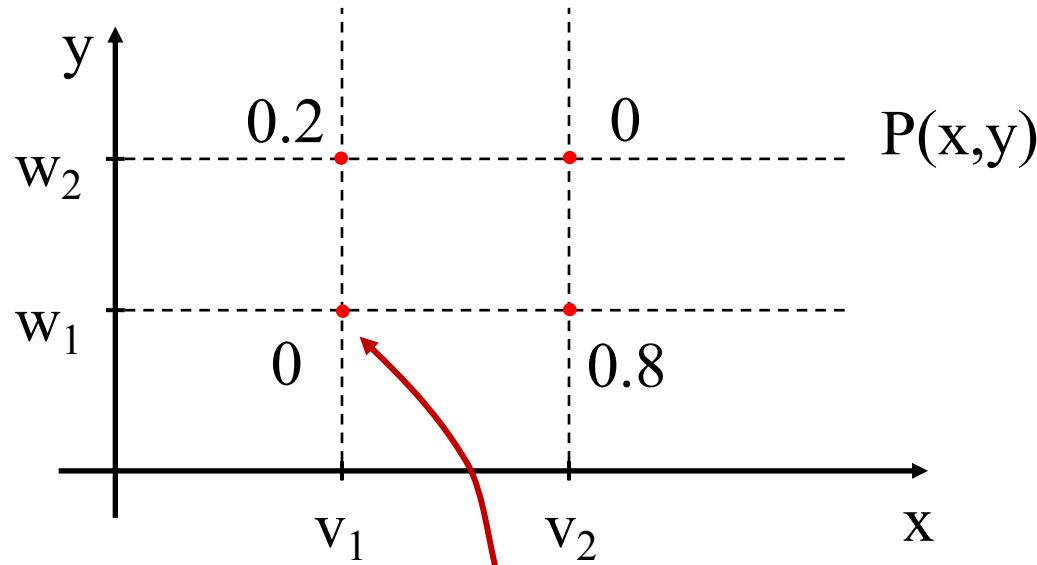


→ x and y are statistically dependent

Knowing the value of x lets us get a better estimate of the values of y and vice versa.

1.5 Statistical independence

In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



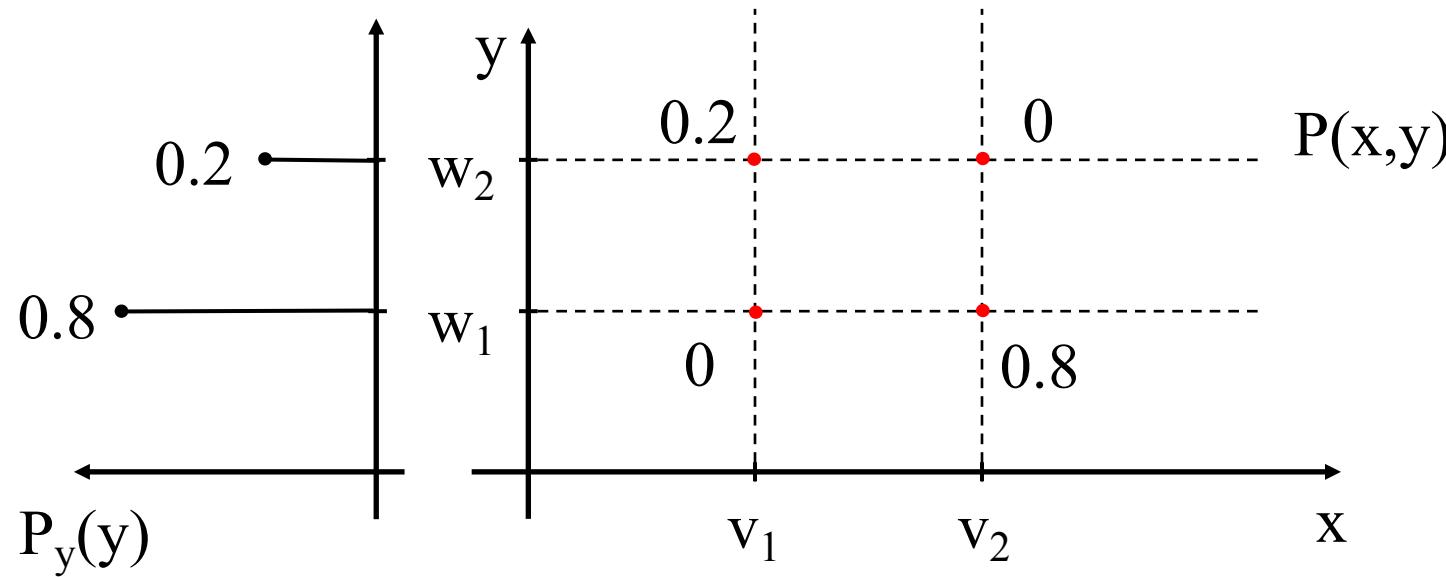
Example:

Let $x = v_1 \rightarrow$ gives us the knowledge, that $y=w_1$ is **impossible!**

Let $x = v_2 \rightarrow$ gives us the knowledge, that $y=w_2$ is **impossible!**

1.5 Statistical independence

In case of statistical independence: $P(x, y) = P_x(x) P_y(y)$



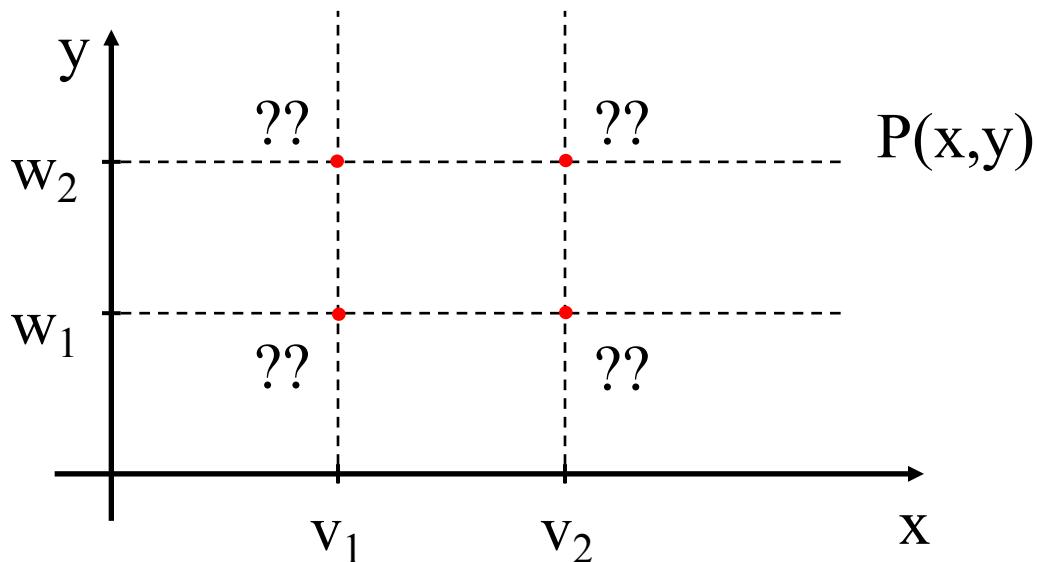
Example:

Without knowledge about the value of x we must use the marginal distribution $P_y(y)$

1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

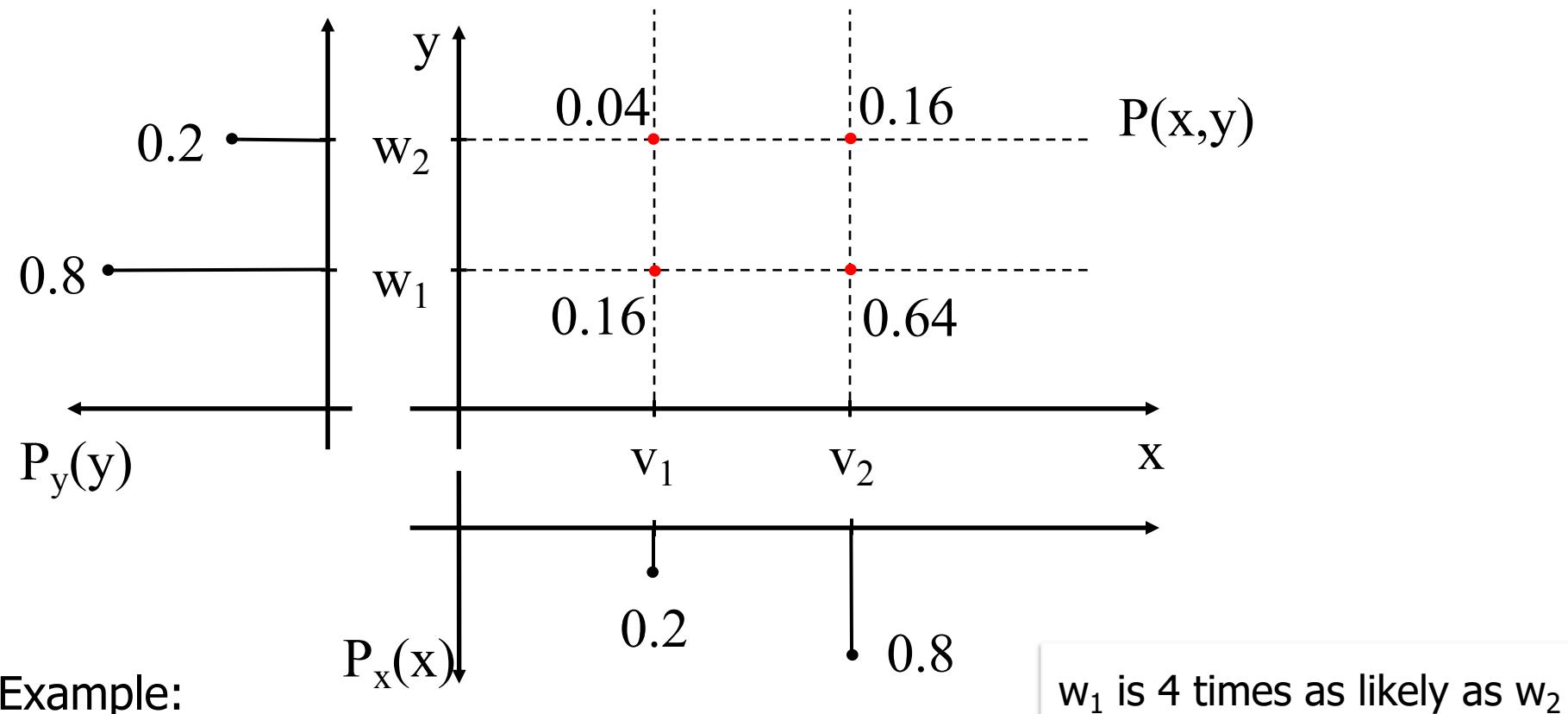
How could a distribution look like in case of statistical independence?



1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

How could a distribution look like in case of statistical independence?

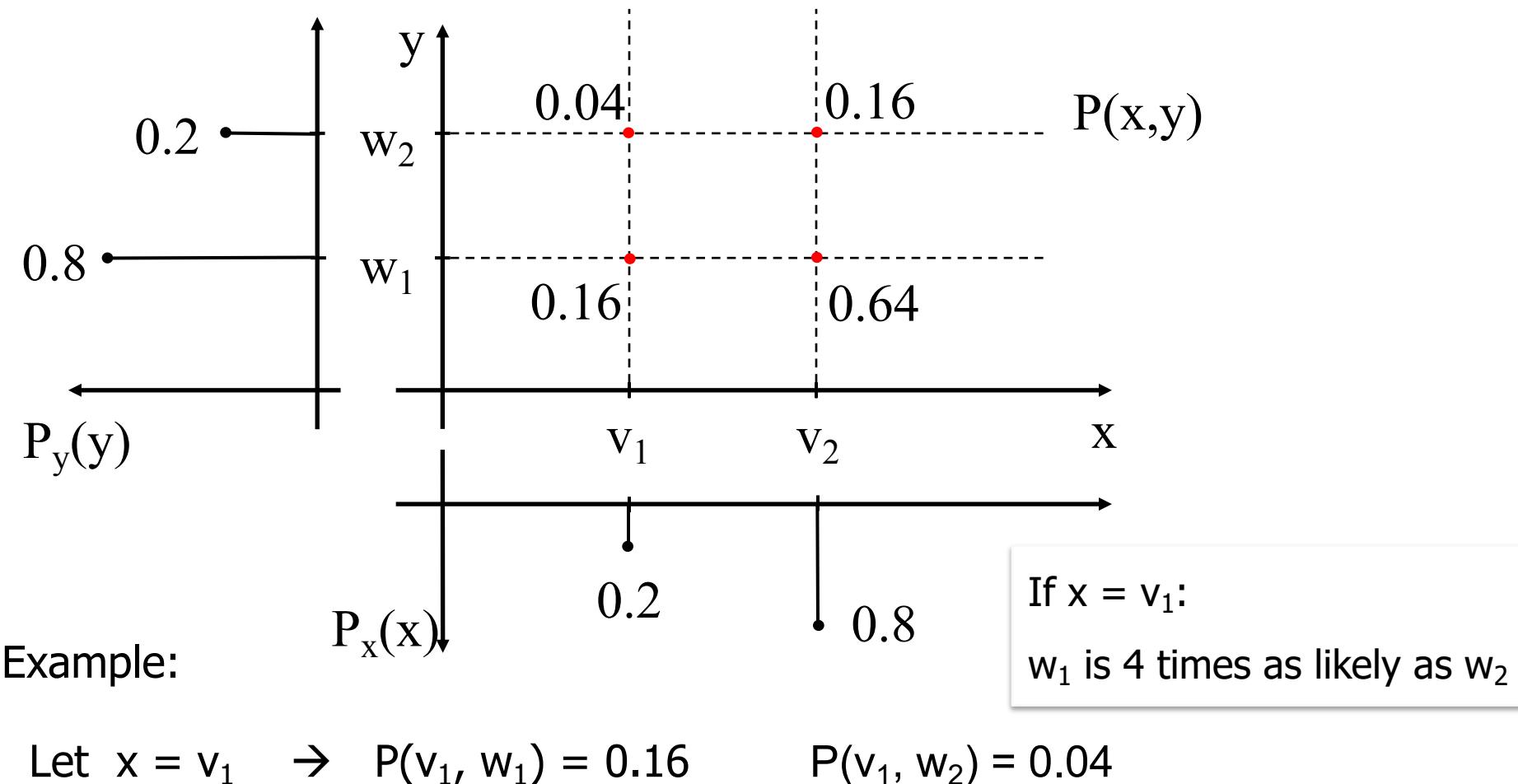


Without knowledge of x value we know that: $P(w_1) = 4 \cdot P(w_2)$

1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

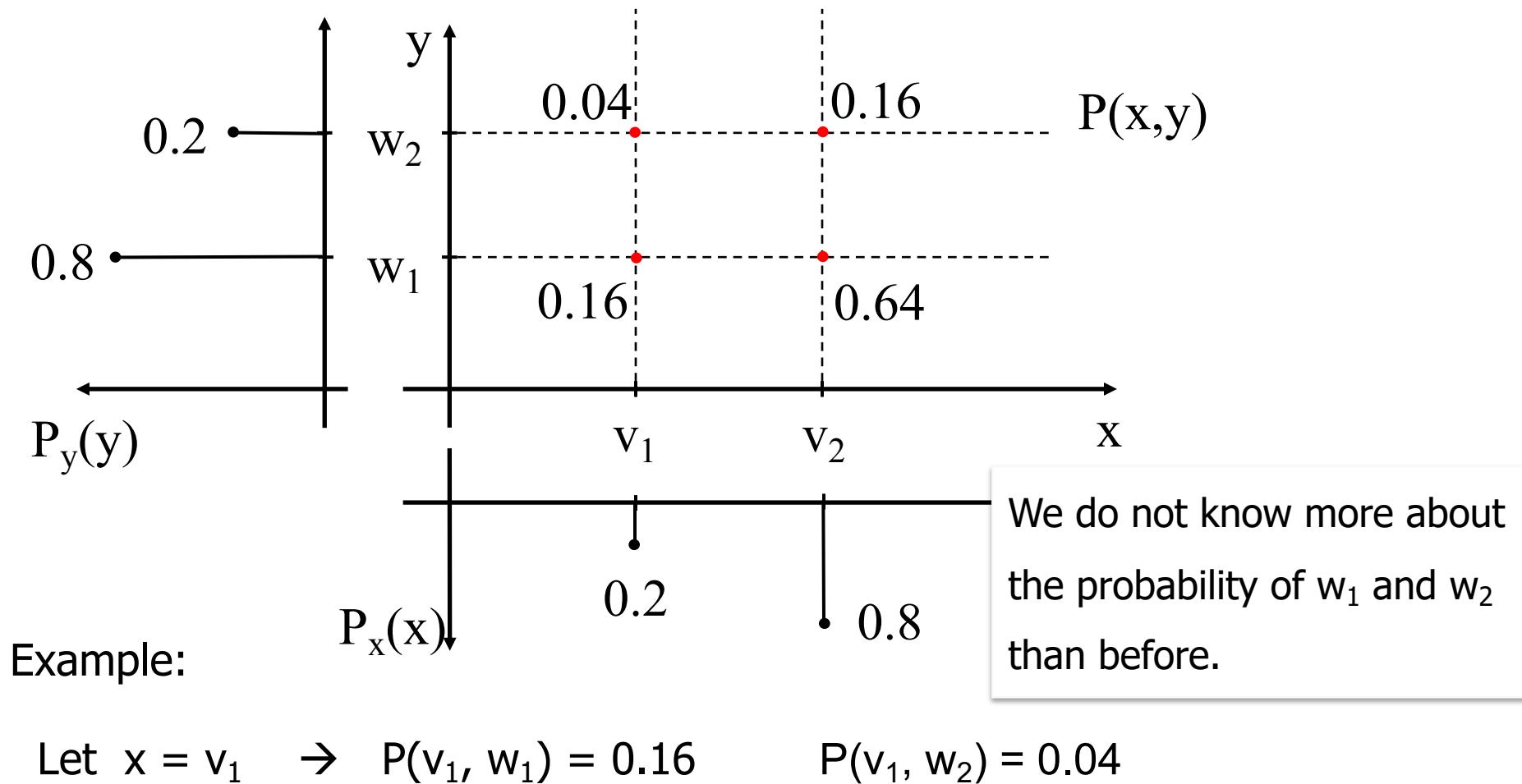
How could a distribution look like in case of statistical independence?



1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

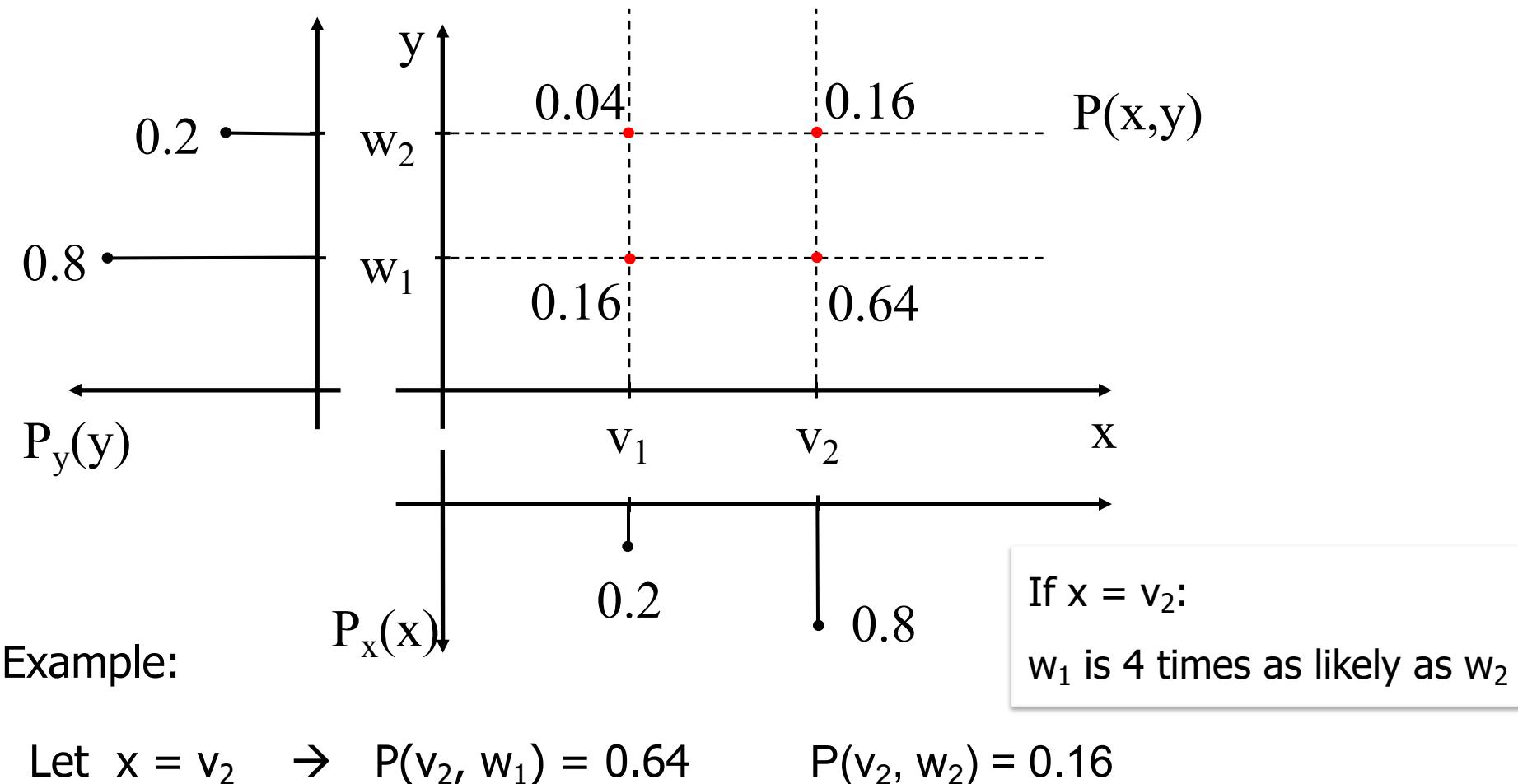
How could a distribution look like in case of statistical independence?



1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

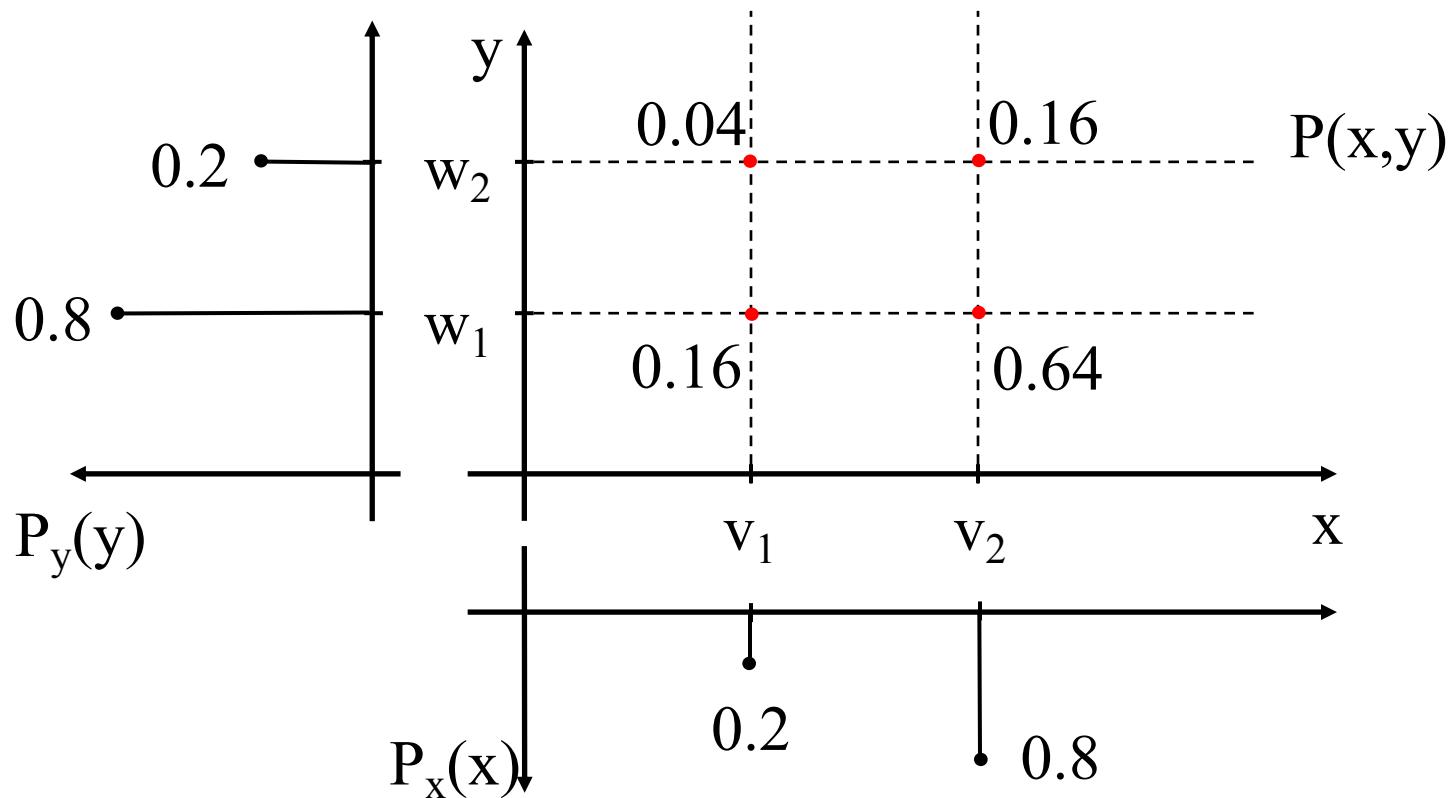
How could a distribution look like in case of statistical independence?



1.5 Statistical independence

$$P(x, y) = P_x(x) P_y(y)$$

How could a distribution look like in case of statistical independence?



Knowledge about x does **not** give us additional information about y and vice versa.

1.6 Expected Values of Functions of Two Variables

The expected value of a **function $f(x, y)$** of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Interpretation 1:

A weighting of the function values at point (x, y) with the probability that this combination of x and y really occurs.

1.6 Expected Values of Functions of Two Variables

The expected value of a **function $f(x, y)$** of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y)P(x, y)$$

Interpretation 2:

Function $f(x, y)$ depends upon random variables.

1. Perform experiment \rightarrow get one observation for x and y
2. Enter values into function and observe corresponding function value
3. Repeat a large number of times, compute mean observed function value
 \rightarrow expected value of function $f(x, y)$

1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y)P(x, y)$$

Remark: Result is one scalar expectation value!

Usual notation is somewhat unclear:

- Here, x and y are random variables
 - Placeholder to show that f depends upon x and y
 - Do not enter concrete values for x and y here
- Here, **concrete values** for x and y are used to determine the result

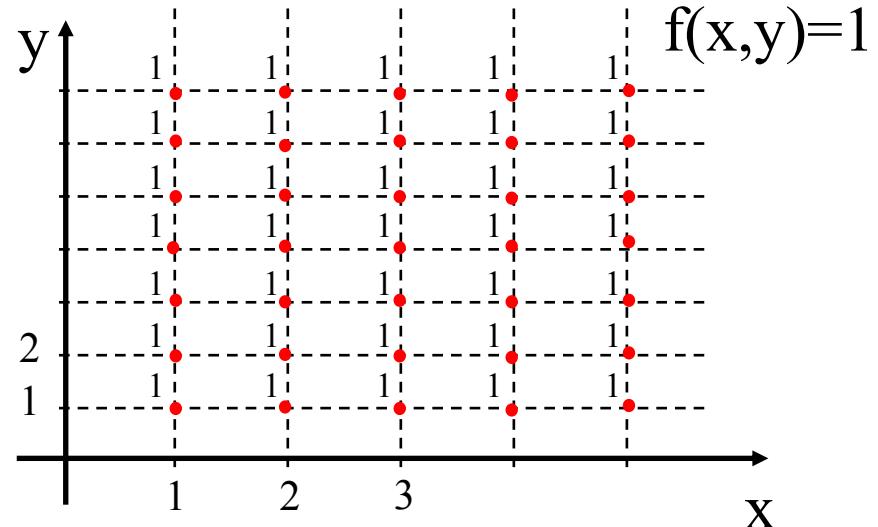
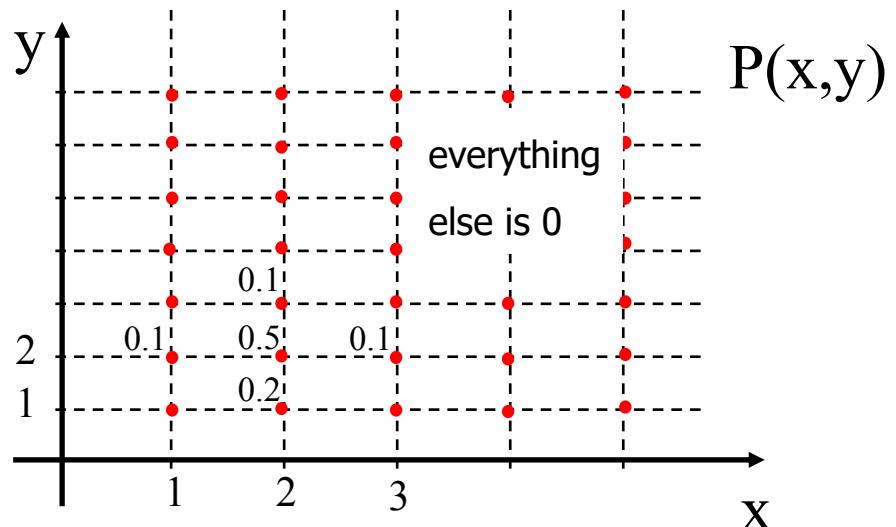
1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

Which value do you expect for function $f(x, y)$ in this case?



1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

$$E[f(x, y)] = 1 \cdot 0.1 + 1 \cdot 0.2 + 1 \cdot 0.5 + 1 \cdot 0.1 + 1 \cdot 0.1 = 1$$

→ Reasonable, since the function has always value 1. Therefore, the probability of (x, y) does not have an influence on the expected value.

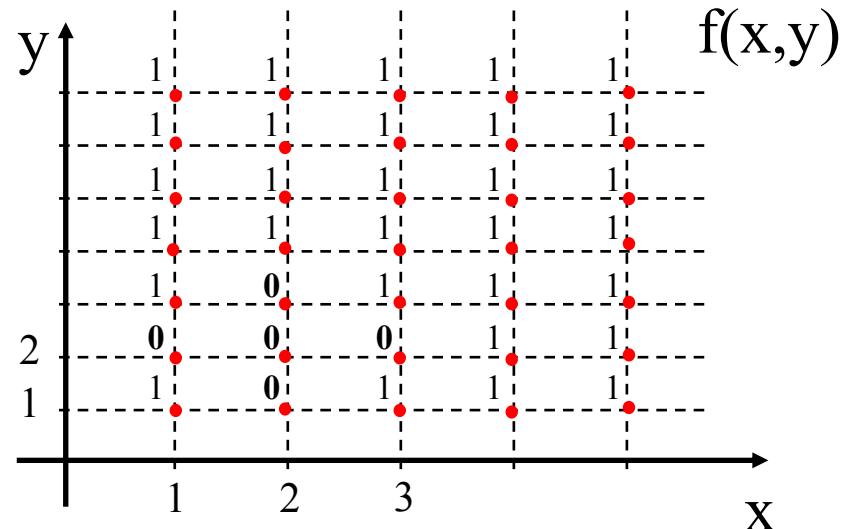
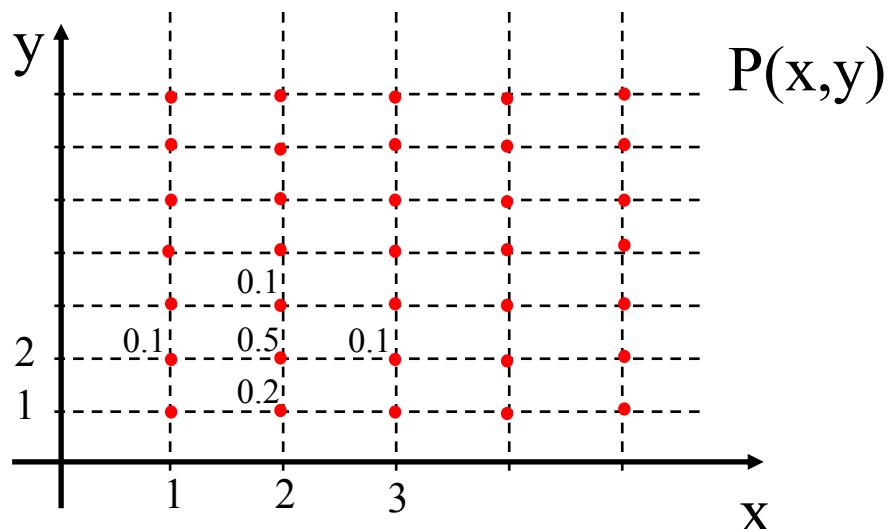
1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

Which value do you expect for function $f(x, y)$ in this case?



1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

$$E [f(x, y)] = 1 \cdot 0 + \dots + 1 \cdot 0 + 0 \cdot 0.1 + 0 \cdot 0.2 + 0 \cdot 0.5 + 0 \cdot 0.1 + 0 \cdot 0.1 = 0$$

- For all **possible** combinations of x and y the function value is 0.
- For impossible combinations, the function value does not matter.

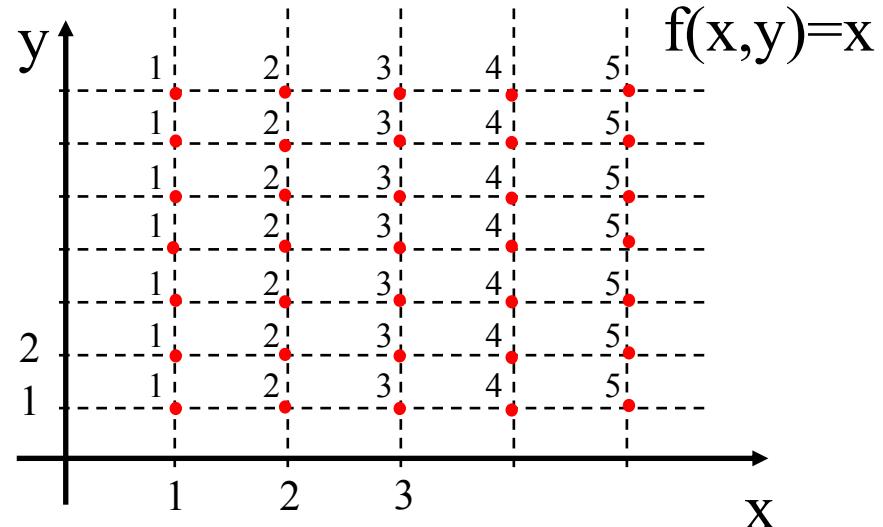
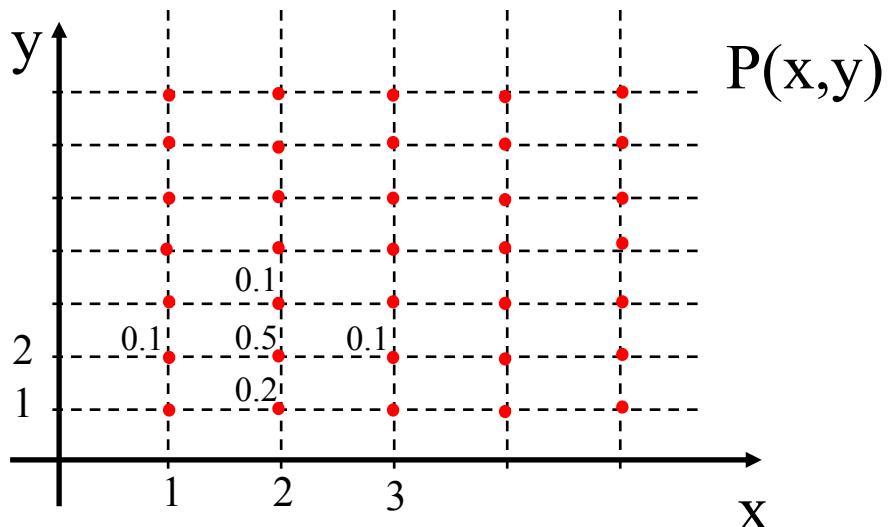
1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

Which value do you expect for function $f(x, y)$ in this case?



1.6 Expected Values of Functions of Two Variables

The expected value of a function $f(x, y)$ of two random variables is given by:

$$E[f(x, y)] = \sum_{x \in X} \sum_{y \in Y} f(x, y) P(x, y)$$

Example

$$E [f(x, y)] = 1 \cdot 0.1 + 2 \cdot 0.2 + 2 \cdot 0.5 + 2 \cdot 0.1 + 3 \cdot 0.1 = 2$$

- Expected function value 'in x-direction'
- Mean or *first moment*

1.6 Expected Values of Functions of Two Variables

Means (first moments) and variances (second moments) are given by:

$$\mu_x = E[x] = \sum_{x \in X} \sum_{y \in Y} x P(x, y)$$

$$\mu_y = E[y] = \sum_{x \in X} \sum_{y \in Y} y P(x, y)$$

$$\sigma_x^2 = E[(x - \mu_x)^2] = \sum_{x \in X} \sum_{y \in Y} (x - \mu_x)^2 P(x, y)$$

$$\sigma_y^2 = E[(y - \mu_y)^2] = \sum_{x \in X} \sum_{y \in Y} (y - \mu_y)^2 P(x, y)$$

1.6 Expected Values of Functions of Two Variables

Means (first moments) and variances (second moments) are given by:

How strong does a random variable
deviate from its mean value.

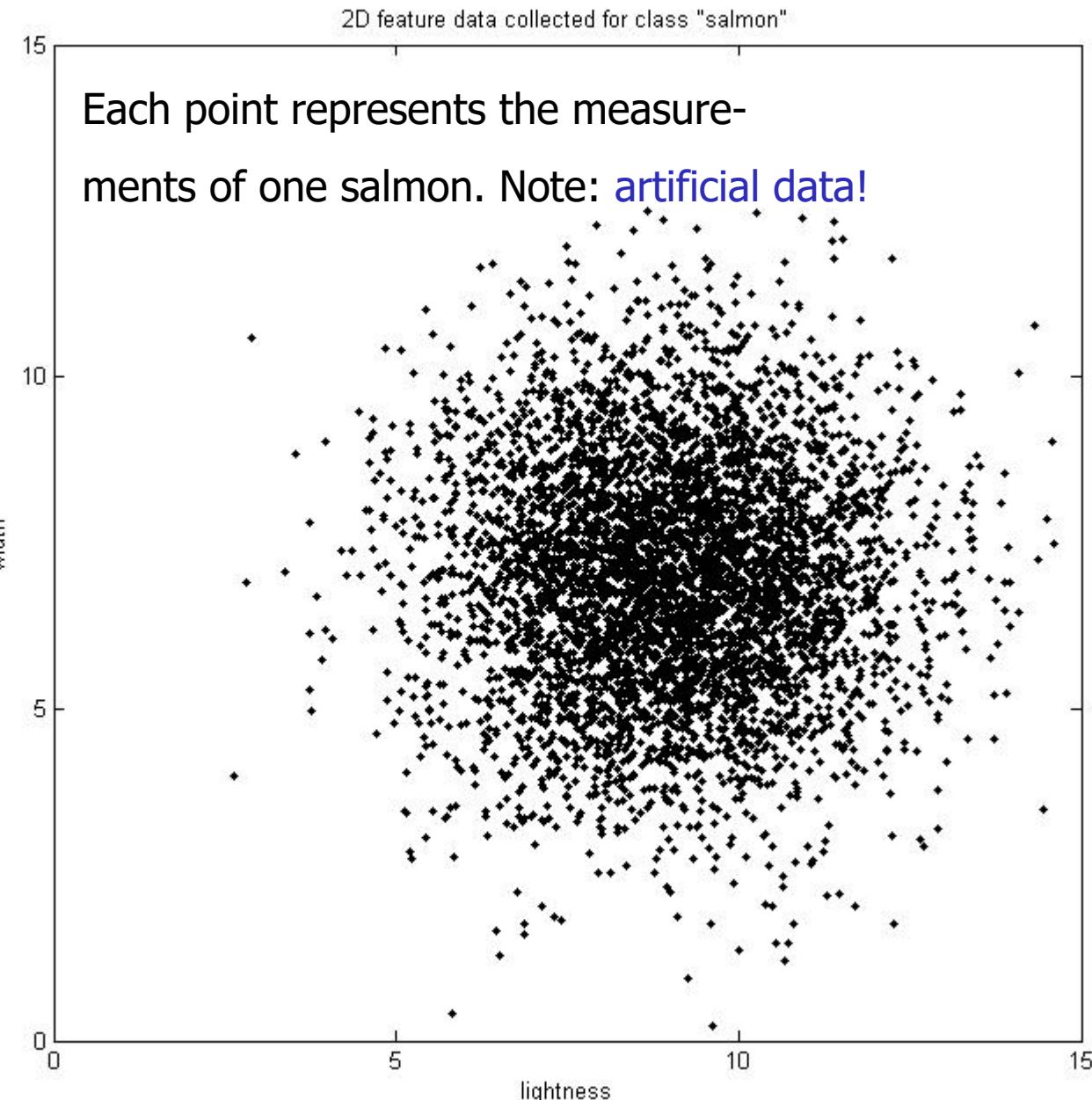
$$\sigma_x^2 = E[(x - \mu_x)^2] = \sum_{x \in X} \sum_{y \in Y} (x - \mu_x)^2 P(x, y)$$

$$\sigma_y^2 = E[(y - \mu_y)^2] = \sum_{x \in X} \sum_{y \in Y} (y - \mu_y)^2 P(x, y)$$

Practical example:

x = lightness

y = width

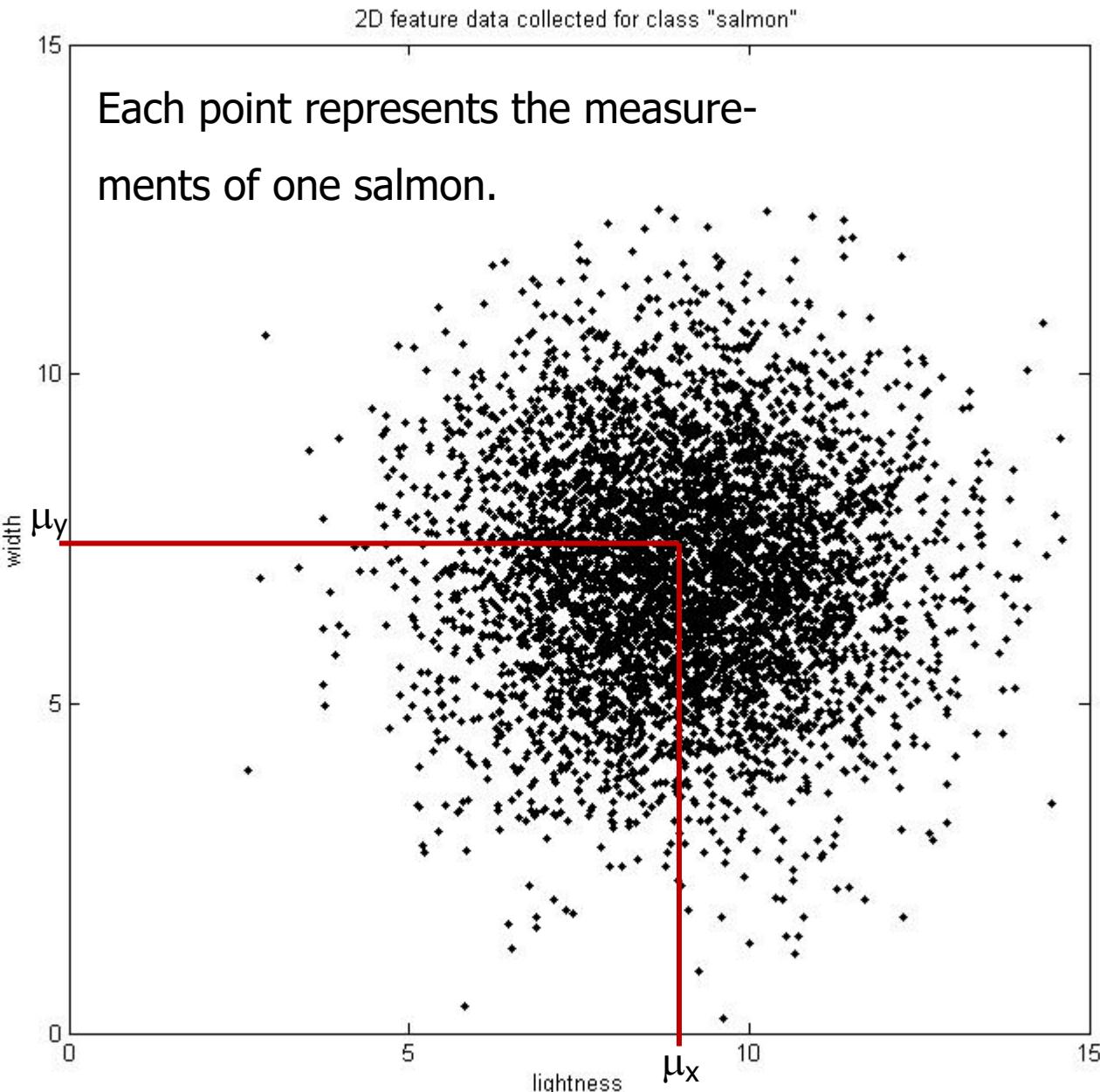


Practical example:

x = lightness

y = width

Means: μ_x, μ_y

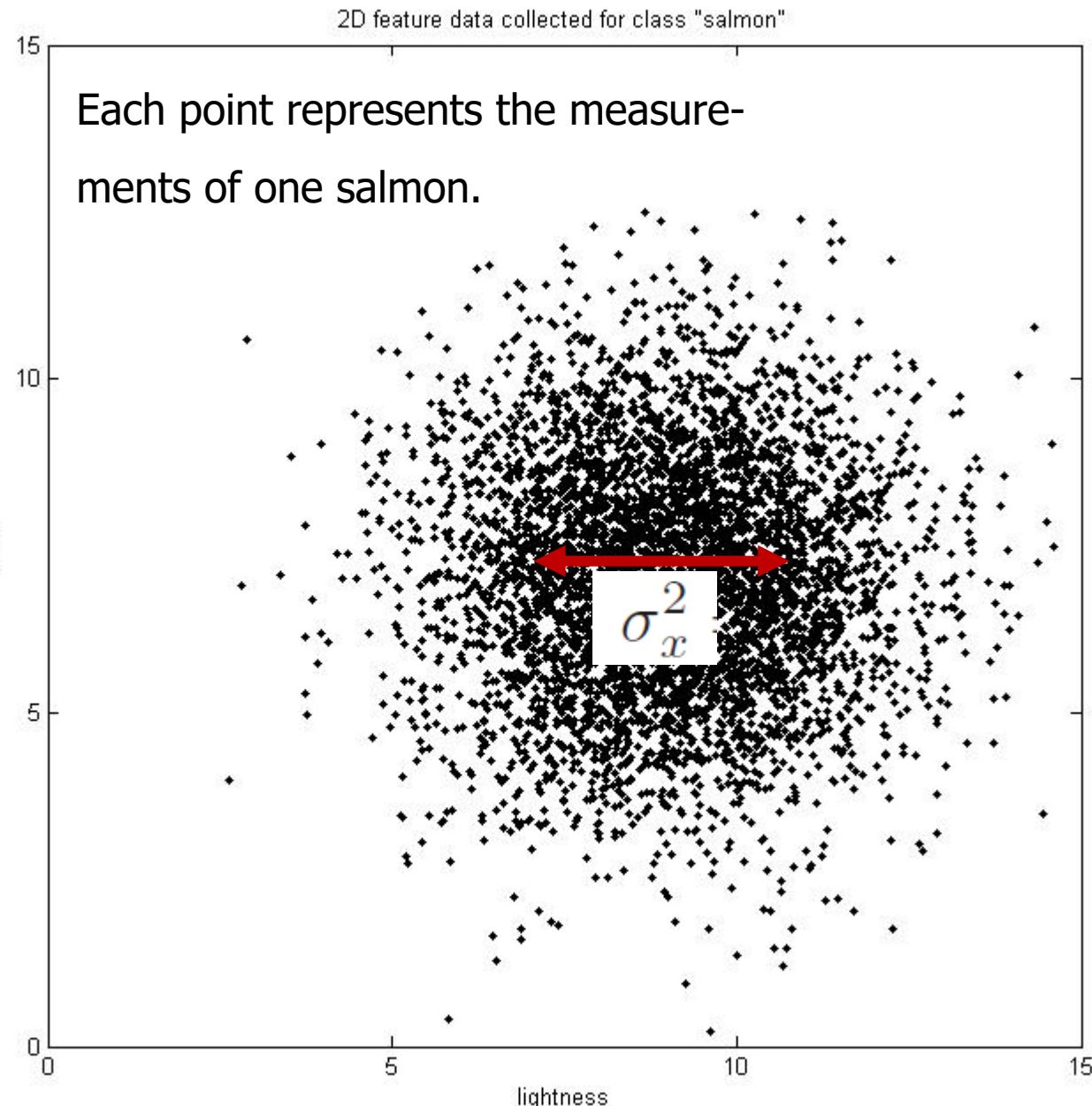


Practical example:

x = lightness

y = width

Variance: σ_x^2



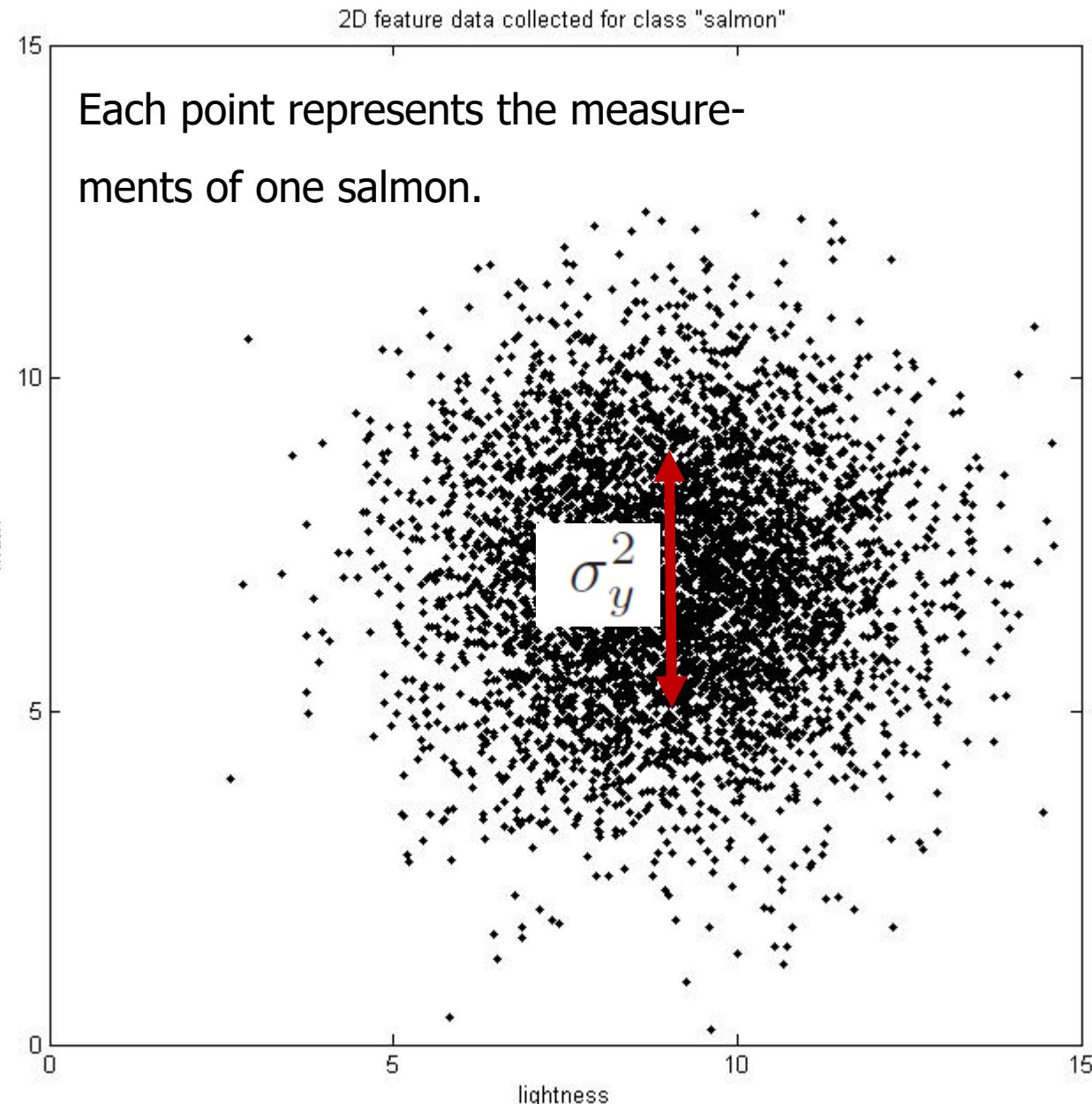
Practical example:

x = lightness

y = width

Variance: σ_y^2

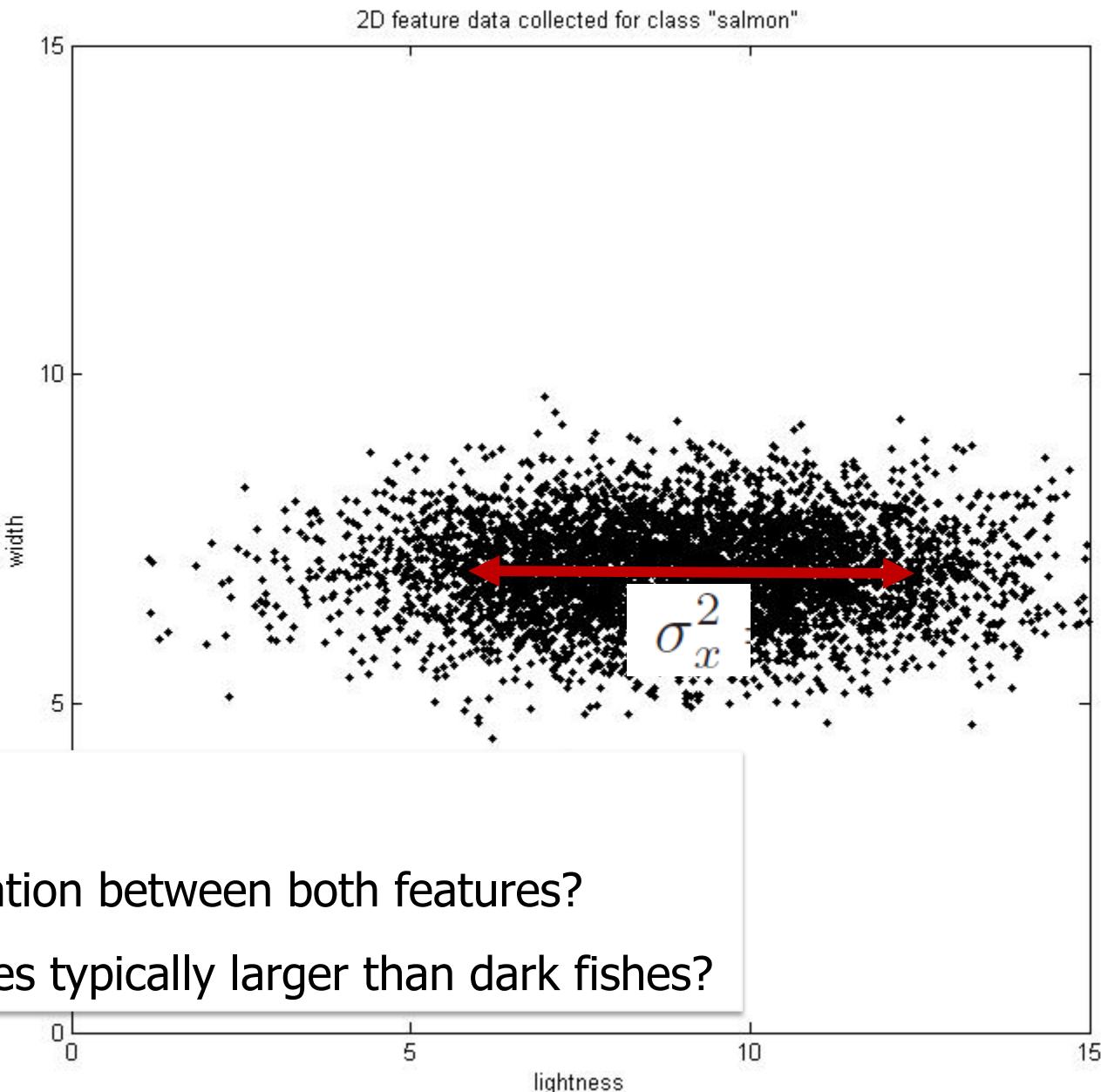
→ Same variance!



Practical example:

Another distribution:

→ stronger variance of
lightness



Question:

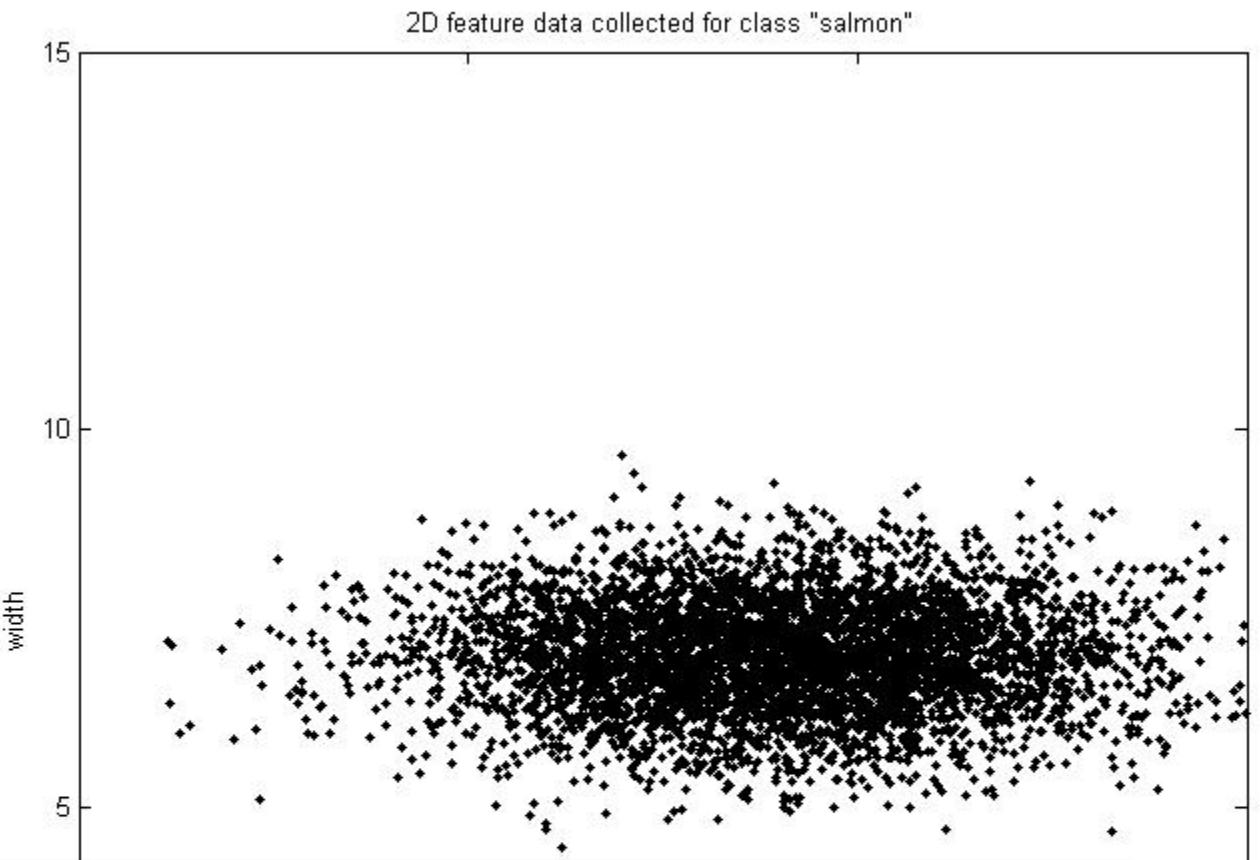
Is there any correlation between both features?

Or: Are bright fishes typically larger than dark fishes?

Practical example:

Another distribution:

→ stronger variance of
lightness



Answer: No!

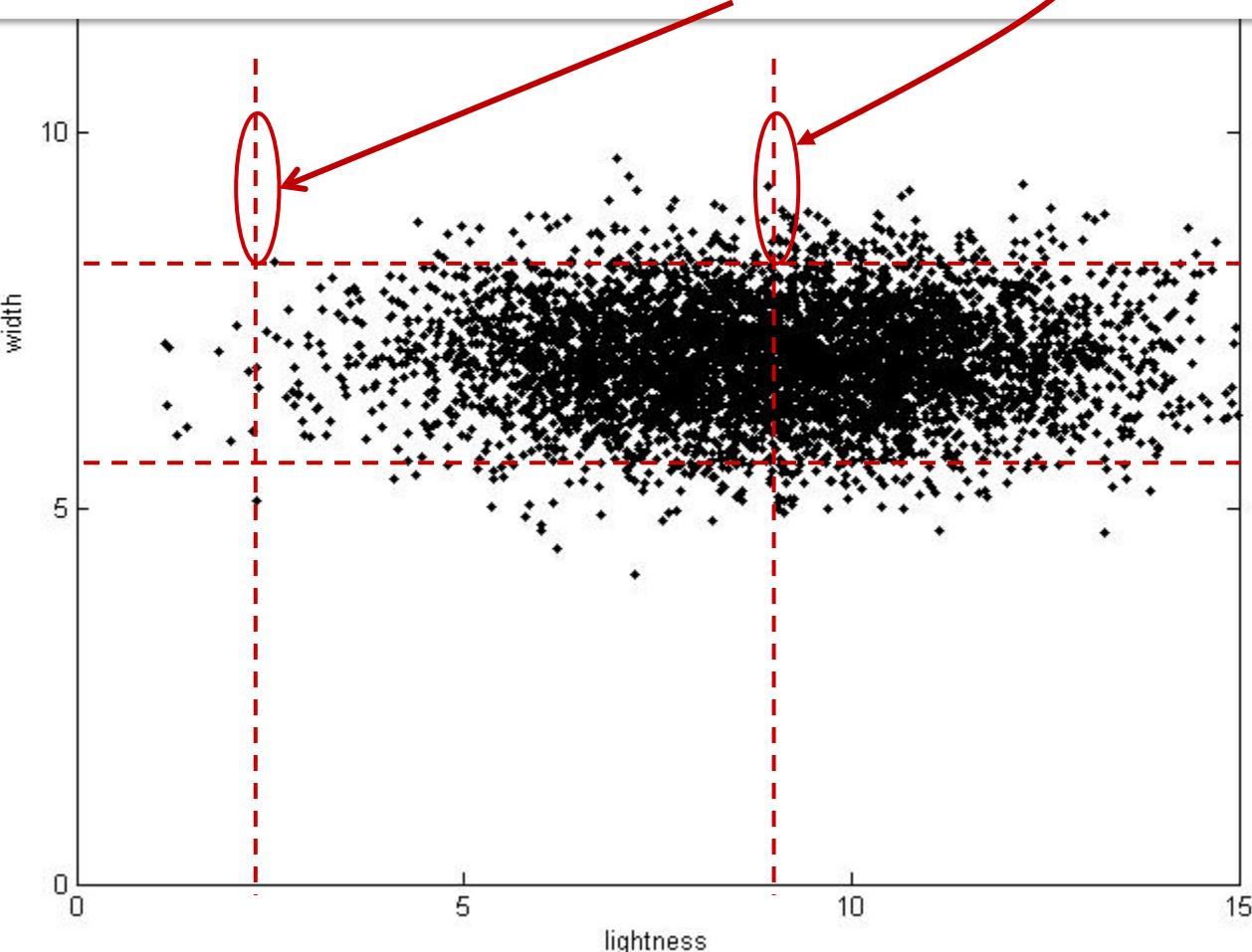
Knowledge of lightness does not improve our knowledge about width



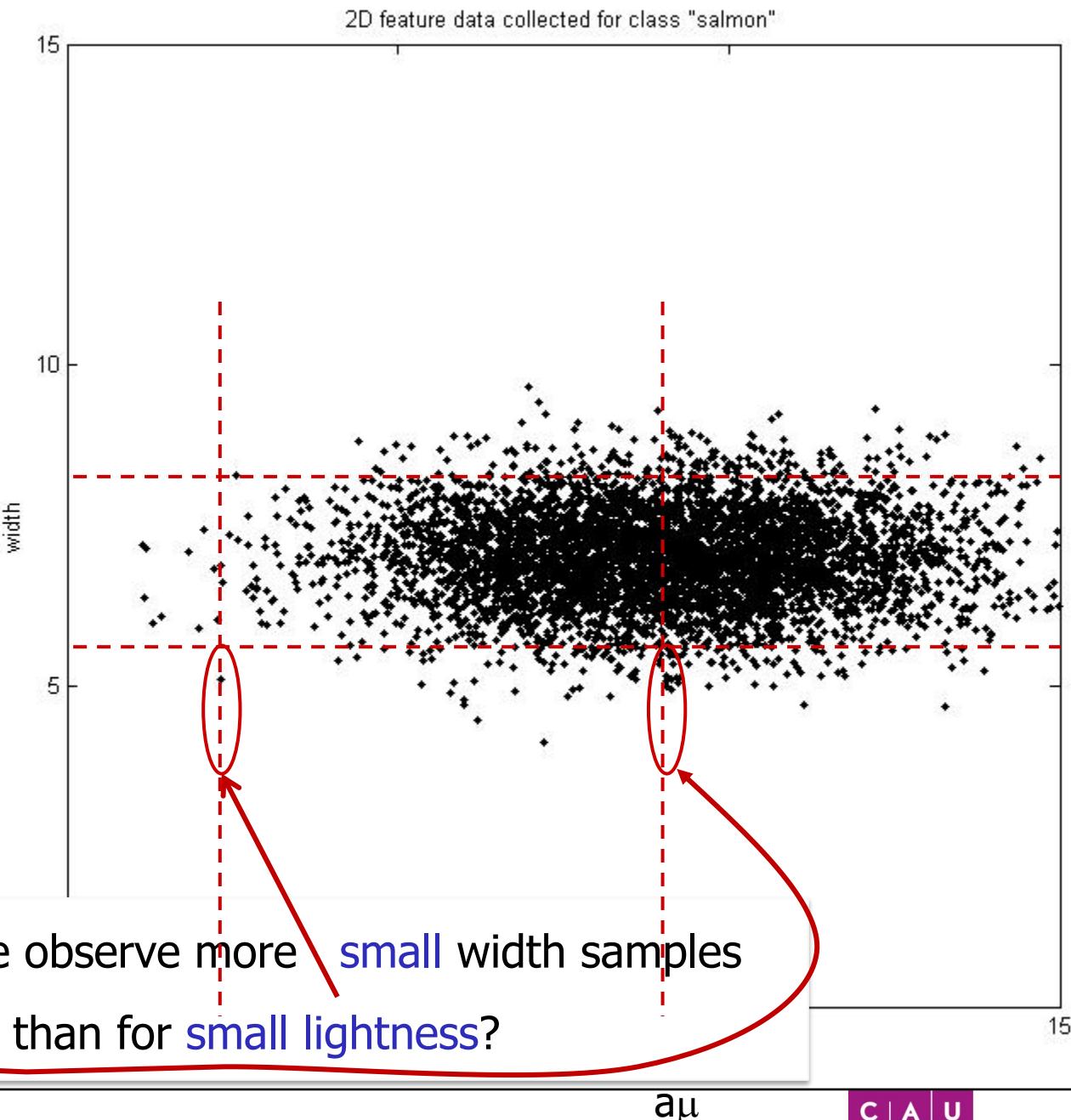
Practical example:

2D feature data collected for class "salmon"

Question: Why do we observe more **large** width samples for **medium** lightness than for **small** lightness?



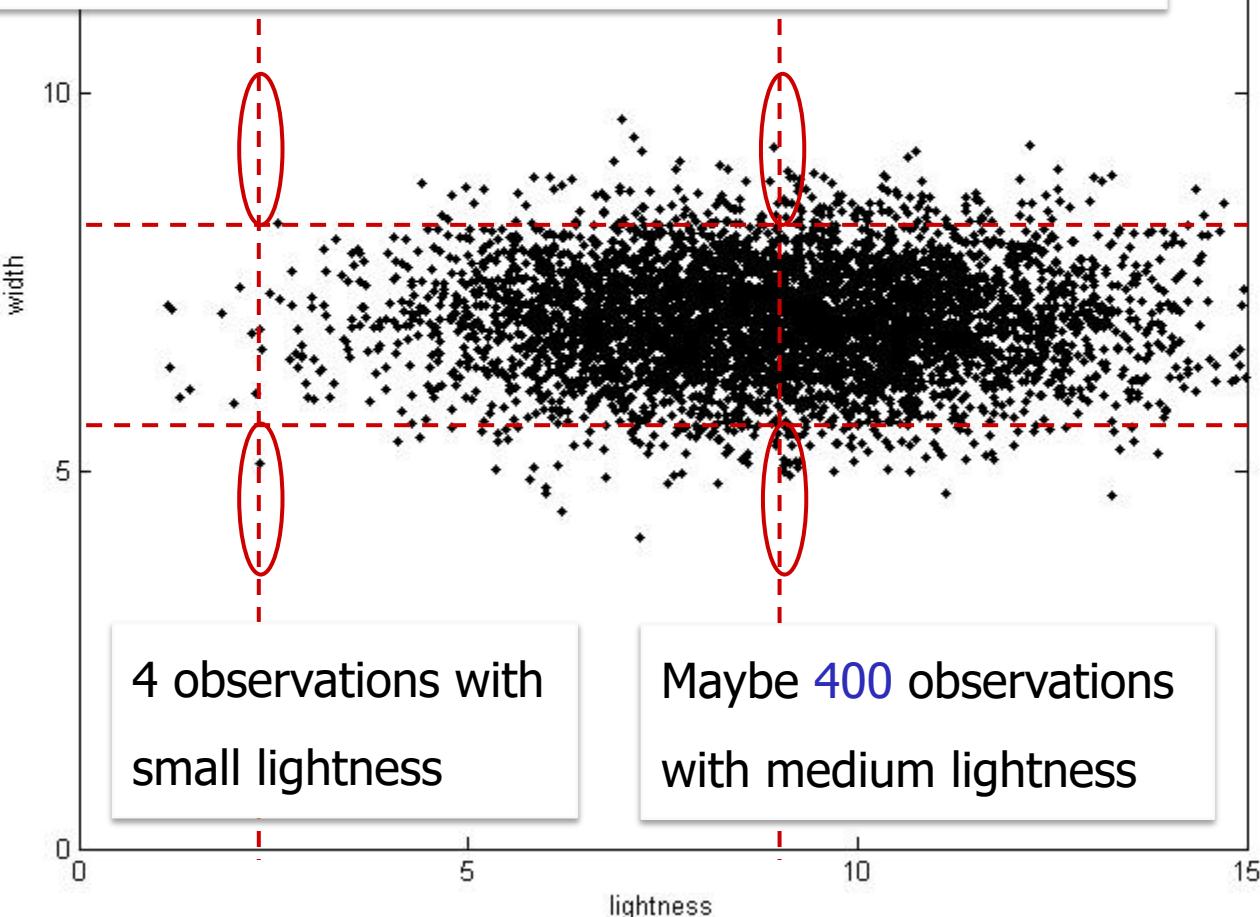
Practical example:



Practical example:

2D feature data collected for class "salmon"

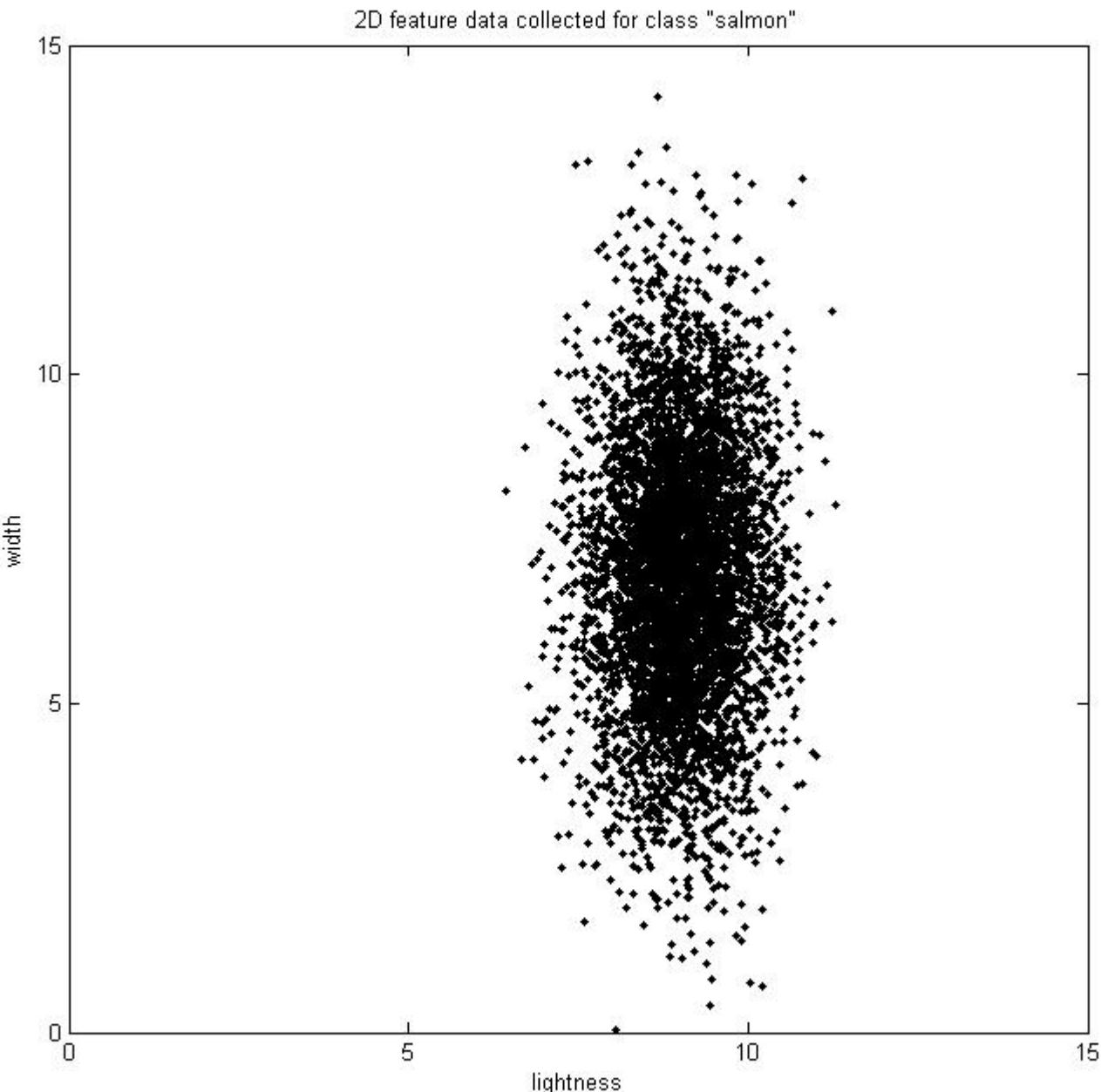
The number of observed **medium** lightness salmons is much higher
→ higher probability of unlikely events (large/small width)



Practical example:

Another distribution:

→ strong variance of
width

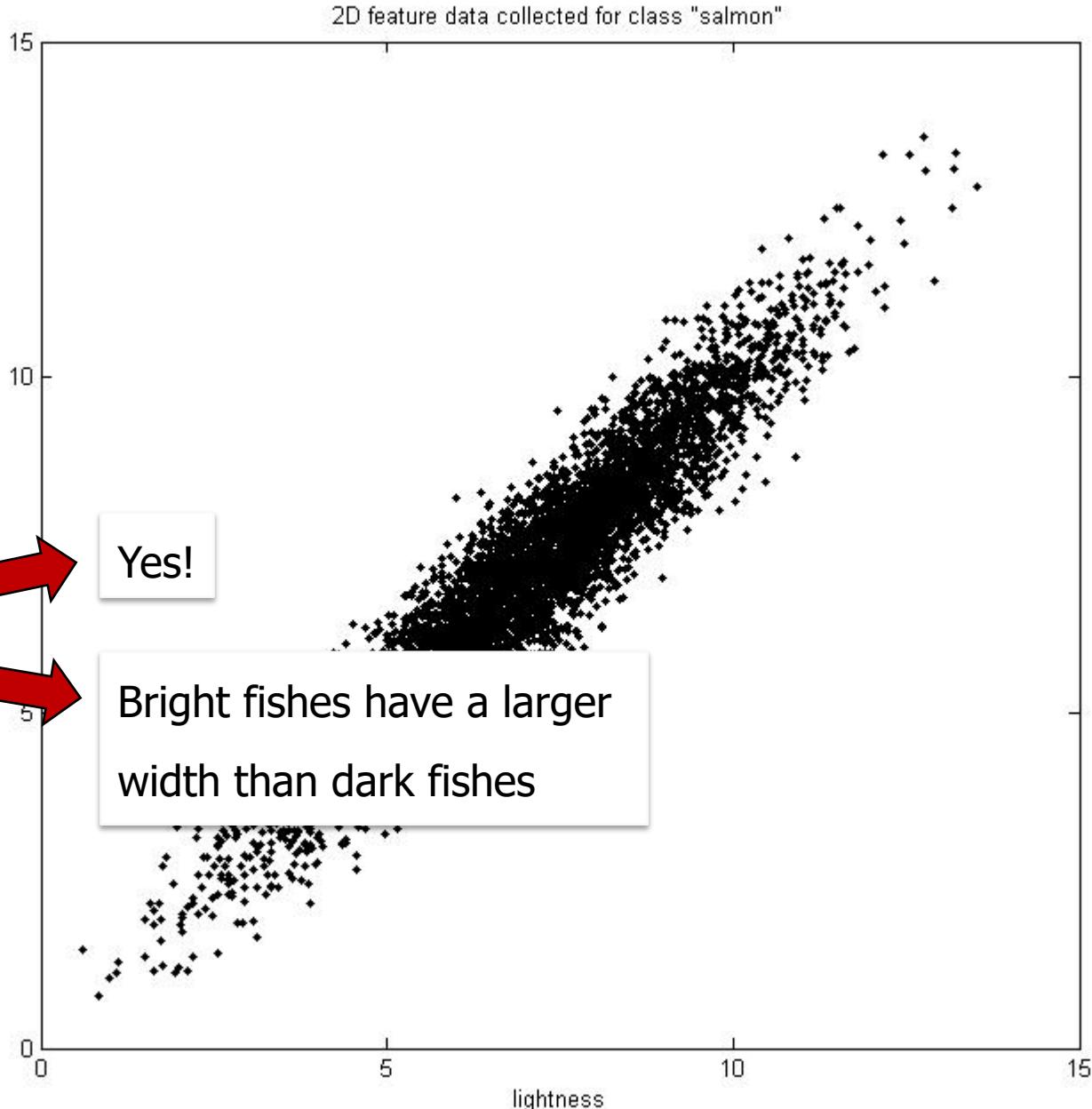


Practical example:

What about this
distribution?

Comments?

- Statistically dependend?  **Yes!**
- Practical interpretation?  Bright fishes have a larger width than dark fishes



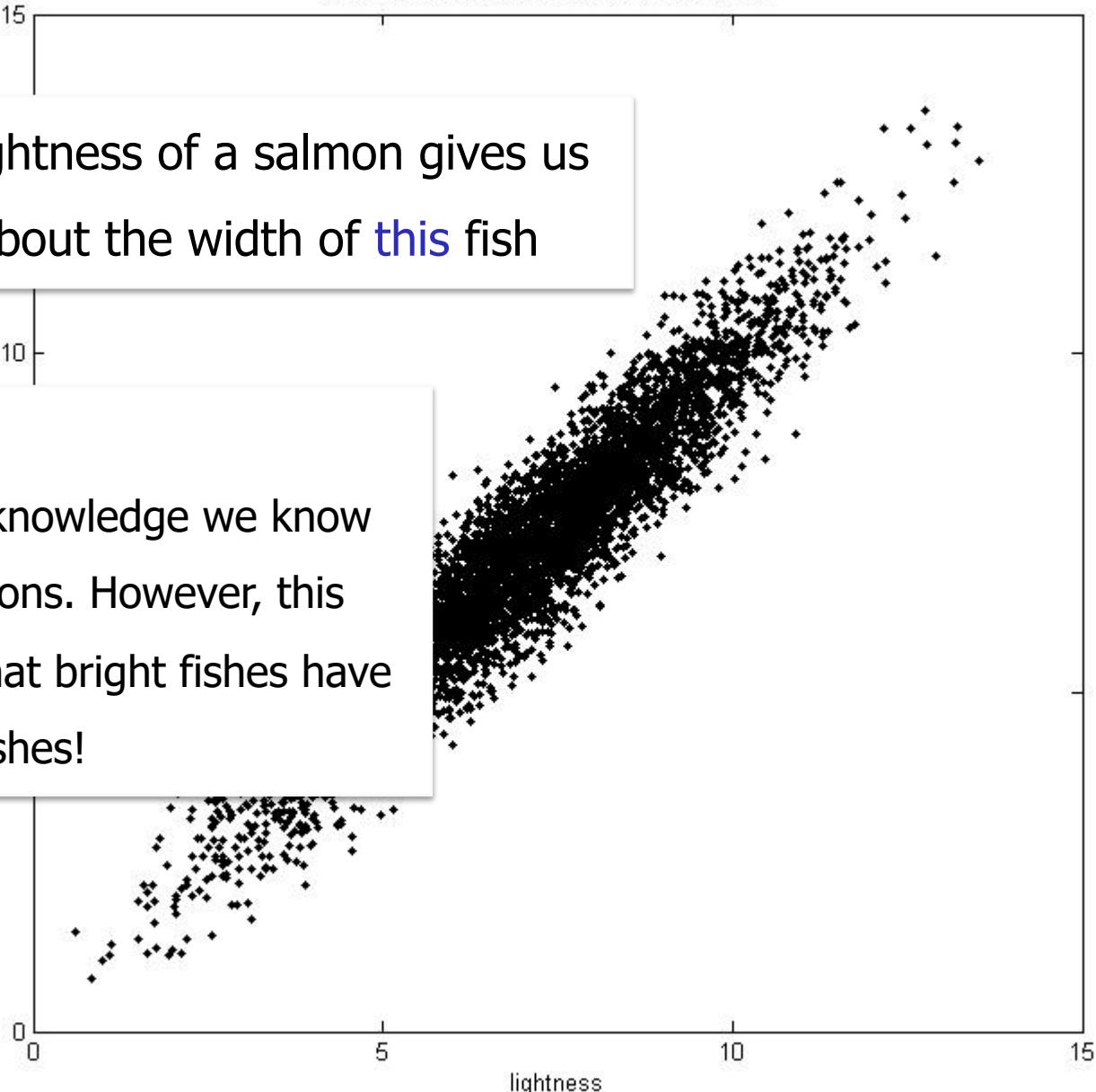
Practical example:

2D feature data collected for class "salmon"

Knowledge about the lightness of a salmon gives us
additional information about the width of this fish

Why additional?

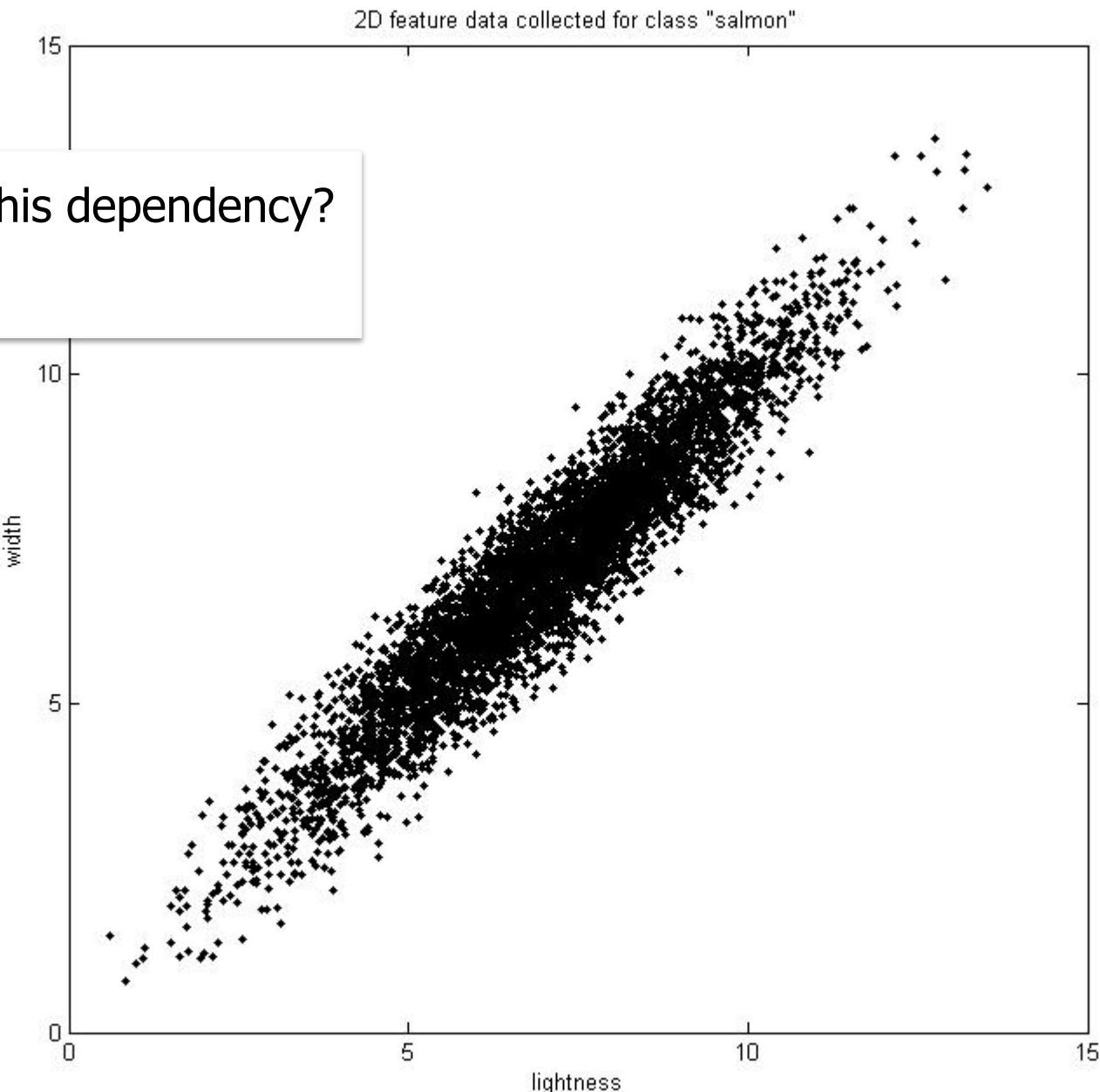
Because also without this knowledge we know
the mean width of all salmons. However, this
does not exploit the fact that bright fishes have
a larger width than dark fishes!



Practical example:

How can we **measure** this dependency?

→ Covariance



1.7 Covariance

The covariance of x and y is a 'cross-moment' which is defined as:

$$\begin{aligned}\sigma_{xy} &= E[(x - \mu_x)(y - \mu_y)] \\ &= \sum_{x \in X} \sum_{y \in Y} (x - \mu_x)(y - \mu_y) P(x, y)\end{aligned}$$

Can be used for measuring the **level of statistical dependency** of x and y .

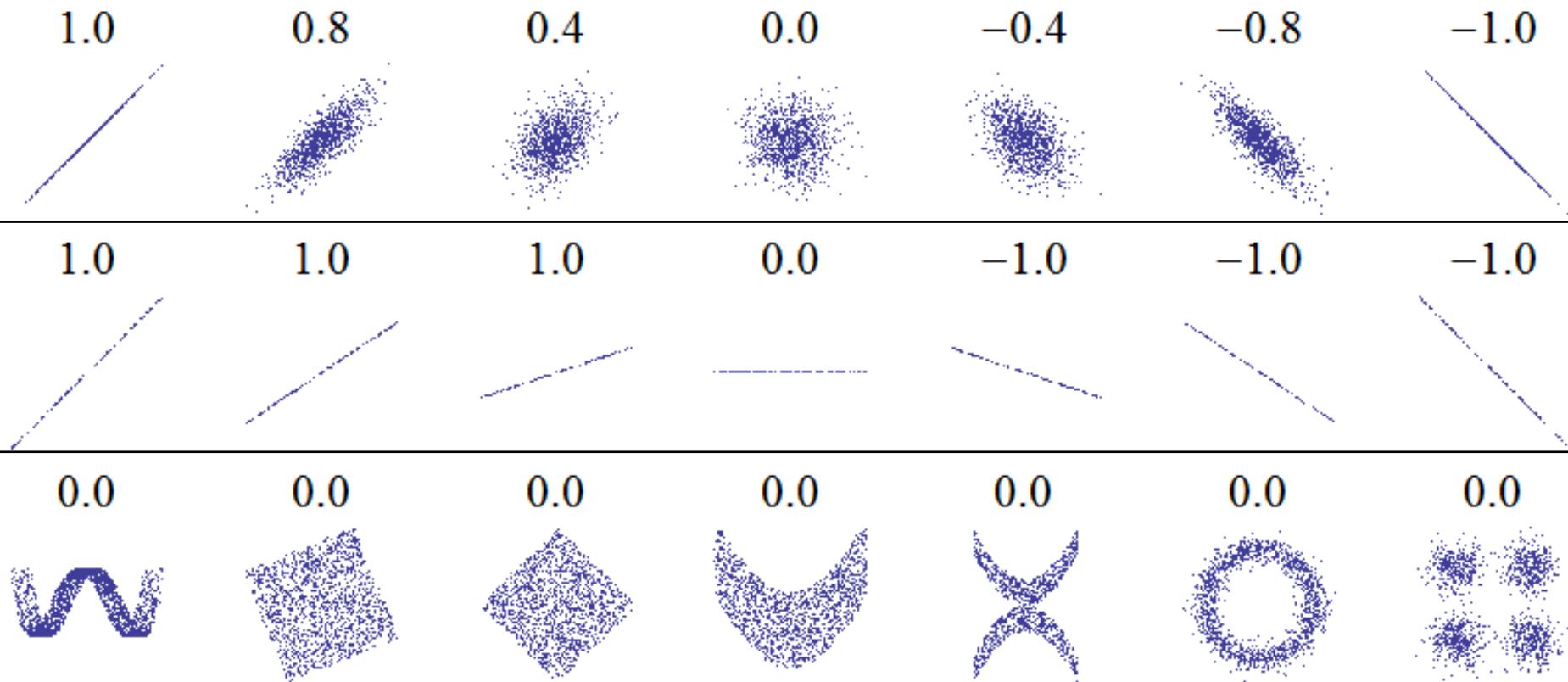
Normalized covariance: **Correlation coefficient**

$$\rho = \frac{\sigma_{xy}}{\sigma_x \cdot \sigma_y}$$

- Values between -1 and 1
- $\rho = +1$ maximally positively related
- $\rho = -1$ maximally negatively related
- $\rho = 0$ uncorrelated

1.7 Covariance

Some examples of 2D feature sample distributions and correlation coefficient:



Source: Wikipedia

Last row: Examples for uncorrelated but statistically dependent distributions

1.7 Covariance

Covariance matrix

$$\Sigma = \begin{bmatrix} \sigma_{xx} & \sigma_{xy} \\ \sigma_{yx} & \sigma_{yy} \end{bmatrix} = \begin{bmatrix} \sigma_x^2 & \sigma_{xy} \\ \sigma_{yx} & \sigma_y^2 \end{bmatrix}$$

- Diagonal elements are just the variances of x and y
- Statistically independent x and y \rightarrow covariances σ_{xy} and σ_{yx} are zero
 \rightarrow matrix is diagonal
- Covariance matrix is symmetric since $\sigma_{xy} = \sigma_{yx}$

1.8 Conditional probability

Let x and y be statistically dependent variables

→ knowing the value of one gives us better estimate of the value of the other

Expressed by definition of **conditional probability**:

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

Explanation:

Probability of value x : $P(x)$

Probability of value x if we know the value of y : $P(x|y)$

1.8 Conditional probability

Let x and y be statistically dependent variables

→ knowing the value of one gives us better estimate of the value of the other

Expressed by definition of **conditional probability**:

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

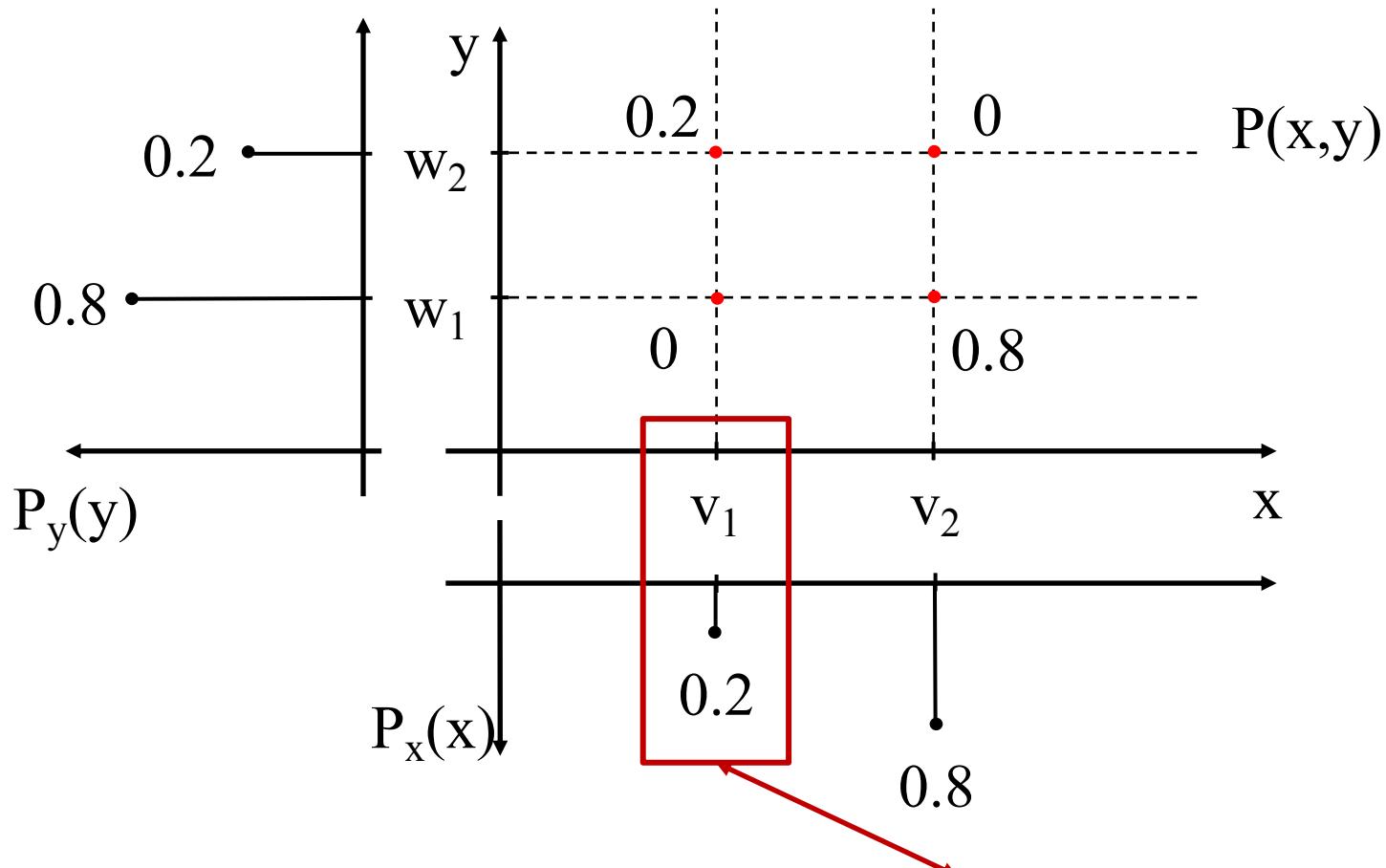
If x and y are statistically independent: $P(x|y) = P(x)$

The probability of x stays the same – independently of knowledge about y .

1.8 Conditional probability

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

Example: x and y are statistically **dependent**

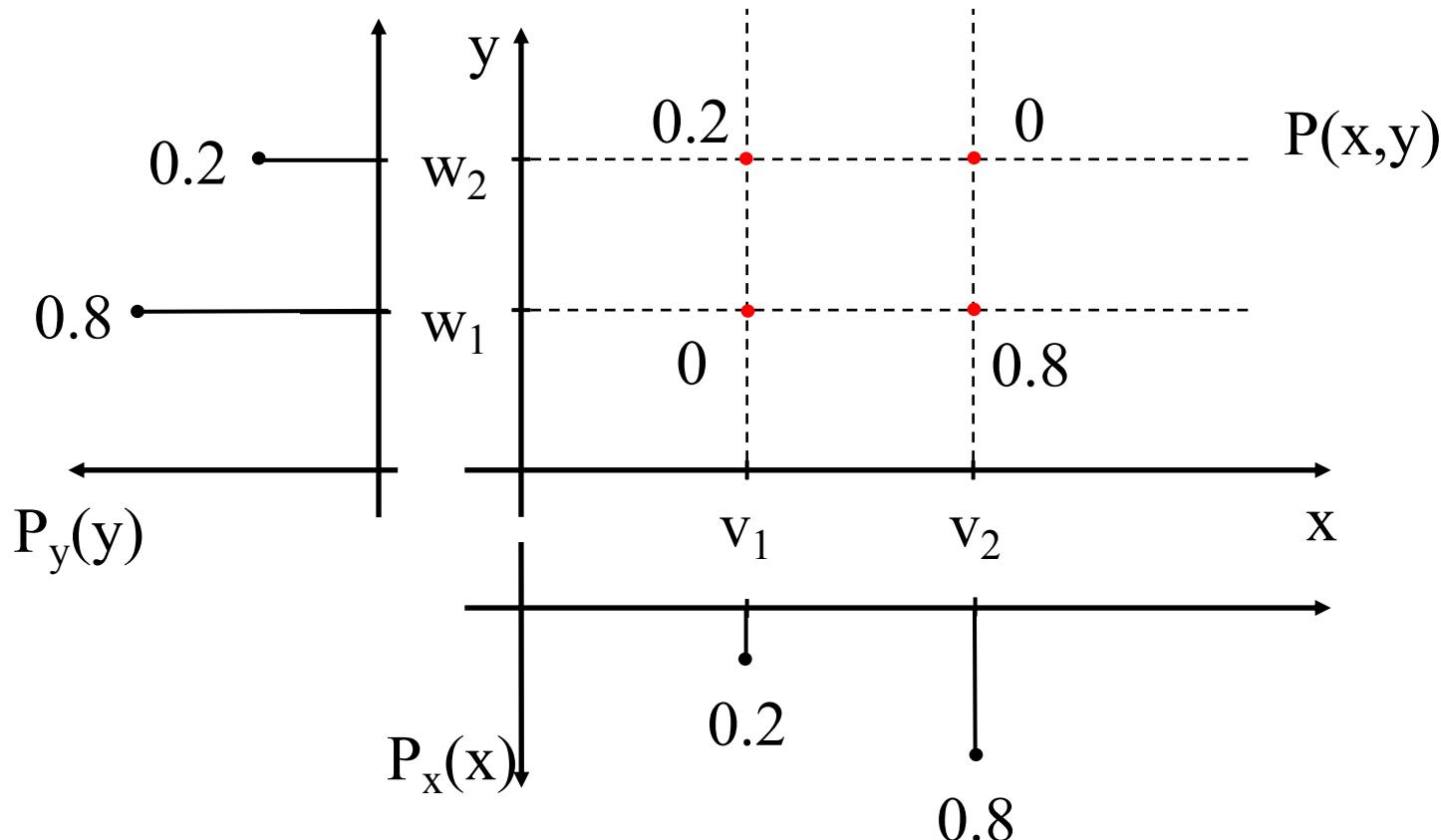


Probability of $x = v_1$ **without** knowledge about y : $P(x=v_1) = 0.2$

1.8 Conditional probability

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

Example: x and y are statistically **dependent**

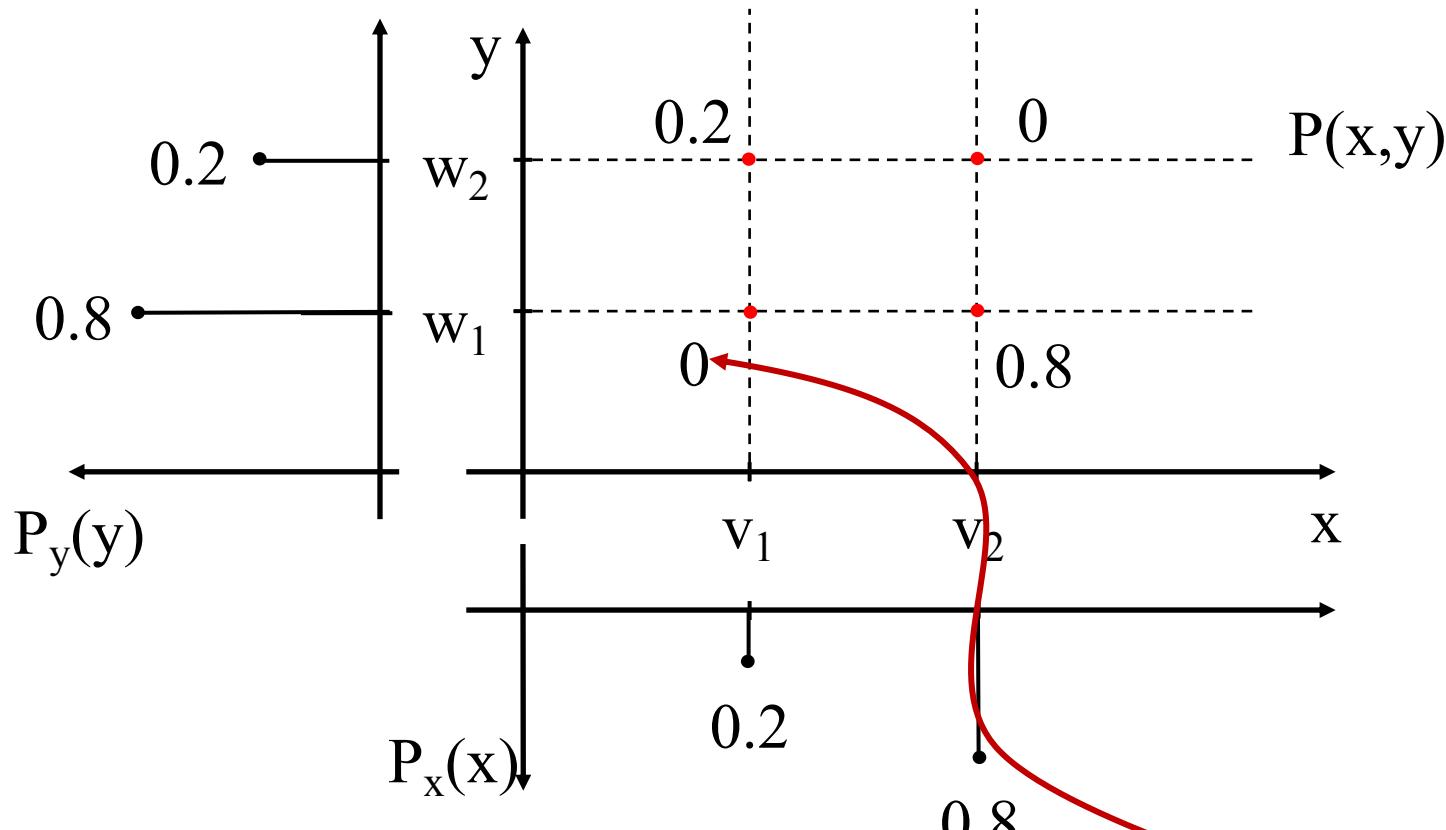


Probability of $x = v_1$ with knowledge about y : $P(x=v_1|y=w_1) = P(v_1, w_1) / P(w_1)$

1.8 Conditional probability

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

Example: x and y are statistically **dependent**

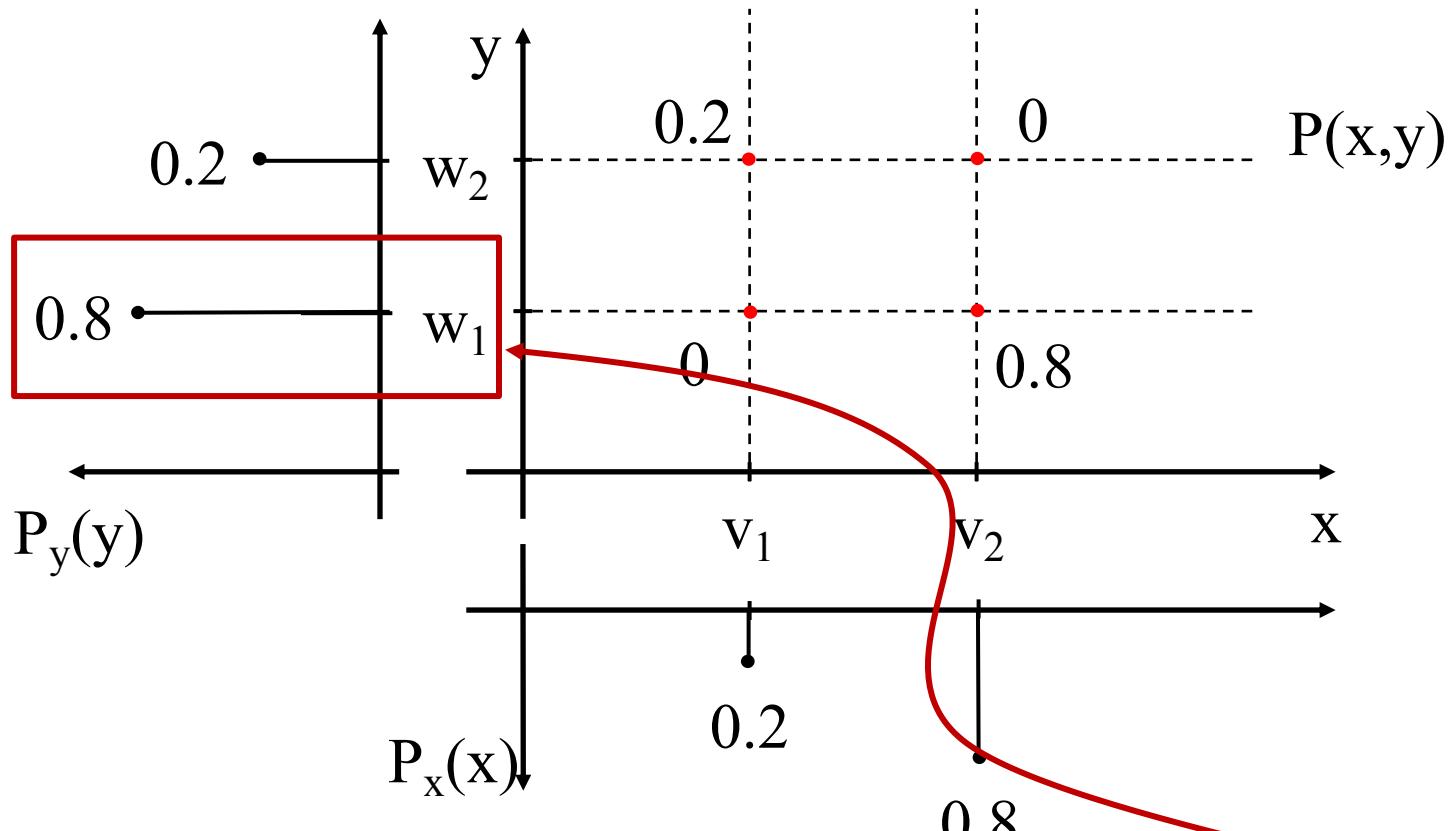


Probability of $x = v_1$ with knowledge about y : $P(x=v_1|y=w_1) = \boxed{P(v_1, w_1)} / P(w_1)$

1.8 Conditional probability

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

Example: x and y are statistically **dependent**

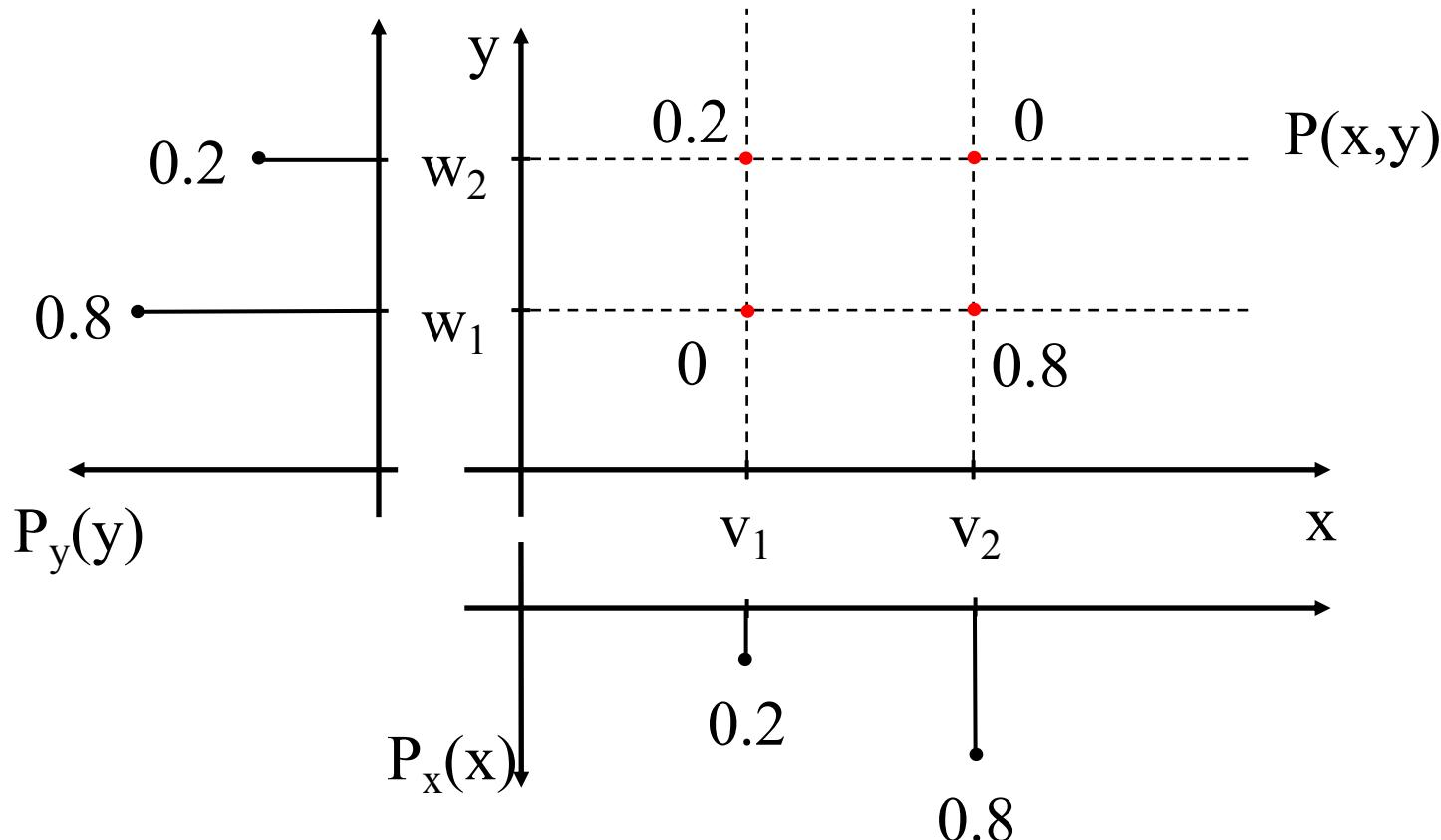


Probability of $x = v_1$ with knowledge about y : $P(x=v_1|y=w_1) = P(v_1, w_1) / P(w_1)$

1.8 Conditional probability

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

Example: x and y are statistically **dependent**

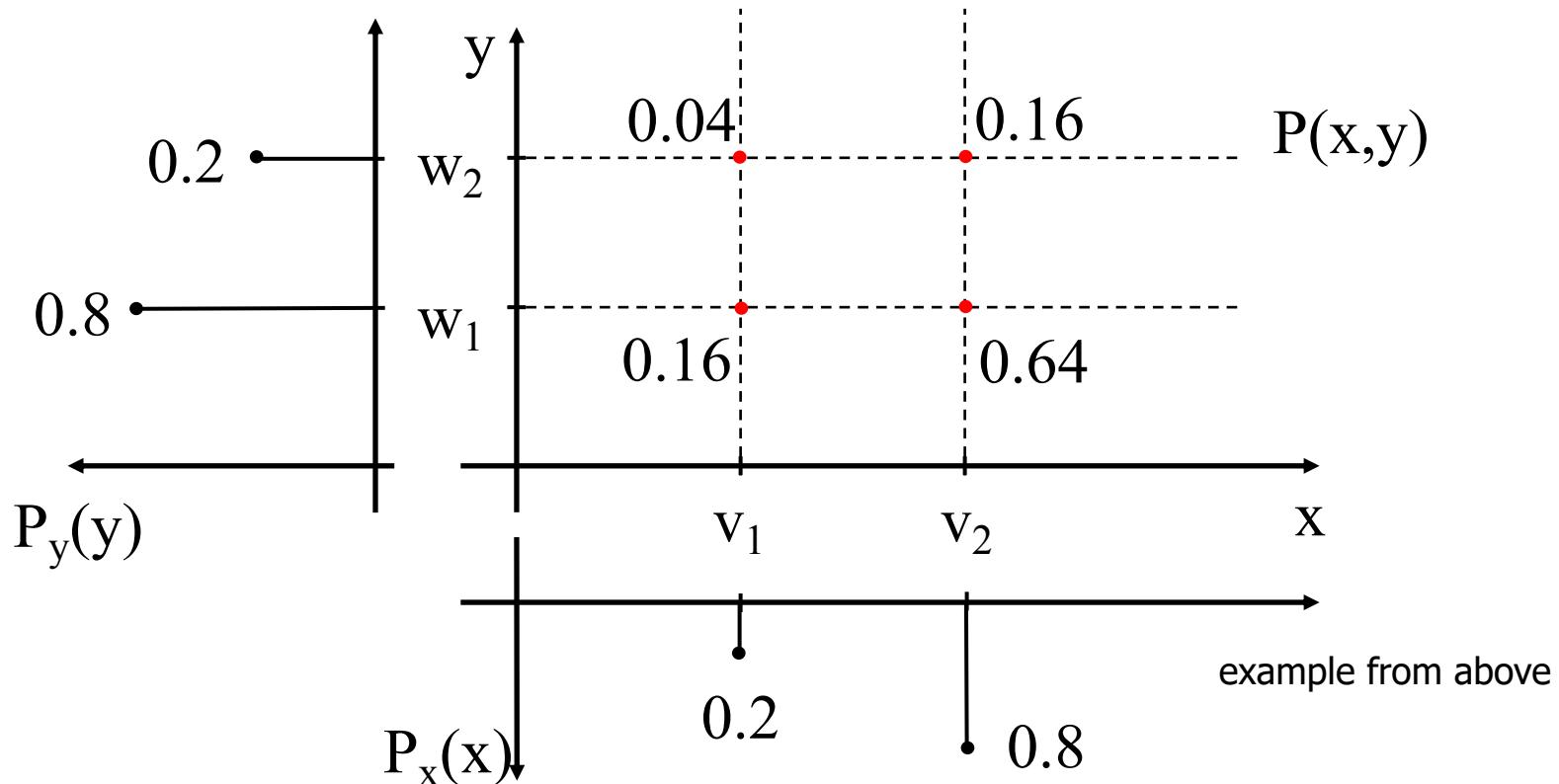


$$P(x=v_1|y=w_1) = P(v_1, w_1) / P(w_1) = 0 / 0.8 = 0 \quad \neq \quad P(x=v_1) = 0.2$$

1.8 Conditional probability

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

Example: x and y are statistically **independent**

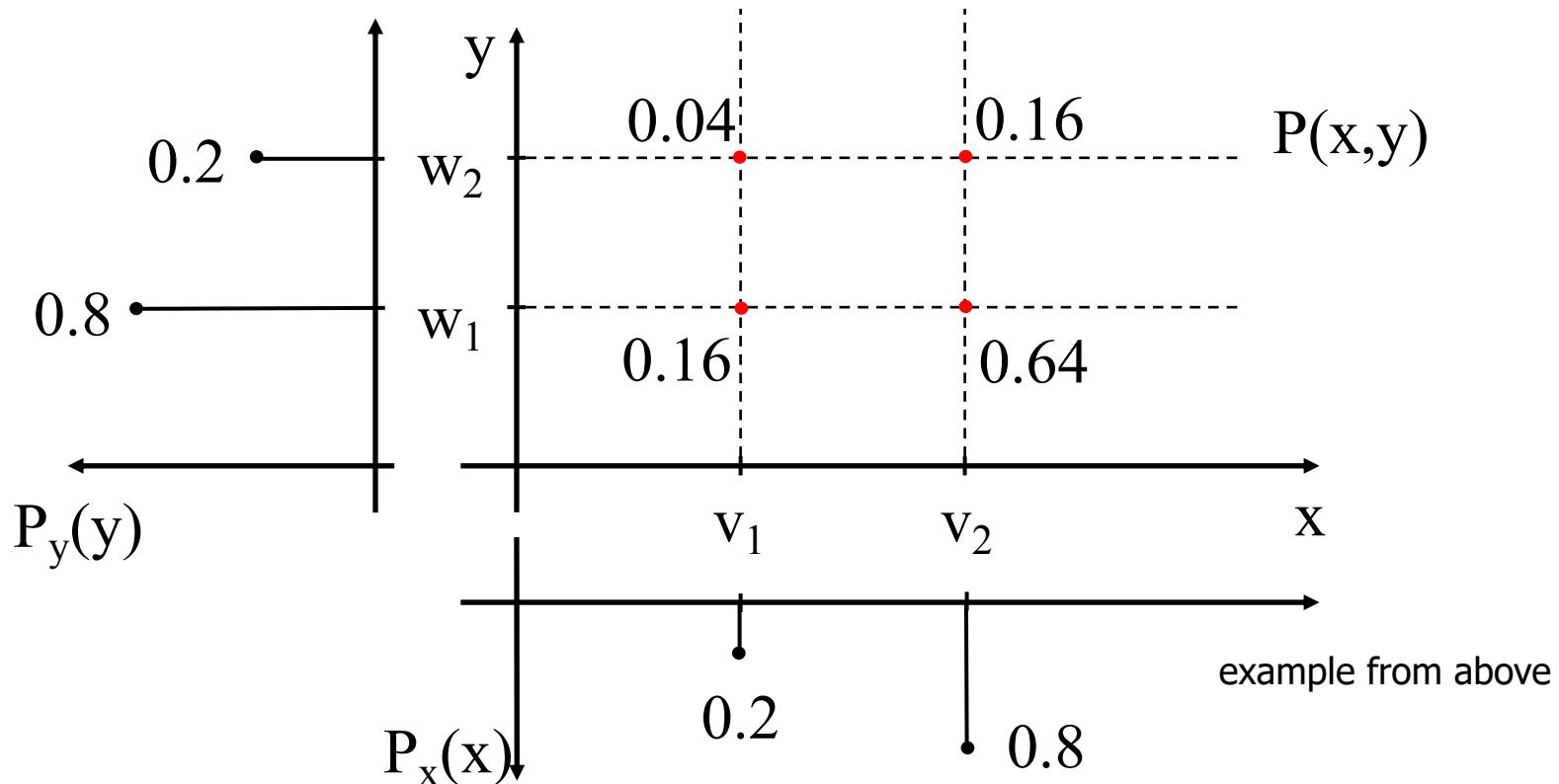


$$P(x=v_1|y=w_1) = P(v_1, w_1) / P(w_1) = 0.16 / 0.8 = 0.2 \quad \equiv \quad P(x=v_1) = 0.2$$

1.8 Conditional probability

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

Example: x and y are statistically **independent**

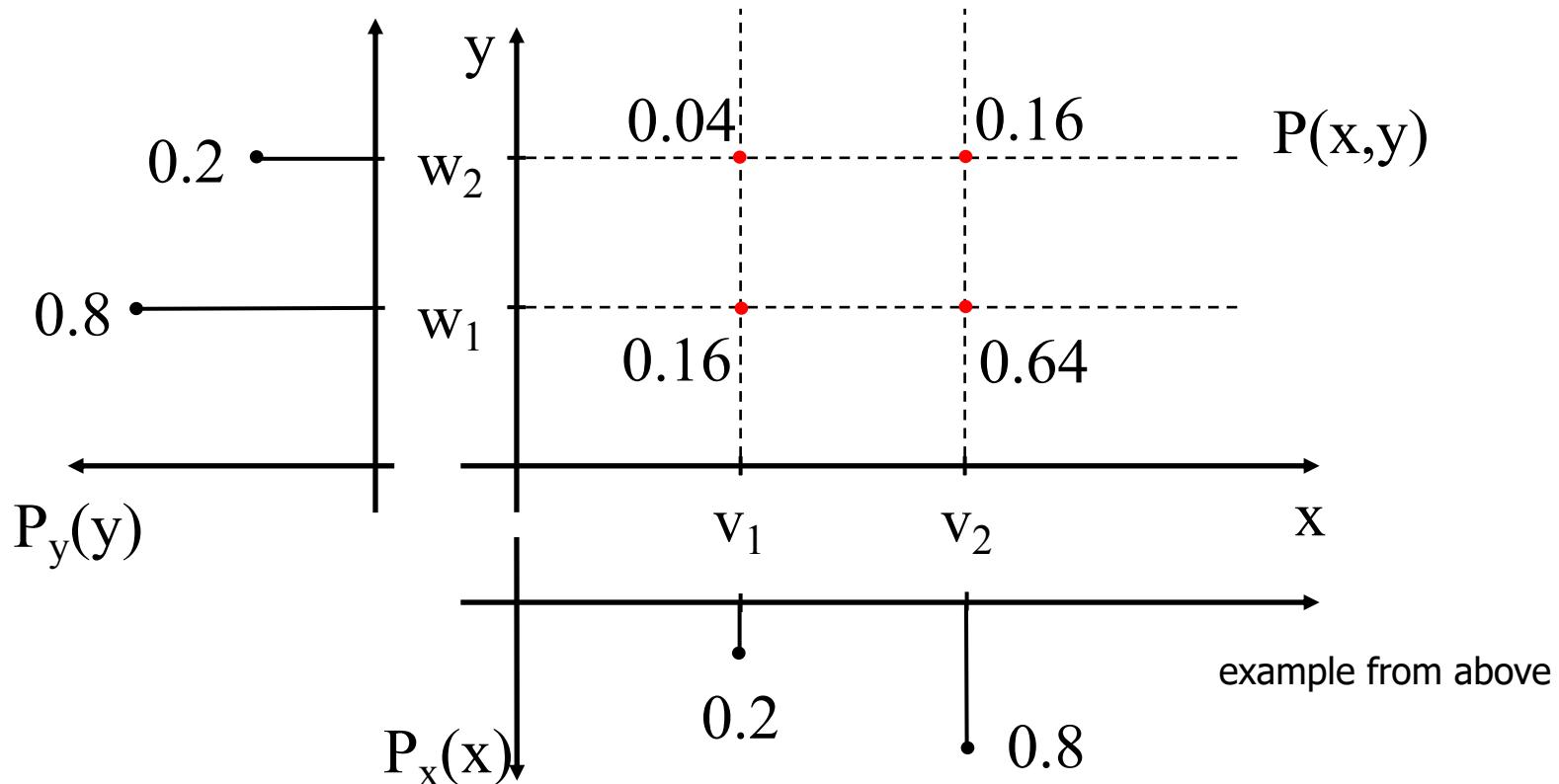


We do not need the knowledge of y to determine the probability of $x=v_1$.

1.8 Conditional probability

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

Example: x and y are statistically **independent**



$$P(x=v_2|y=w_2) = P(v_2, w_2) / P(w_2) = 0.16 / 0.2 = 0.8 \quad \equiv \quad P(x=v_2) = 0.8$$

1.9 Vector random variables

In pattern recognition applications we typically need more than one feature

→ Feature vector $\mathbf{x} = [x_1, x_2, \dots, x_d]^T$

Probability mass function satisfies

$$P(\mathbf{x}) \geq 0$$

$$\sum_{\mathbf{x}} P(\mathbf{x}) = 1$$

sum extends over all
possible values for the
vector \mathbf{x}

1.9 Vector random variables

Probability mass function can be expressed by conditional probabilities:

$$\begin{aligned} P(x_1, x_2, x_3, x_4, x_5) &= P(x_1, x_2, x_3, x_4|x_5)P(x_5) \\ &= P(x_1, x_2, x_3|x_4, x_5)P(x_4|x_5)P(x_5) \\ &= \dots \end{aligned}$$

Remark: We will use this in the context of spam filtering.

Simplification, if the random variables x_i are statistically independent:

$$P(\mathbf{x}) = P_{x_1}(x_1) P_{x_2}(x_2) \cdots P_{x_d}(x_d) = \prod_{i=1}^d P_{x_i}(x_i)$$



subscripts are used to denote that the $P(x_i)$ will in general have a different form

1.9 Vector random variables

Mean vector and Covariance matrix

Mean
vector

$$\mu = E[\mathbf{x}] = \begin{bmatrix} E[x_1] \\ E[x_2] \\ \vdots \\ E[x_d] \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_d \end{bmatrix} = \sum_{\mathbf{x}} \mathbf{x} P(\mathbf{x})$$

1.9 Vector random variables

Covariance
matrix

covariance of

variables x_1, x_2 variance of
variable x_2

$$\Sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{11} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d1} & \dots & \sigma_{dd} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \boxed{\sigma_{12}} & \dots & \sigma_{1d} \\ \sigma_{21} & \boxed{\sigma_2^2} & \dots & \sigma_{11} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d1} & \dots & \sigma_d^2 \end{bmatrix}$$

with covariances $\sigma_{ij} = \sigma_{ji} = E[(x_i - \mu_i)(x_j - \mu_j)]$ $i, j = 1, \dots, d$

1.9 Vector random variables

The covariance matrix is diagonal in case of statistically independent

variables x_1, \dots, x_d :

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{11} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d1} & \dots & \sigma_d^2 \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & & & 0 \\ & \sigma_2^2 & & \dots \\ 0 & & \ddots & \sigma_d^2 \end{bmatrix}$$

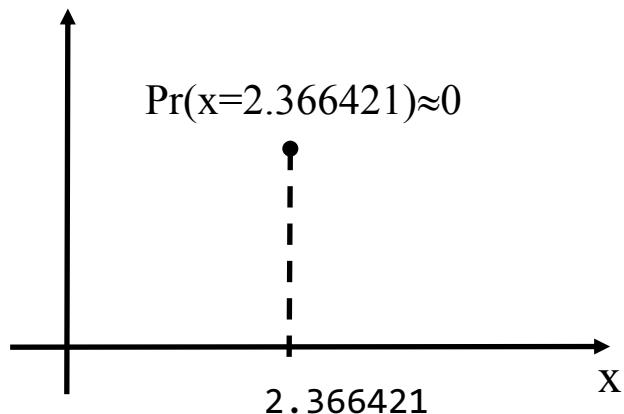
1.10 Continuous random variables

Now: random variable x can take **values in the continuum** (e.g. sensor signal)

It does not make sense to talk about probability that x has a **particular value**

Why?

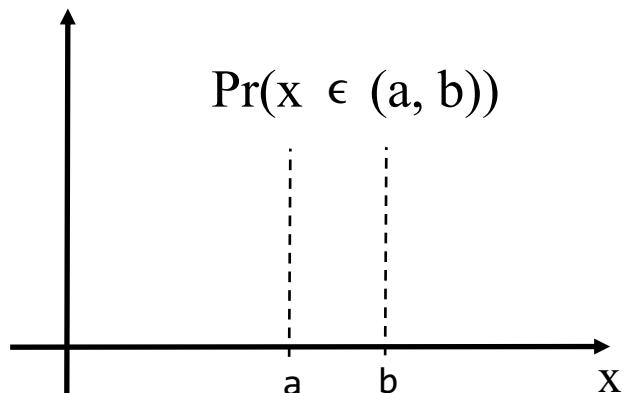
- Probability of any particular exact value (e.g. 2.366421) will almost always be zero.



Because of the infinite number of different real number values hardly anyone of these values will be observed in experiments.

1.10 Continuous random variables

We rather talk about probability that x falls into some interval (a, b)

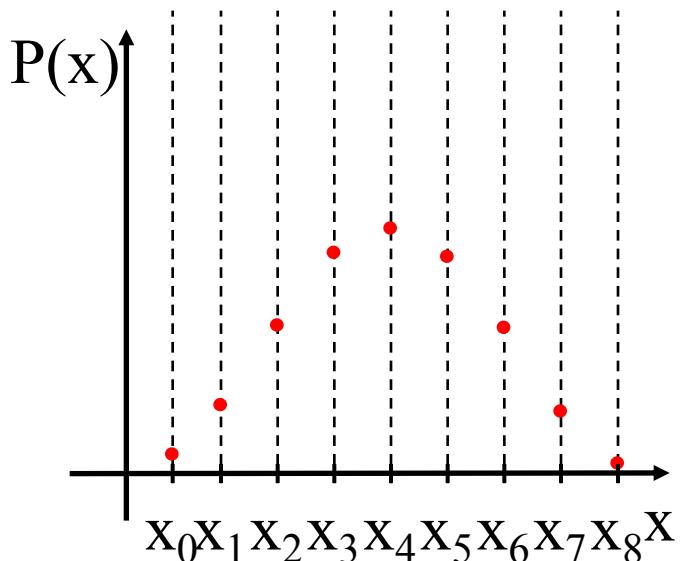


1.10 Continuous random variables

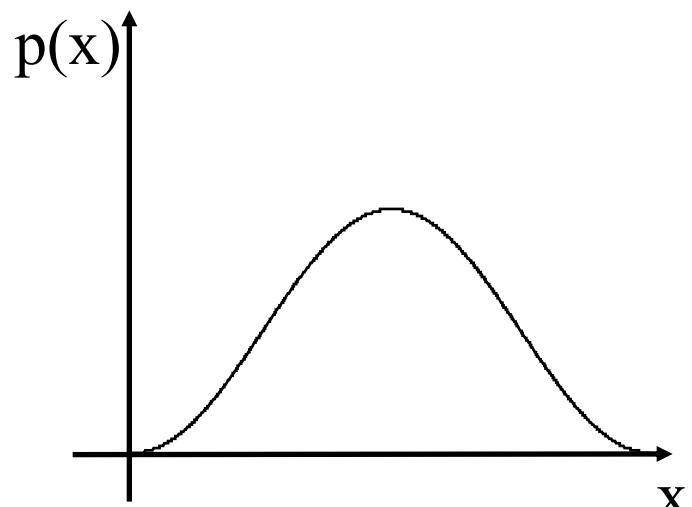
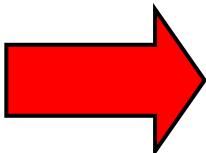
We rather talk about probability that x falls into some interval (a, b)

→ Probability Mass Function becomes a **Probability Density Function**

↓
analogy with material density



Probability Mass Function



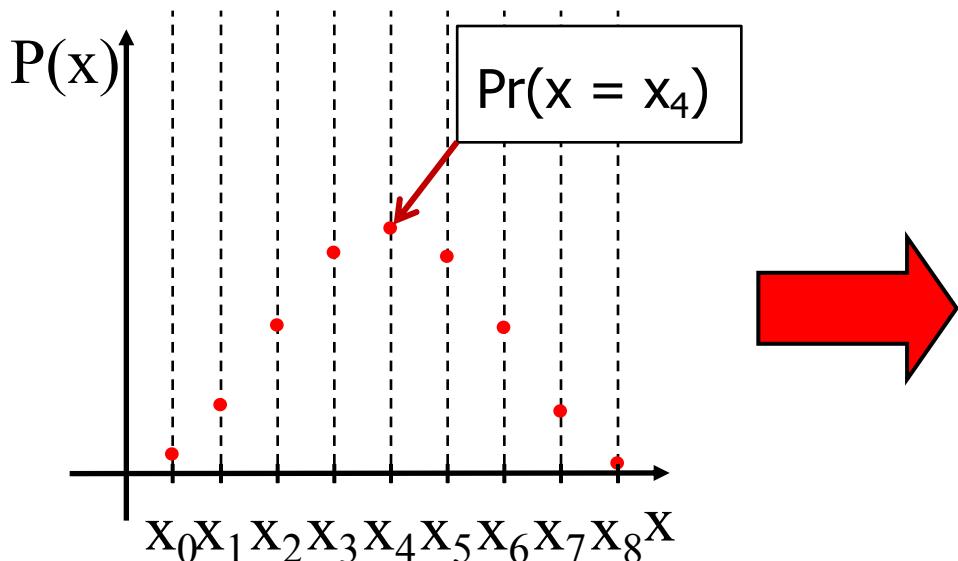
Probability Density Function

1.10 Continuous random variables

We rather talk about probability that x falls into some interval (a, b)

→ Probability Mass Function becomes a **Probability Density Function**

↓
analogy with material density



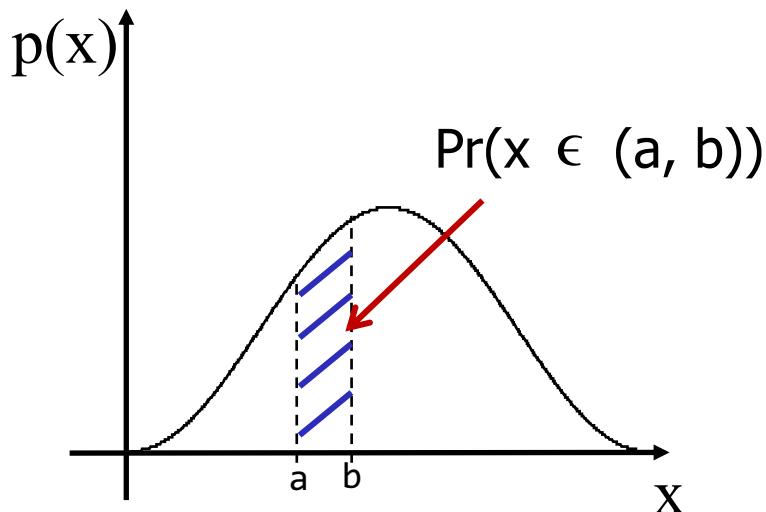
Probability Mass Function

Probability Density Function

1.10 Continuous random variables

Correspondence between probability $\Pr()$ and probability density $p()$:

$$\Pr(x \in (a, b)) = \int_a^b p(x)dx$$

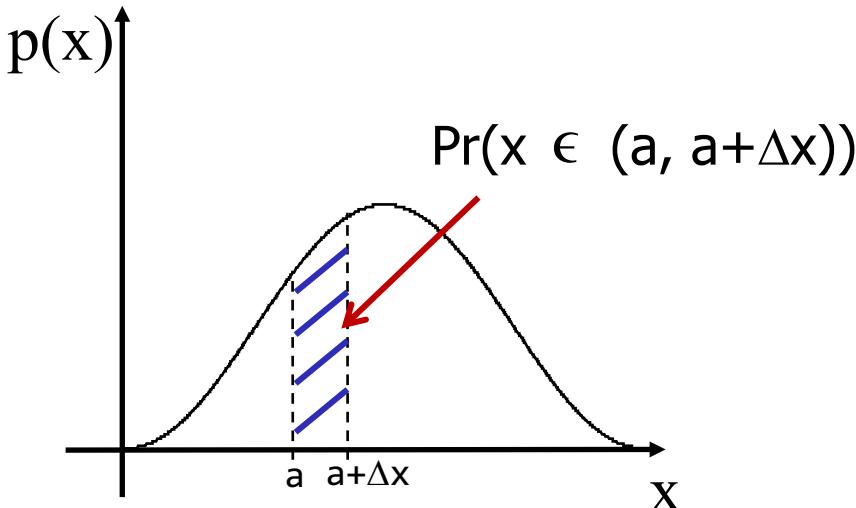


Probability Density Function

1.10 Continuous random variables

Correspondence between probability $\Pr()$ and probability density $p()$:

$$p(x) = \lim_{\Delta x \rightarrow 0} \frac{\Pr\{x \in (a, a + \Delta x)\}}{\Delta x}$$

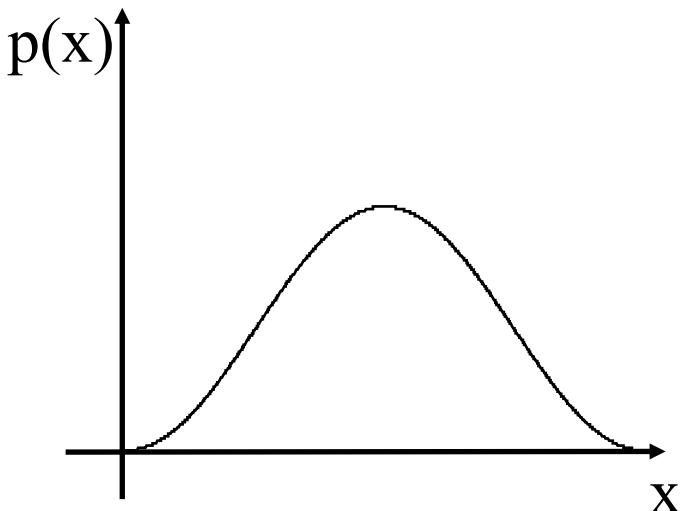


Probability Density Function

1.10 Continuous random variables

Properties of the probability density function:

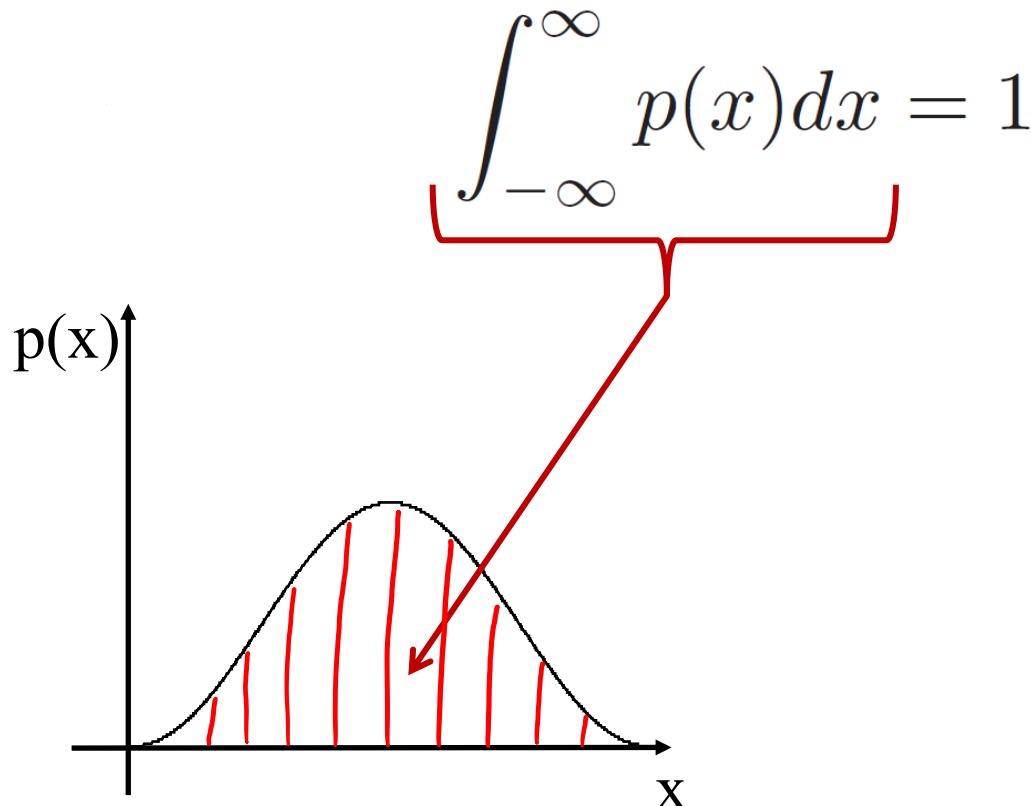
$$p(x) \geq 0$$



Probability Density Function

1.10 Continuous random variables

Properties of the probability density function:



Probability Density Function

1.10 Continuous random variables

Most of the definitions and formulas for discrete random variables carry over to continuous random variables with sums replaced by integrals:

**Expected value
of a function $f(x)$:**

$$E[f(x)] = \int_{-\infty}^{\infty} f(x) p(x) dx$$

Mean

$$\mu = E[x] = \int_{-\infty}^{\infty} x p(x) dx$$

Variance

$$\sigma^2 = E[(x - \mu)^2] = \int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx$$

1.11 Normal distributions

Important result of the probability theory:

Central Limit Theorem:

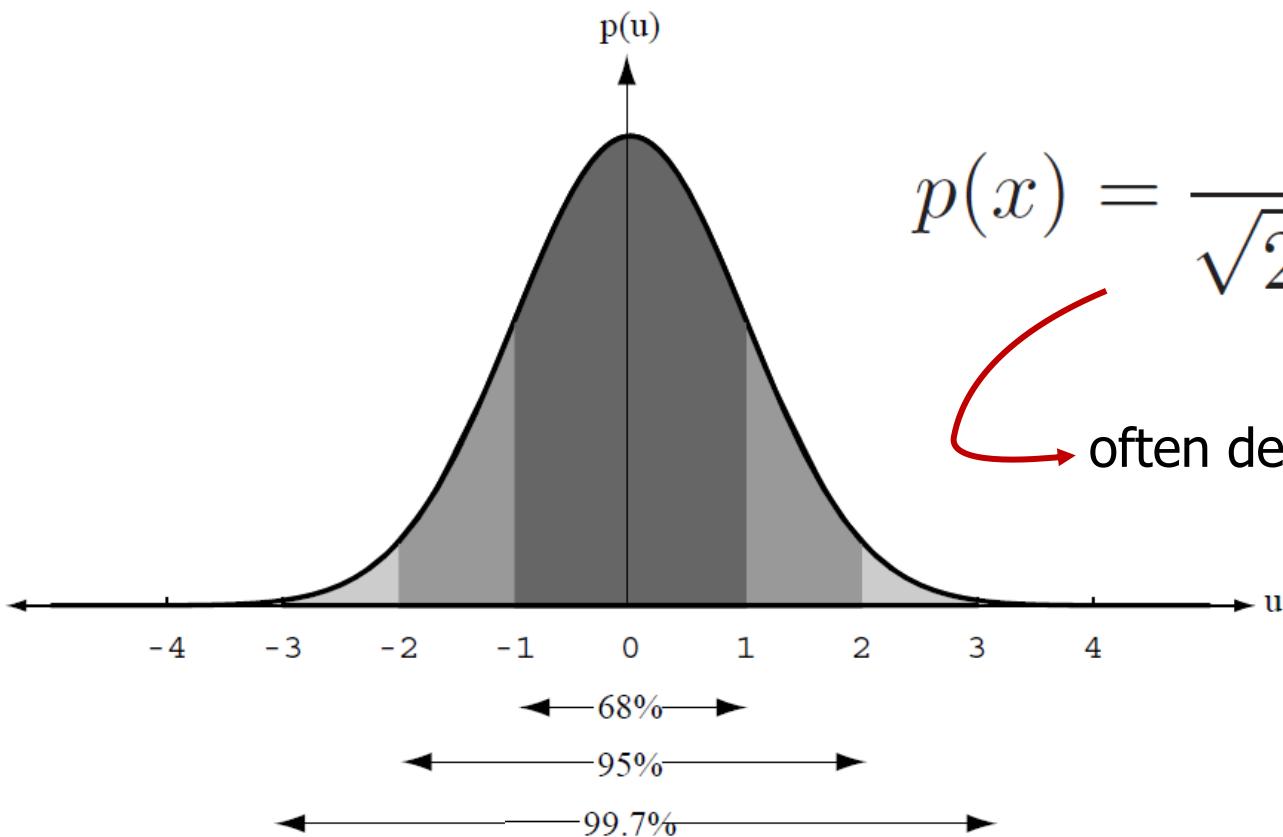
Random variables generated under a large number of statistically independent influences can be considered *normally distributed*:

Representation by a Normal or **Gaussian**
Probability Density Function

1.11 Normal distributions

Univariate Gaussian Distribution

1-dimensional random variable x

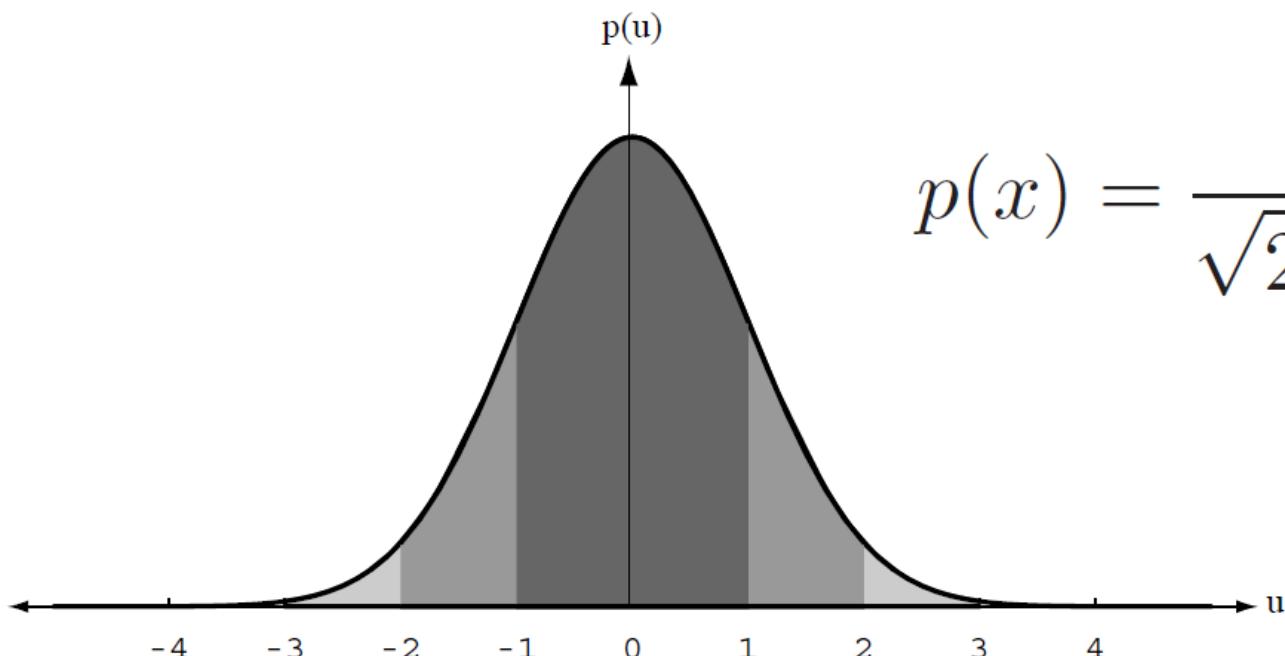


$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

often denoted as $N(\mu, \sigma^2)$

1.11 Normal distributions

Univariate Gaussian Distribution



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

$$\Pr(|x - \mu| \leq \sigma) \approx 0.68$$

$$\Pr(|x - \mu| \leq 2\sigma) \approx 0.95$$

$$\Pr(|x - \mu| \leq 3\sigma) \approx 0.997$$

1.12 Mahalanobis Distance

Measure of the distance between an **observation x** and a **data set** with mean μ and standard deviation σ . „Is x similar to the data set?“

Based on the distance $|x - \mu|$ measured **in units of the standard deviation σ** :

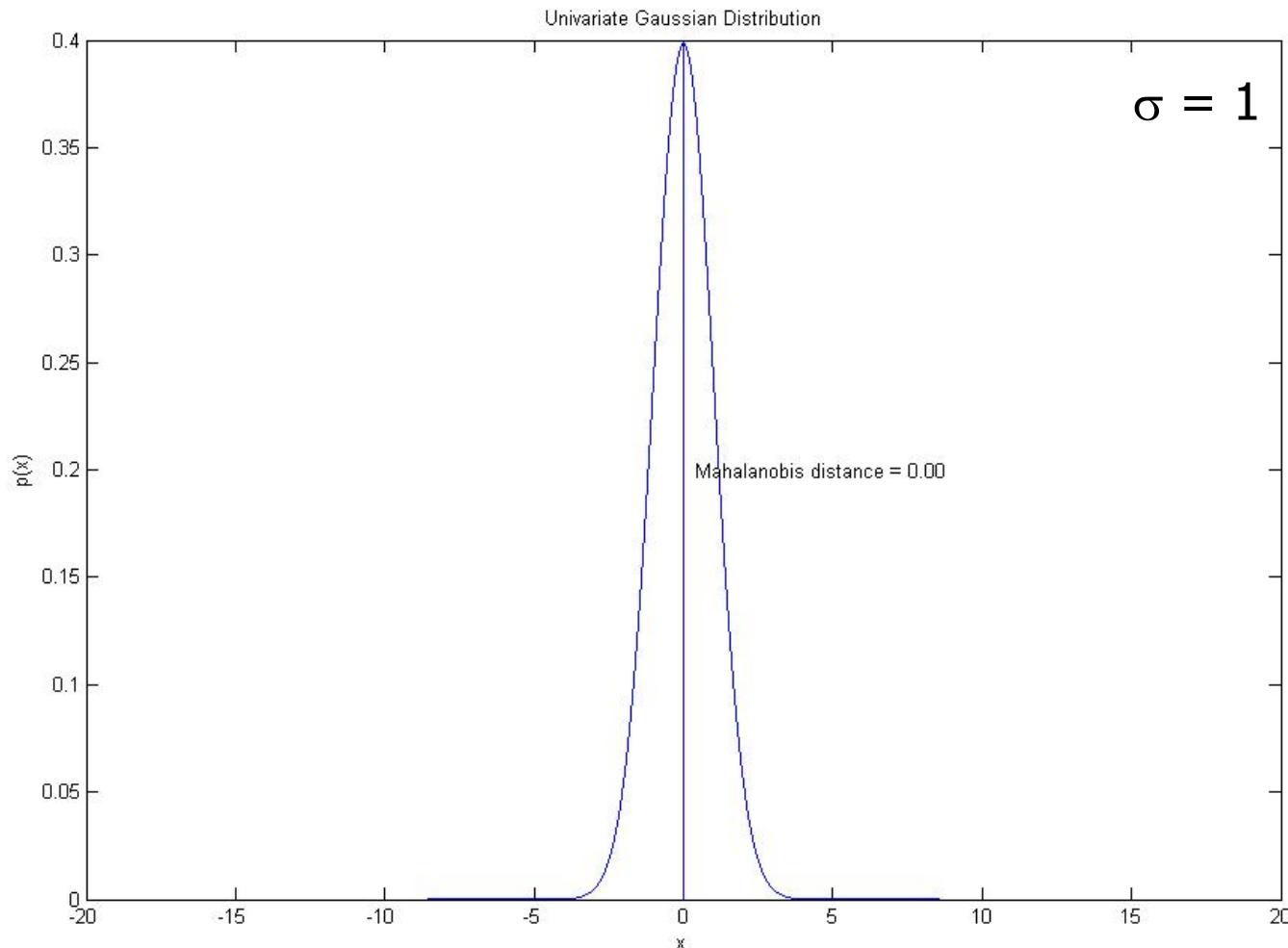
Mahalanobis distance:

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$

→ smaller distance in case of larger variance

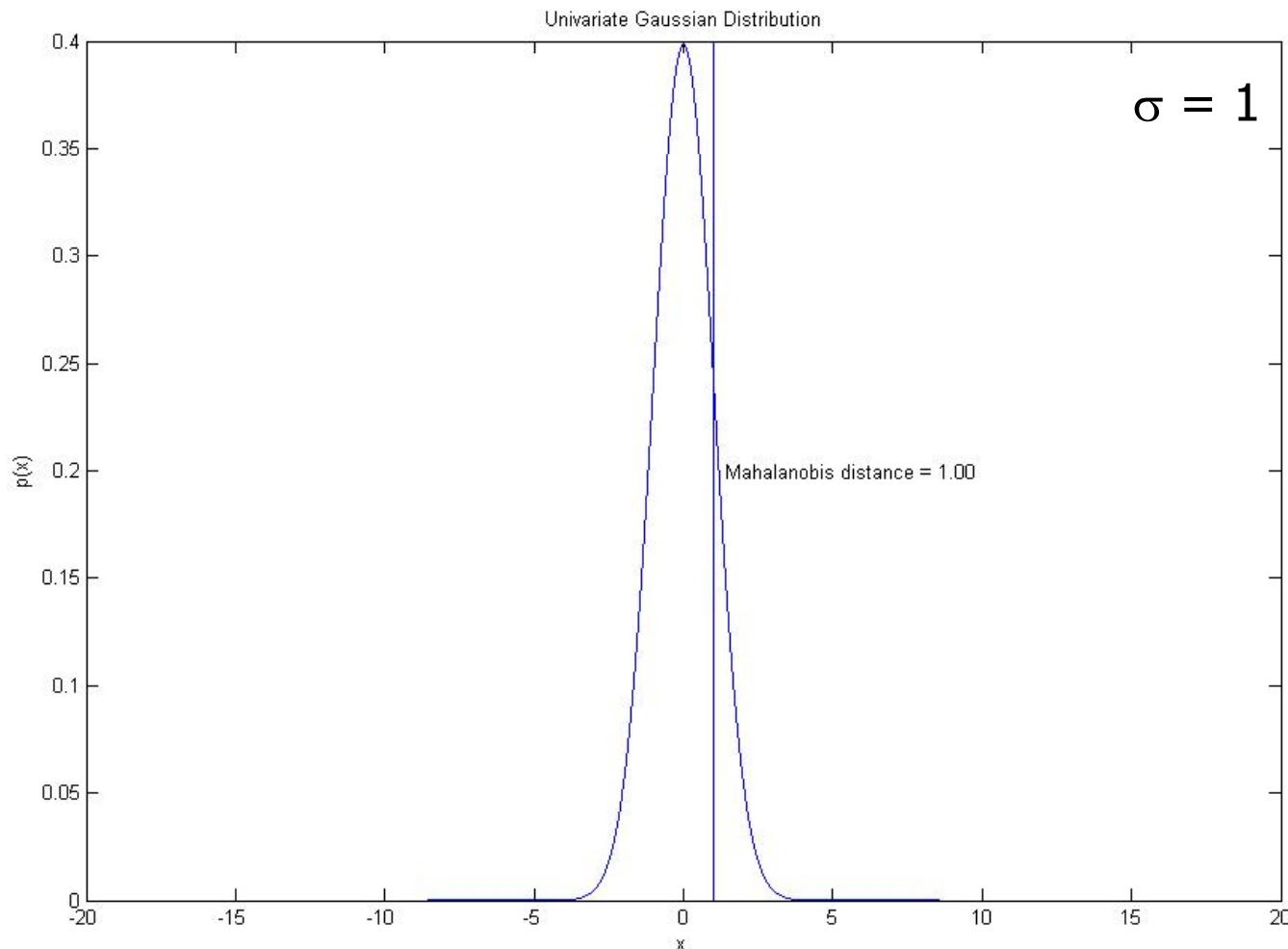
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



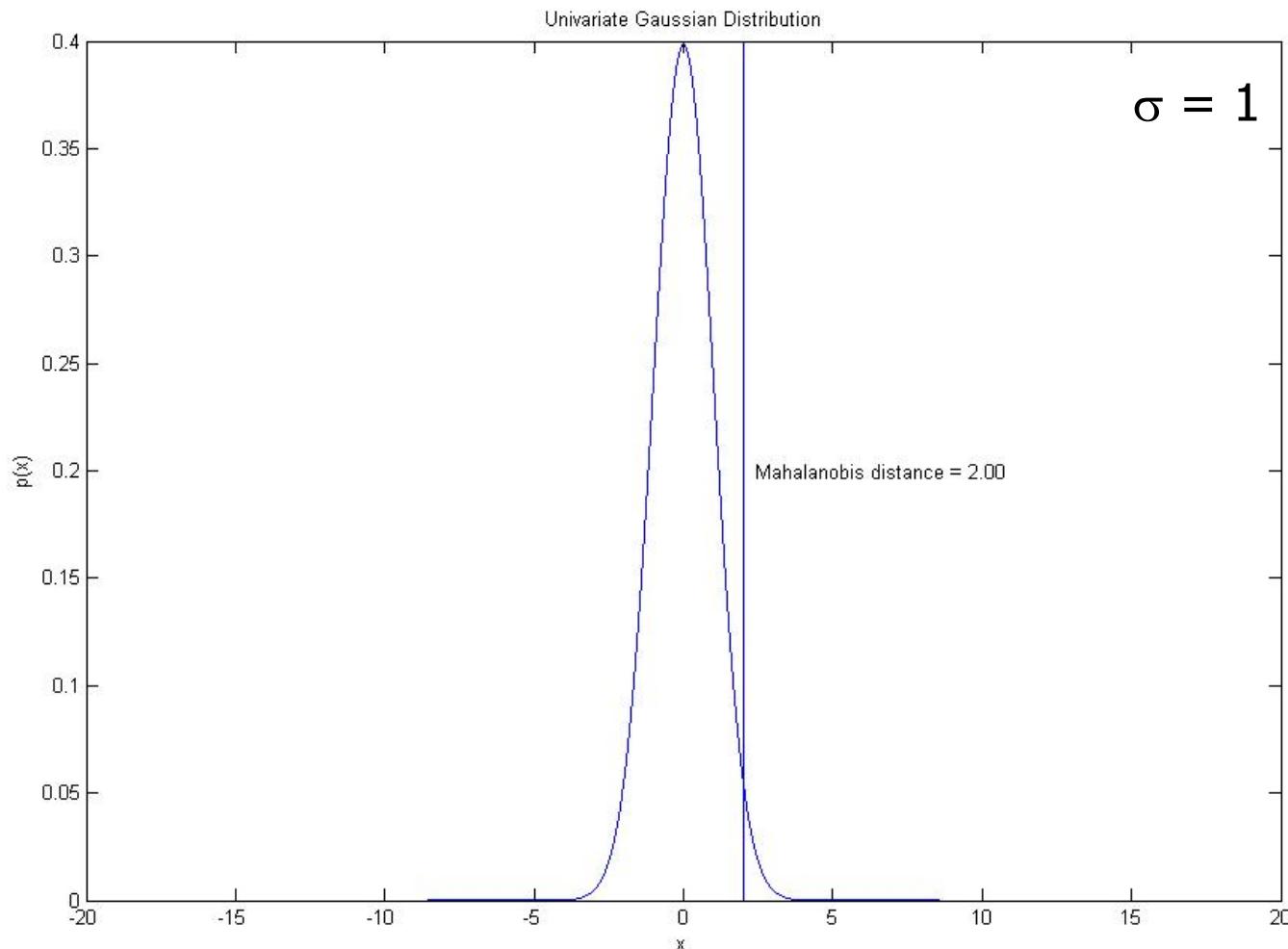
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



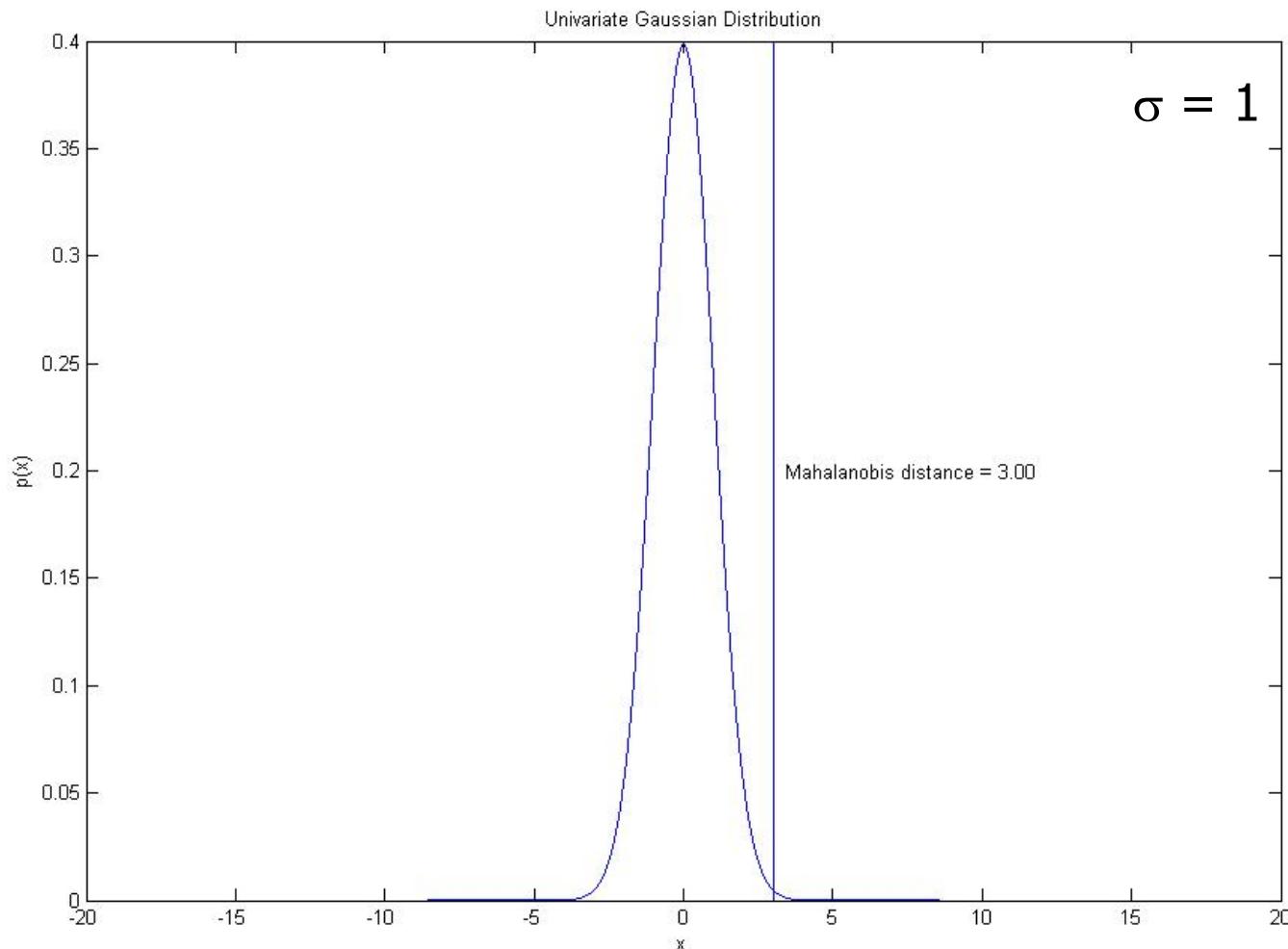
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



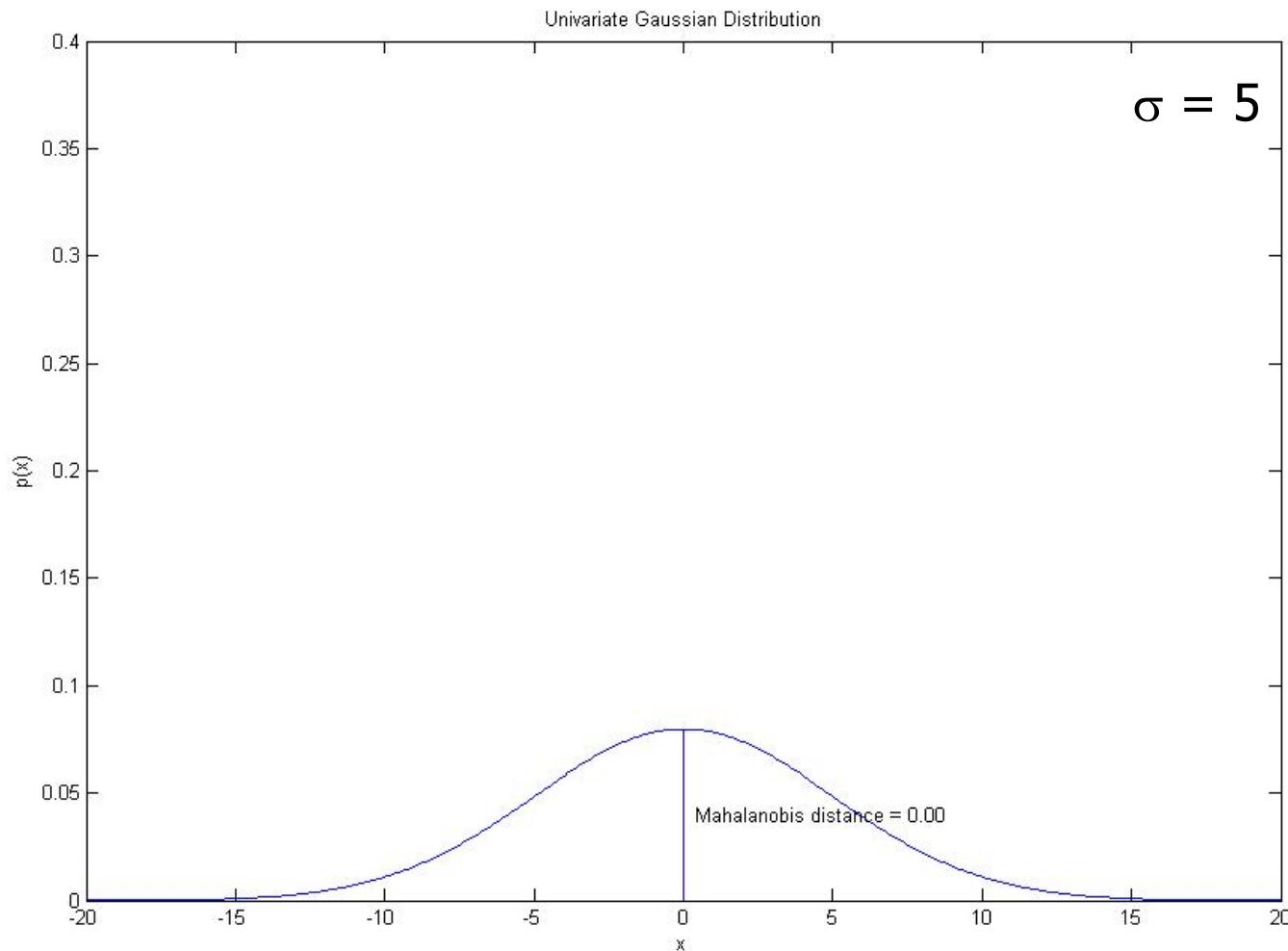
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



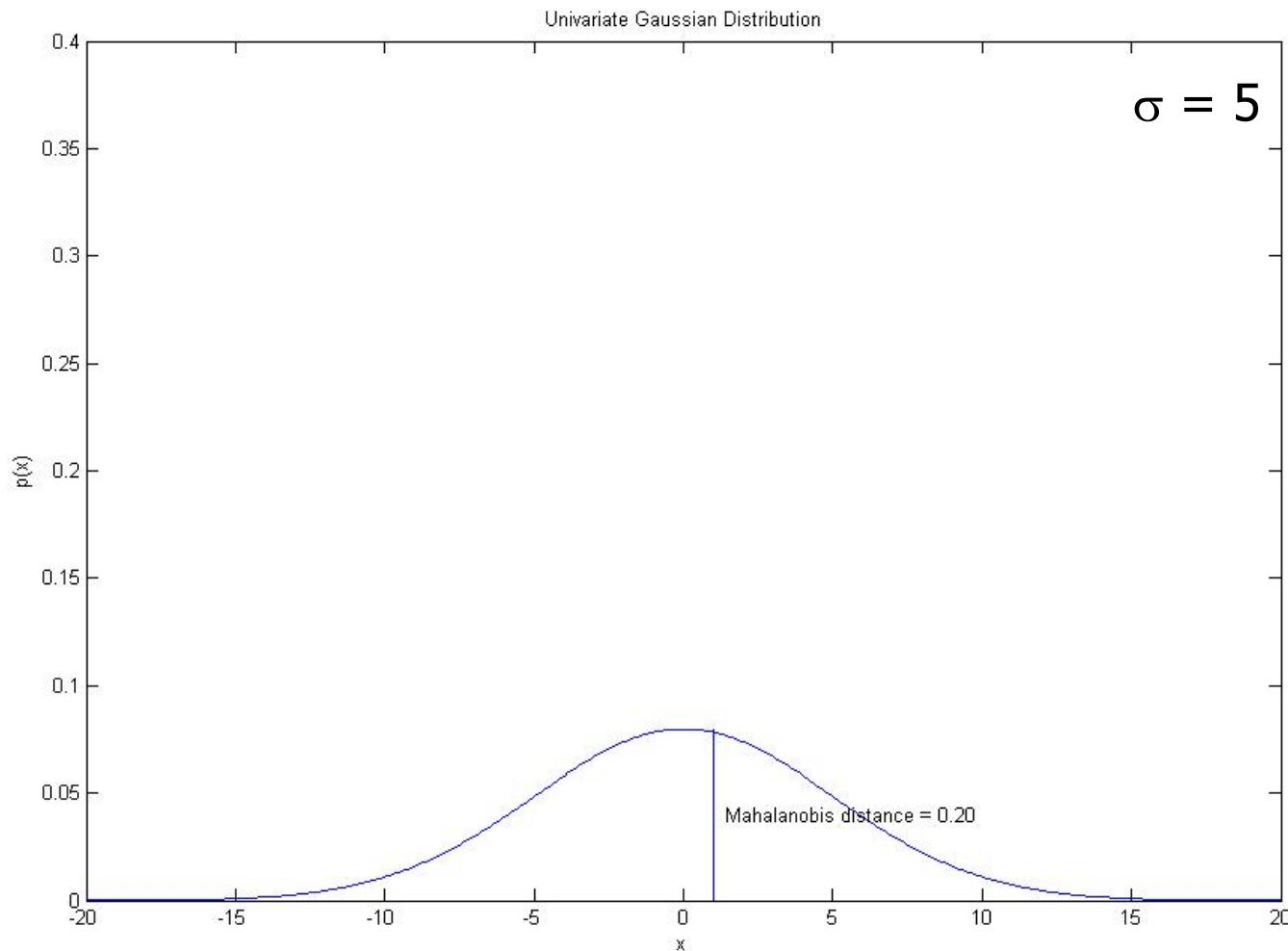
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



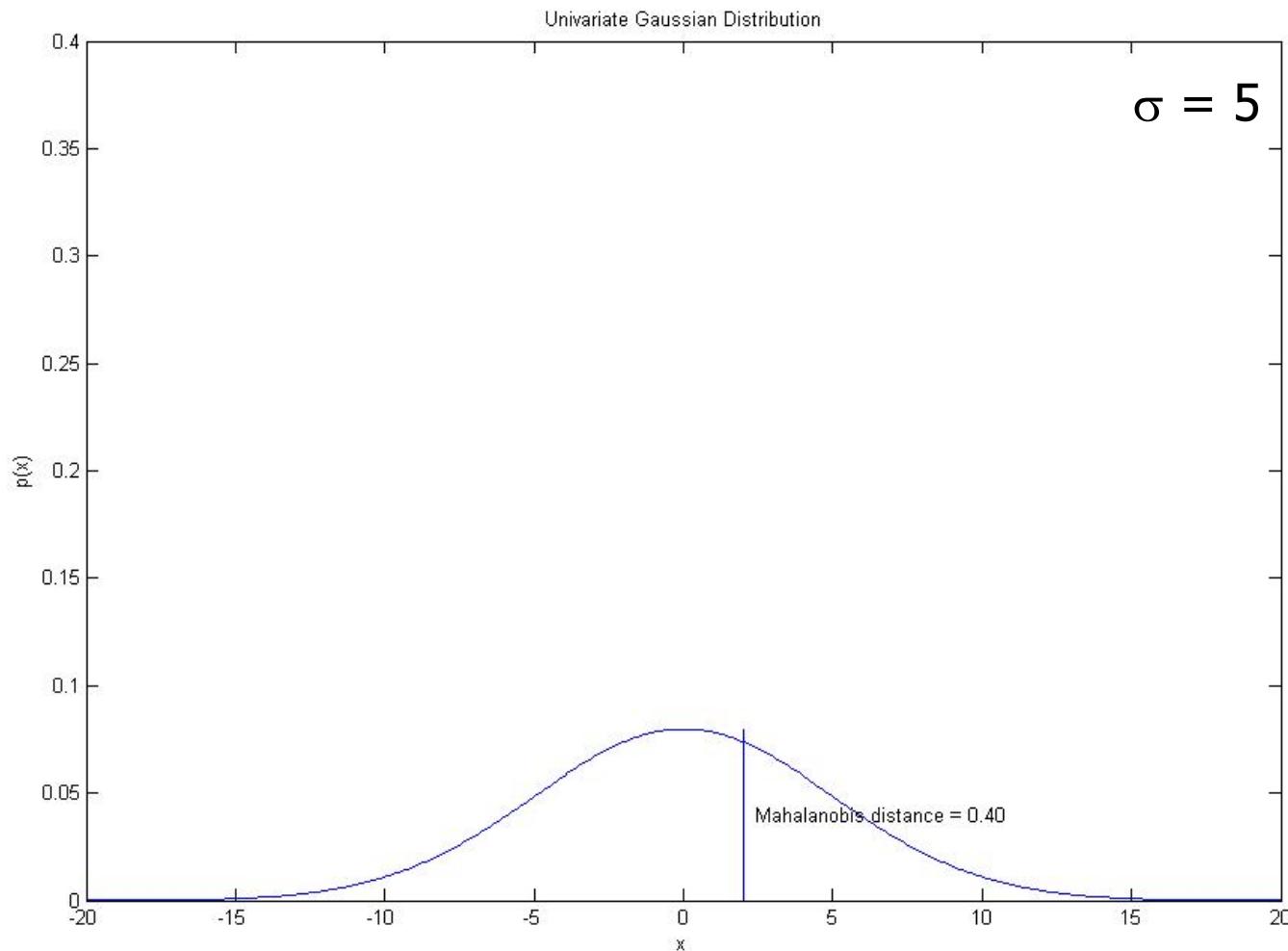
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



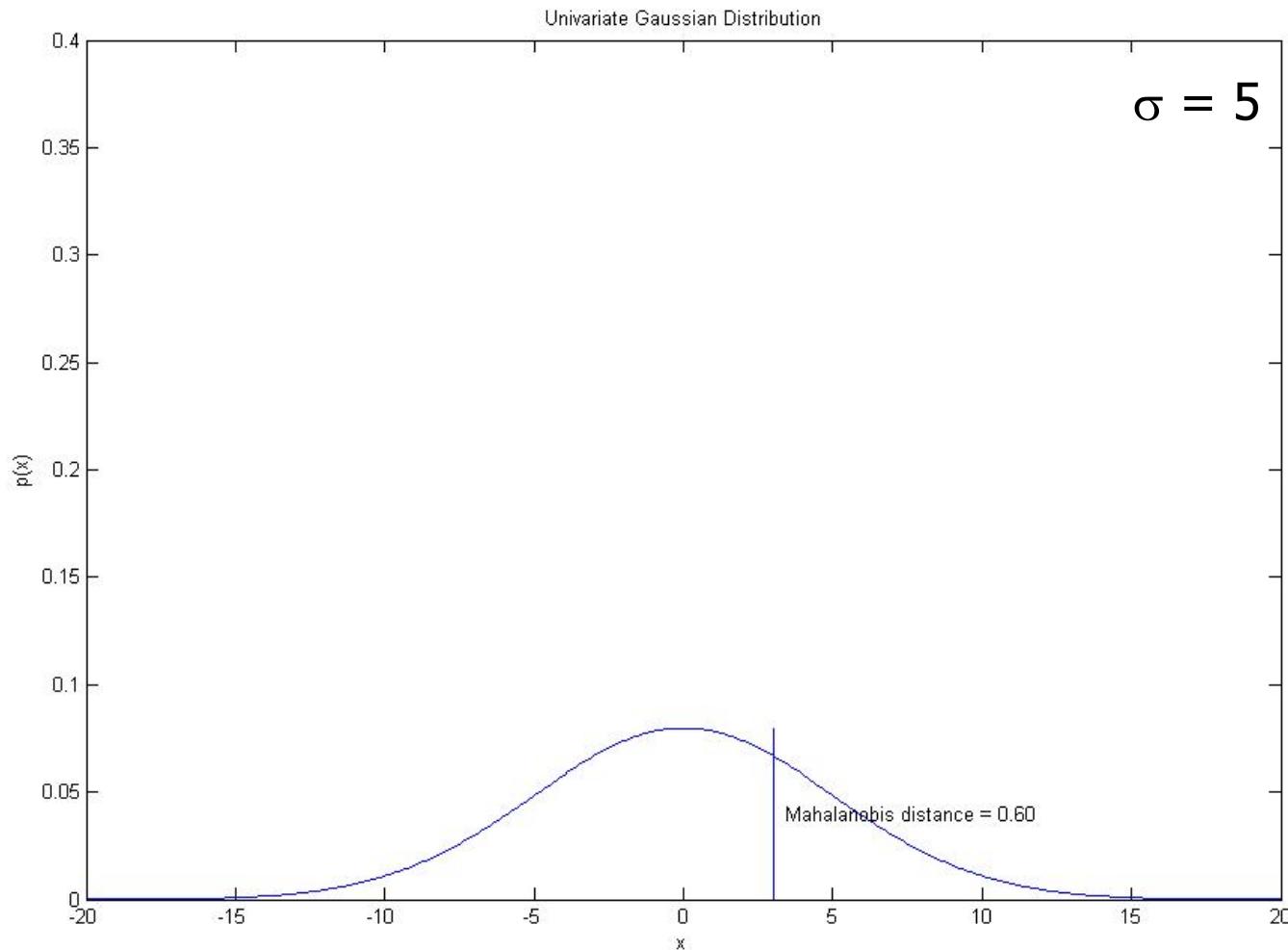
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



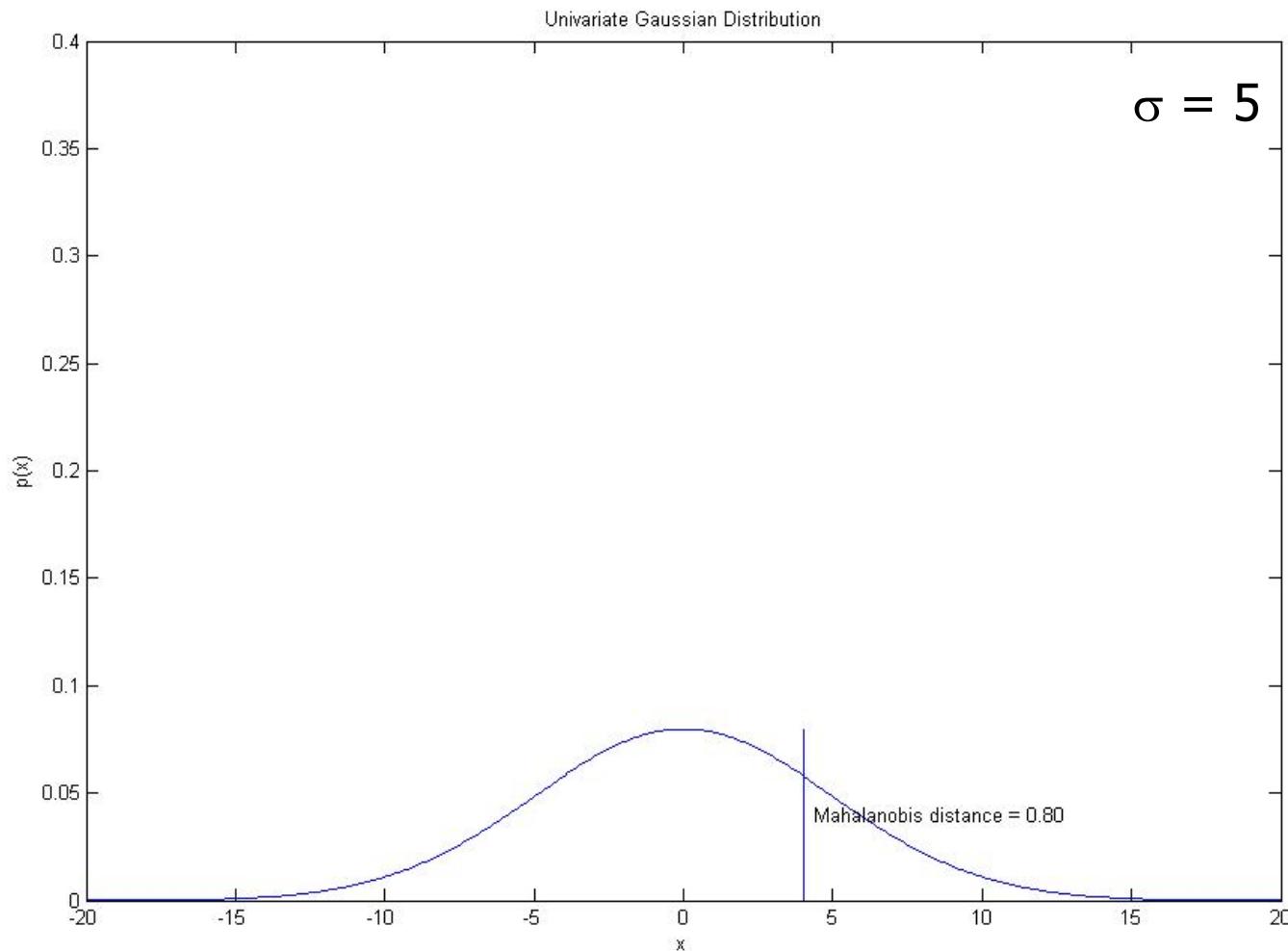
1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$

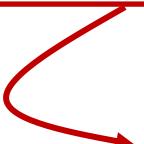


1.12 Mahalanobis Distance

$$d_M(x, \mu) = \frac{|x - \mu|}{\sigma}$$



1.13 **Multivariate Normal Densities**

 more than one random variable

1.13 Multivariate Normal Densities

We consider d normally distributed random variables x_i , each with its own

mean μ_i and variance σ_i^2 : $p(x_i) \sim N(\mu_i, \sigma_i^2)$

1.13 Multivariate Normal Densities

We consider d normally distributed random variables x_i , each with its own

mean μ_i and variance σ_i^2 : $p(x_i) \sim N(\mu_i, \sigma_i^2)$

Description of the joint probability density function?

$$p(x_1, x_2, \dots, x_d)$$

1.13 Multivariate Normal Densities

We consider d normally distributed random variables x_i , each with its own

mean μ_i and variance σ_i^2 : $p(x_i) \sim N(\mu_i, \sigma_i^2)$

Description of the joint probability density function?

$$p(x_1, x_2, \dots, x_d)$$

Case 1: Statistically independent variables

$$p(x_1, x_2, \dots, x_d) = \prod_{i=1}^d p(x_i) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2}$$

Multivariate Gaussian Density With Independent Components

1.13 Multivariate Normal Densities

We consider d normally distributed random variables x_i , each with its own

mean μ_i and variance σ_i^2 : $p(x_i) \sim N(\mu_i, \sigma_i^2)$

Description of the joint probability density function?

$$p(x_1, x_2, \dots, x_d)$$

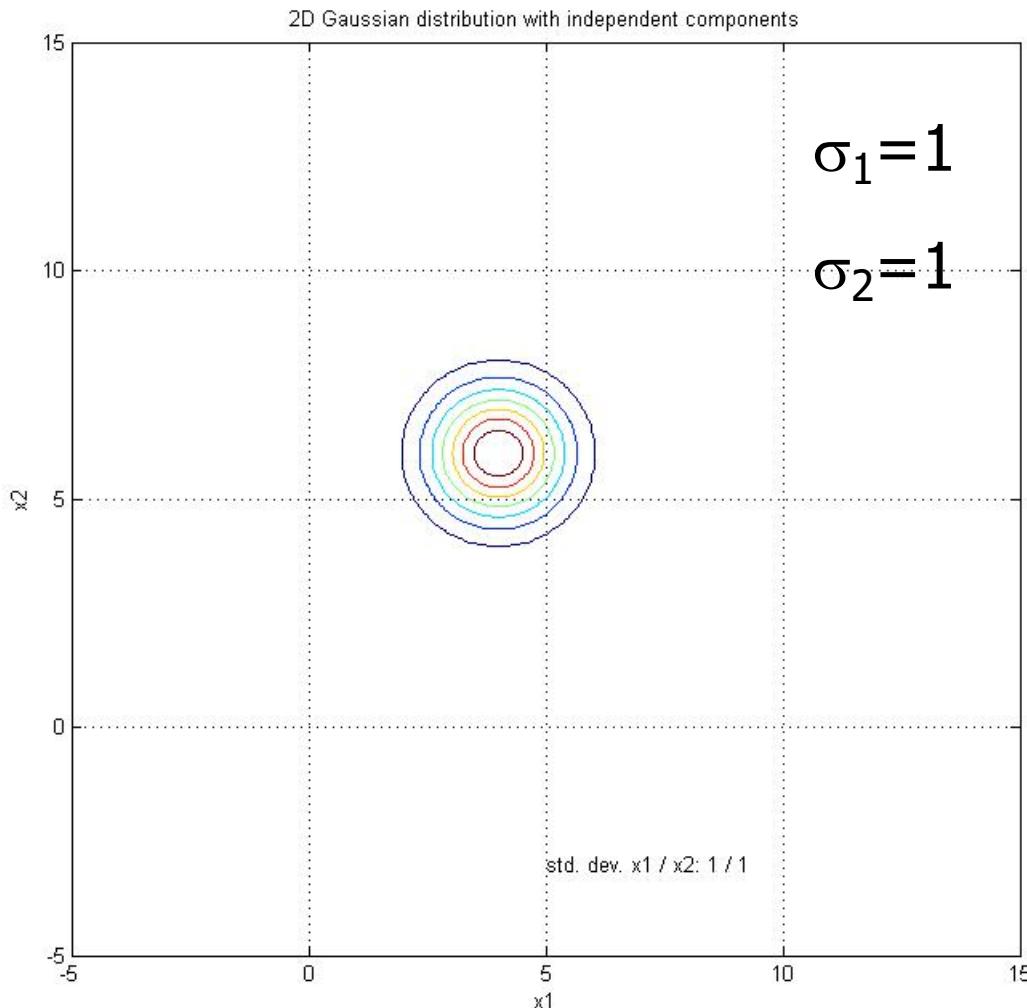
Case 1: Statistically independent variables

$$p(x_1, x_2, \dots, x_d) = \prod_{i=1}^d p(x_i) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_i^2}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2}$$

2-D case: $p(x_1, x_2) = p(x_1) \cdot p(x_2)$

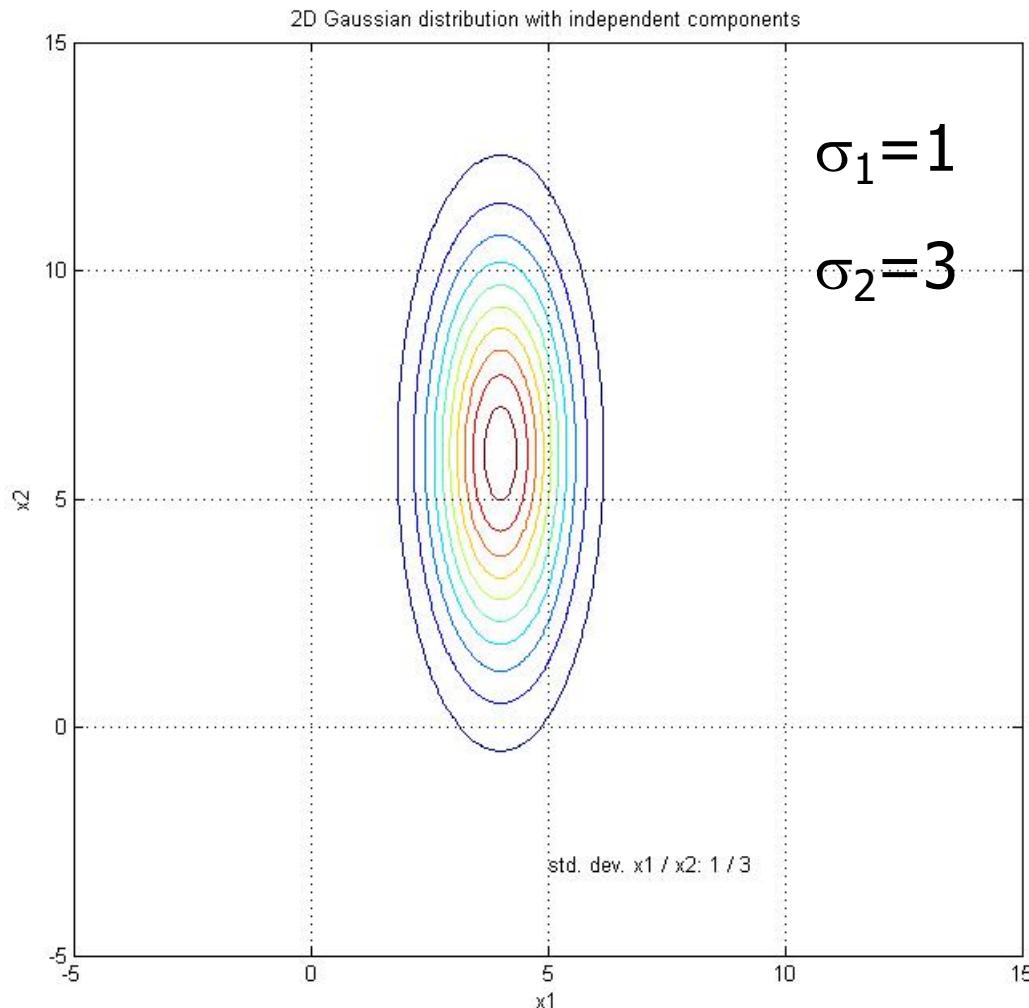
1.13 Multivariate Normal Densities

2D-Example



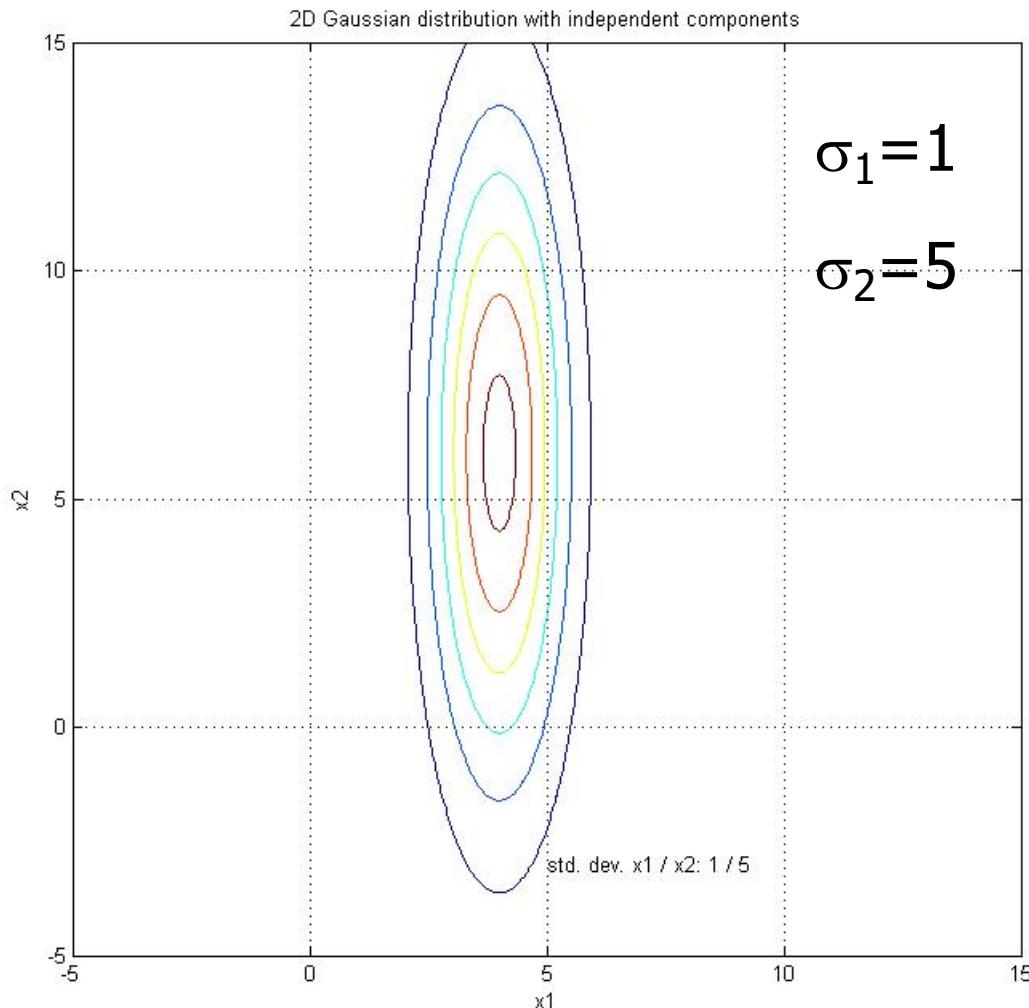
1.13 Multivariate Normal Densities

2D-Example



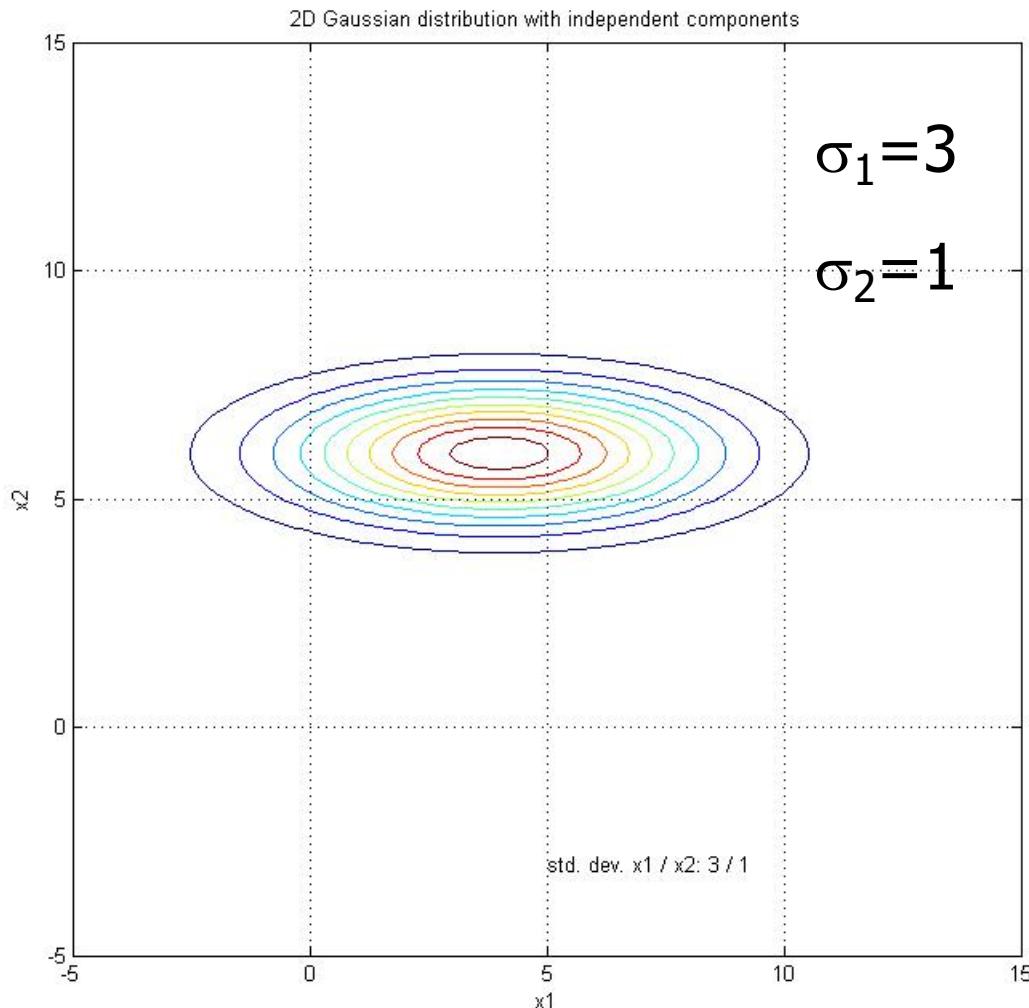
1.13 Multivariate Normal Densities

2D-Example



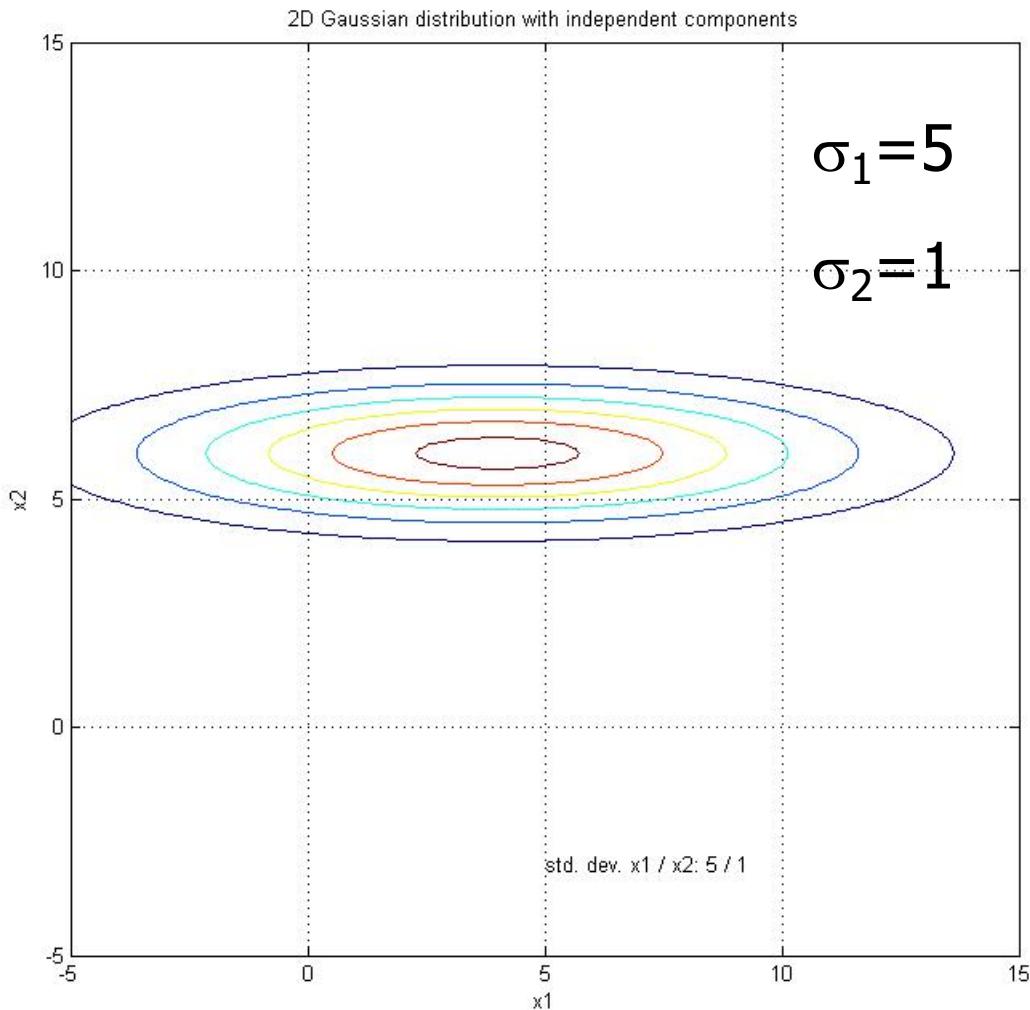
1.13 Multivariate Normal Densities

2D-Example



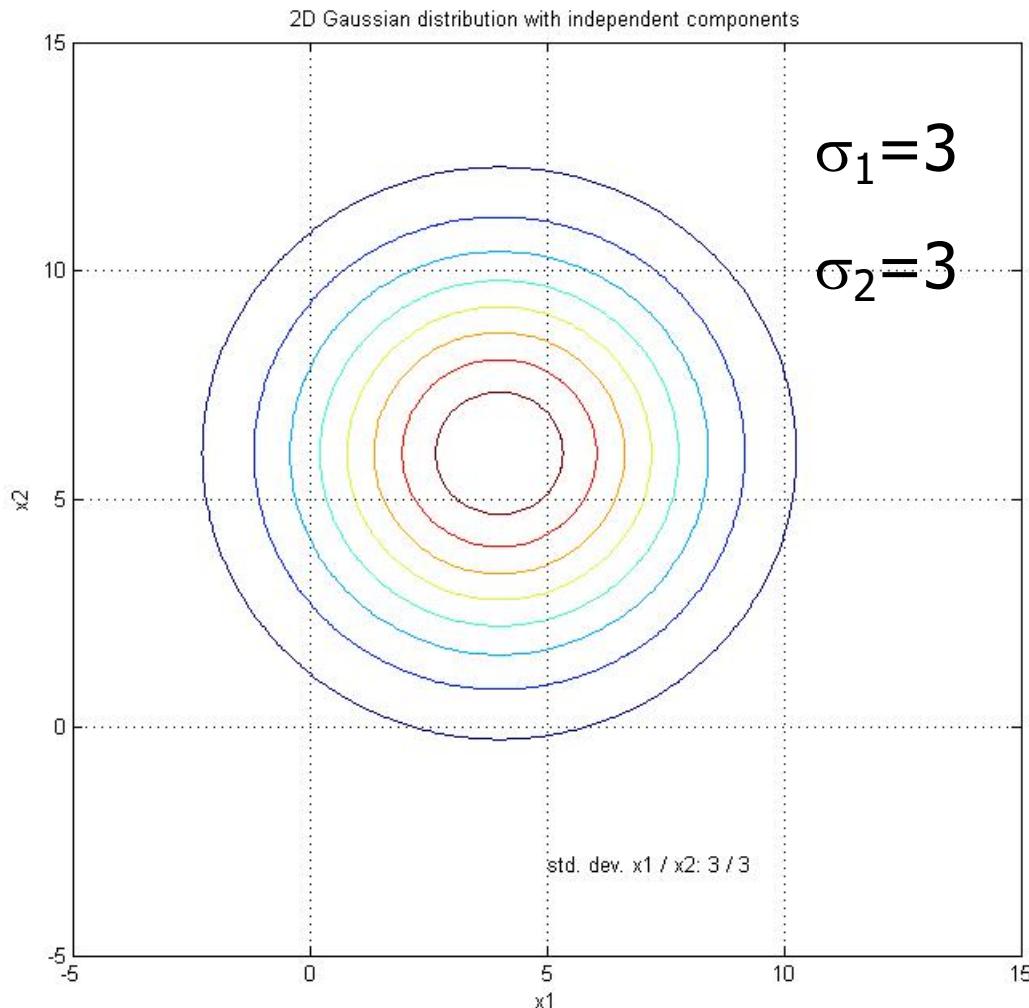
1.13 Multivariate Normal Densities

2D-Example



1.13 Multivariate Normal Densities

2D-Example



1.13 Multivariate Normal Densities

Why do we need this?

→ Application example: **Fish classification**

Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Practical realization: Data is organized in a Matlab array

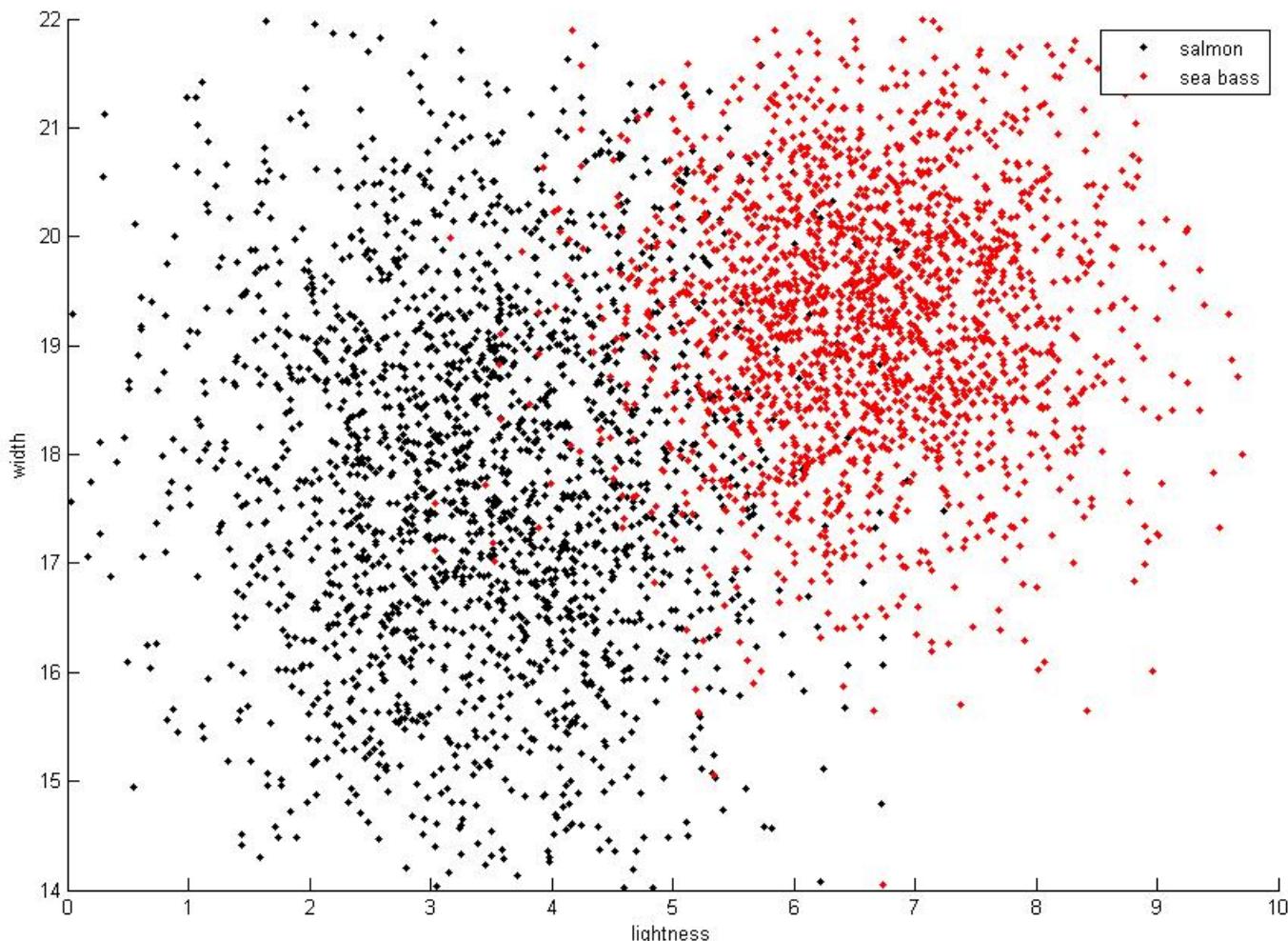
```
>> salmon(1:5,:)
```

ans =

3.6642	20.0939
4.9162	15.3323
3.9818	16.6113
5.3453	15.4816
2.5478	18.4773

First salmon: Lightness = 3.6642, Width = 20.0939

Example: Representation of fish samples using Gaussian Densities



Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2a. Compute mean and variance of each random variable:

```
>> load salmon_artificial.mat          % available in OLAT material folder  
>> m1 = mean(salmon(:,1));           % lightness mean  
>> s1 = sqrt(var(salmon(:,1)));       % lightness standard deviation  
>> m2 = mean(salmon(:,2));           % width mean  
>> s2 = sqrt(var(salmon(:,2)));       % width standard deviation
```

Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2b. Compute marginal distributions

```
% Compute lightness values
```

```
>> pts1 = min(salmon(:,1)):0.1:max(salmon(:,1));
```

```
% Compute width values
```

```
>> pts2 = min(salmon(:,2)):0.1:max(salmon(:,2));
```

```
% Compute marginals: px1 is p(lightness), px2 is p(width)
```

```
>> px1 = exp(-0.5*((pts1-m1)./s1).^2)./(sqrt(2*pi)*s1);
```

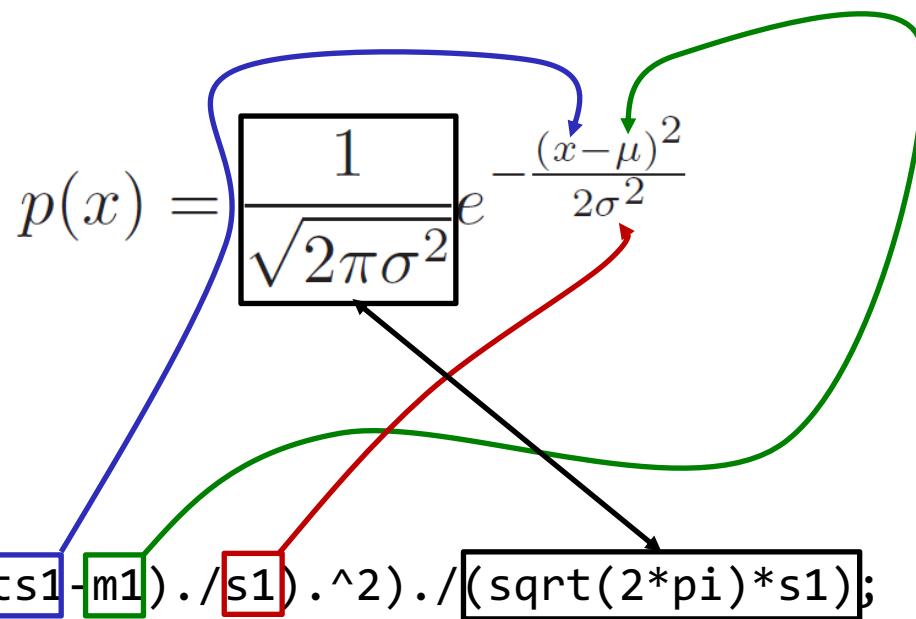
```
>> px2 = exp(-0.5*((pts2-m2)./s2).^2)./(sqrt(2*pi)*s2);
```

Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2b. Compute marginal distributions



```
>> px1 = exp(-0.5*((pts1-m1)./s1).^2)./(sqrt(2*pi)*s1);  
>> px2 = exp(-0.5*((pts2-m2)./s2).^2)./(sqrt(2*pi)*s2);
```

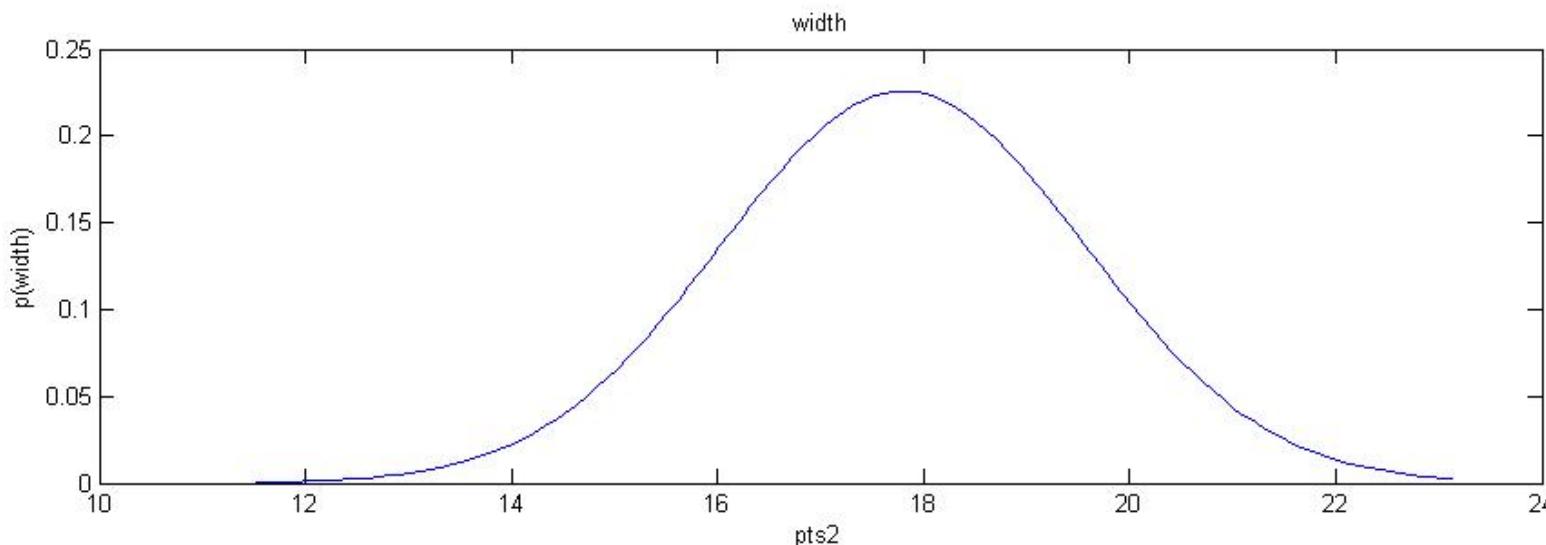
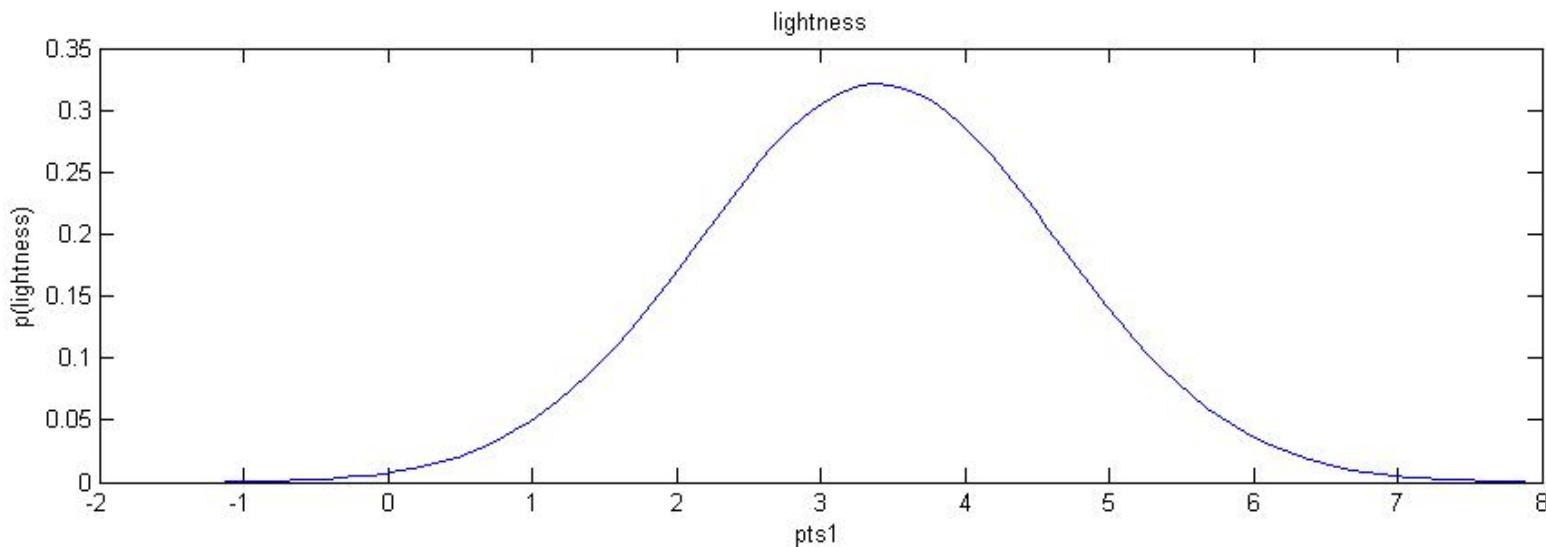
Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2b. Visualize marginal distributions

```
>> subplot(311);
>> plot(pts1, px1)
>> title('lightness')
>> xlabel('pts1')
>> ylabel('p(lightness)')
>> subplot(312);
>> plot(pts2, px2)
>> title('width')
>> xlabel('pts2')
>> ylabel('p(width)')
>> xlabel('pts2')
```



Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2c. Compute joint probability $p(\text{lightness}, \text{width})$ and plot contour lines

% joint probability

```
>> p_x1x2_sal = px2'*px1; % column vector * row vector = matrix
```

$$P(x, y) = P_x(x) P_y(y)$$



Example:

```
>> a = [1, 2]
>> b = [3, 4]
>> a * b
3 4
6 8
```

Example: Representation of fish samples using Gaussian Densities

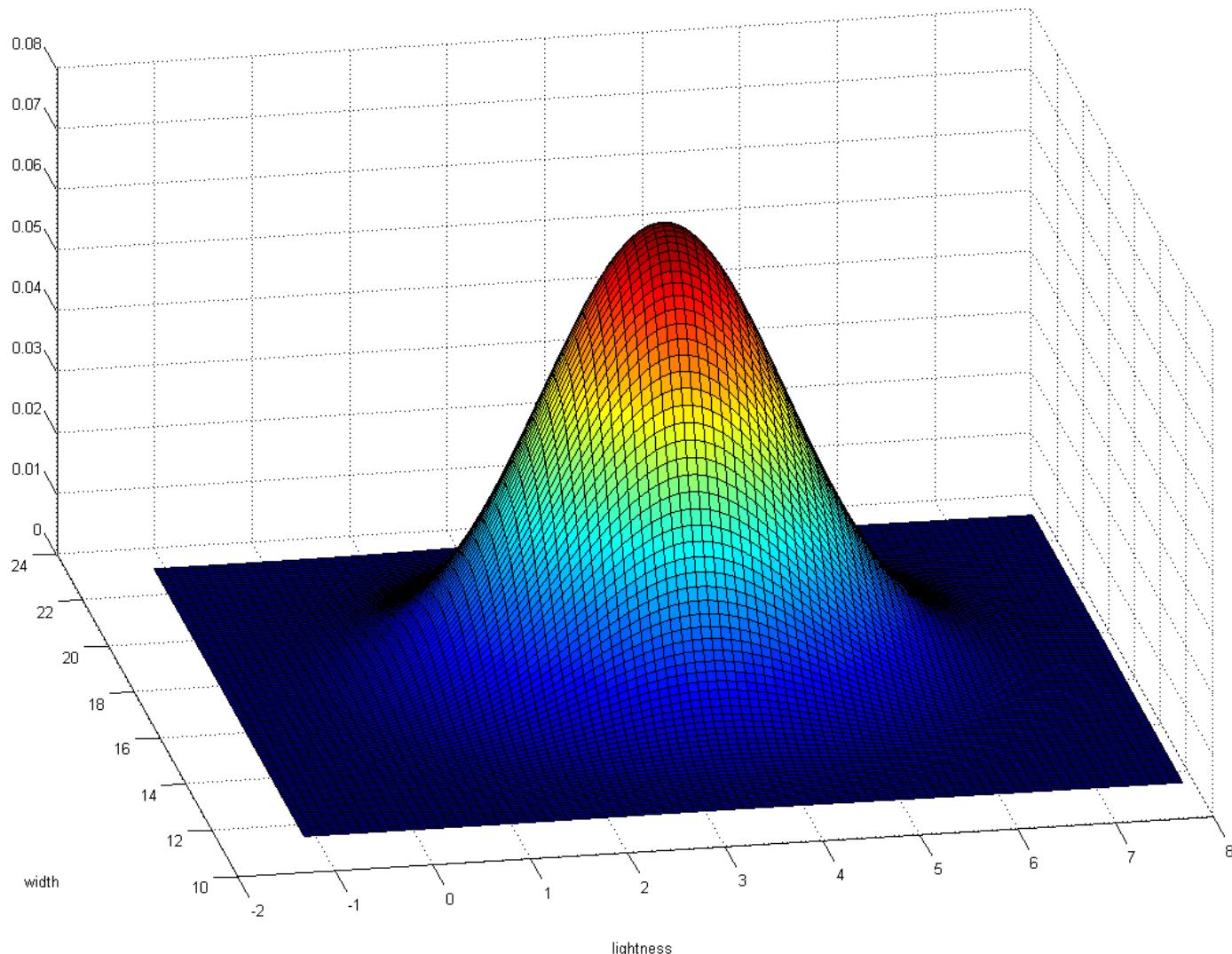
Step 1: Collection of 2D sensor data (lightness, width) for large training set

Step 2: Modeling of the data points in each class with a Gaussian Density

2c. Compute joint probability $p(\text{lightness}, \text{width})$ and plot contour lines

```
% visualize joint probability
>> subplot(313);
>> surf(pts1,pts2,p_x1x2_sal);
>> title('2-D Gaussian density with independent components');
>> xlabel('x1');
>> ylabel('x2');
>> xlabel('lightness');
>> ylabel('width');
```

2-D Gaussian density with independent components



Example: Representation of fish samples using Gaussian Densities

Step 1: Collection of 2D sensor data (lightness, width) for large training set

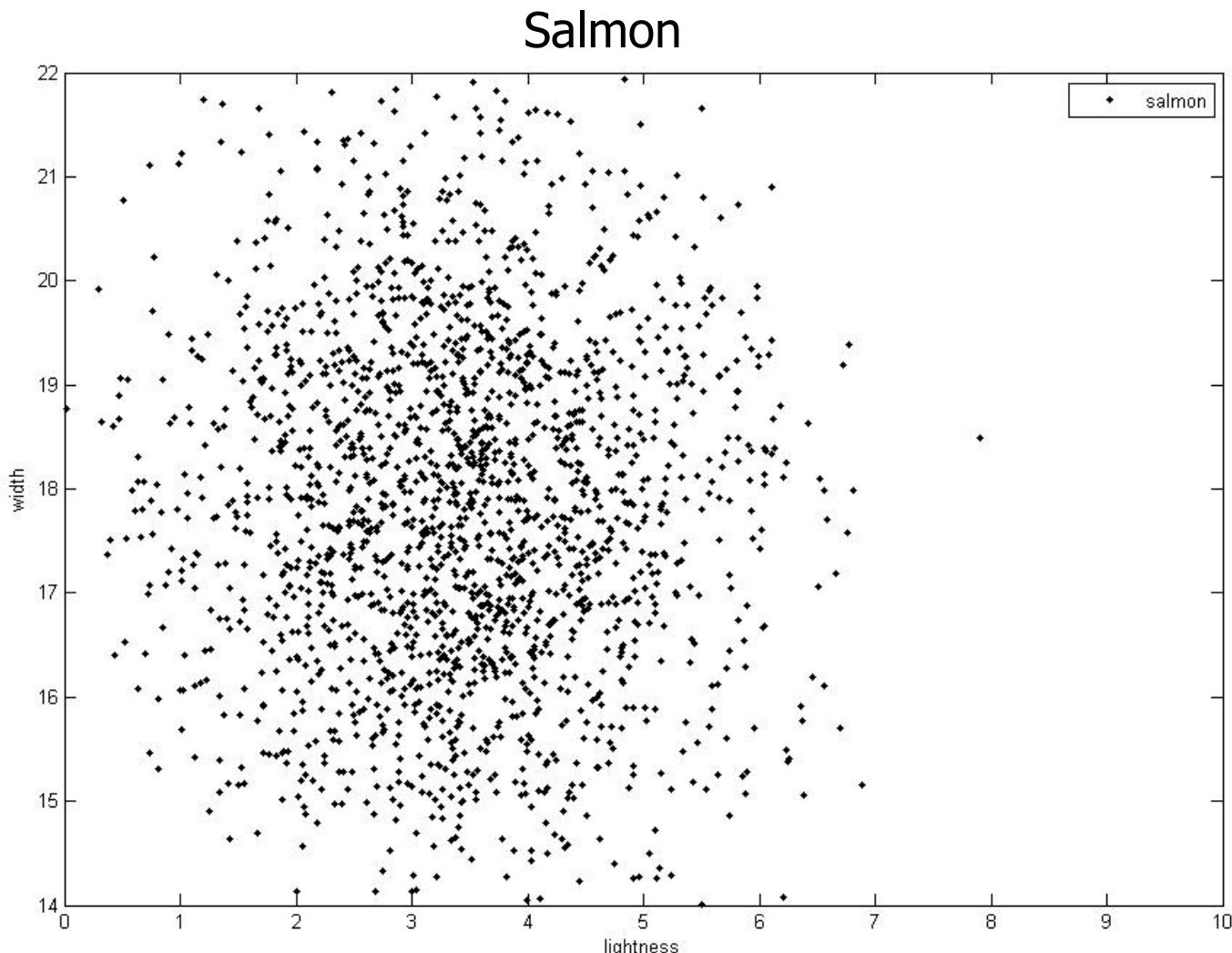
Step 2: Modeling of the data points in each class with a Gaussian Density

2c. Compute joint probability $p(\text{lightness}, \text{width})$ and plot contour lines

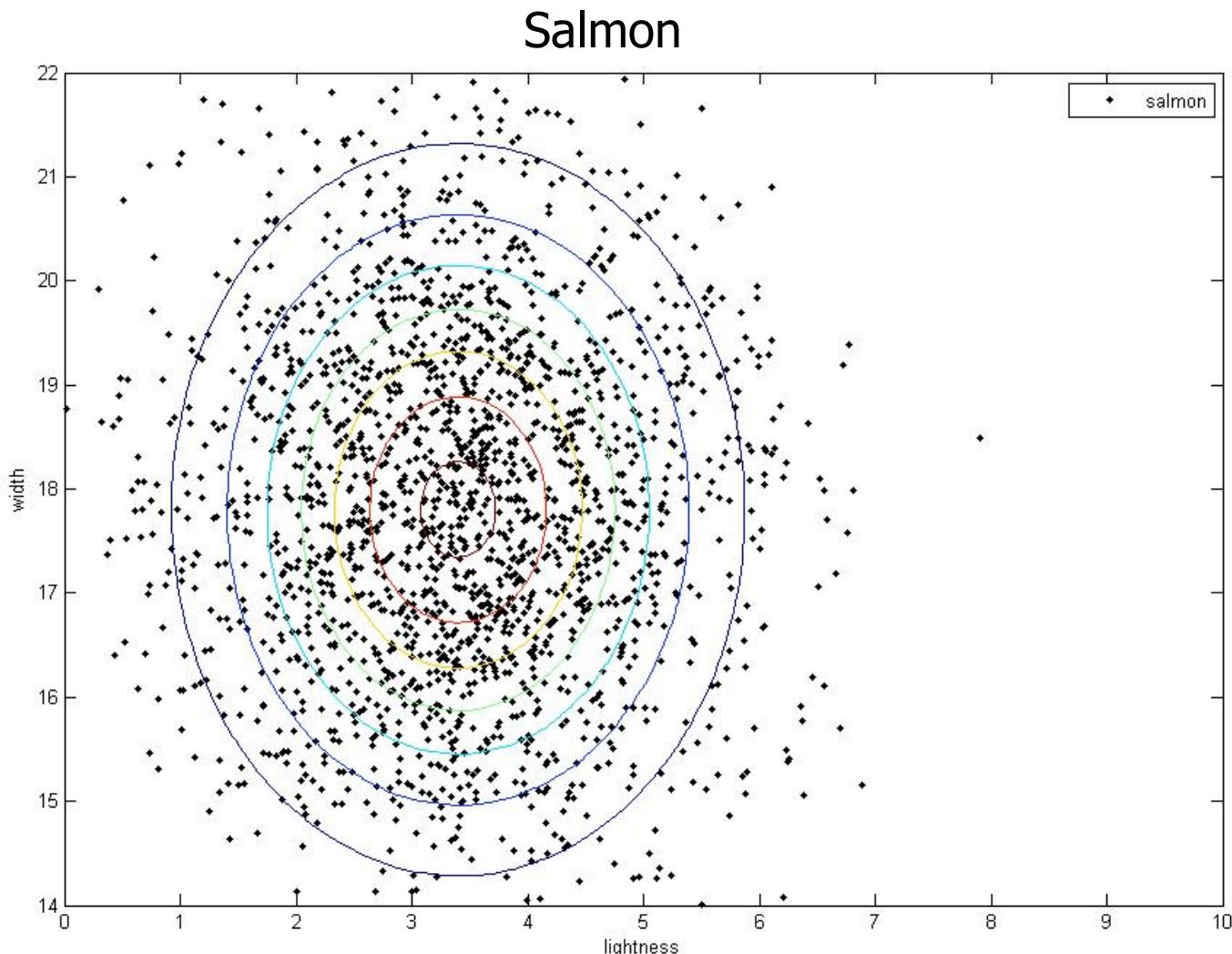
% contour lines

```
>> contour(pts1, pts2, p_x1x2_sal);
```

Example: Representation of fish samples using Gaussian Densities

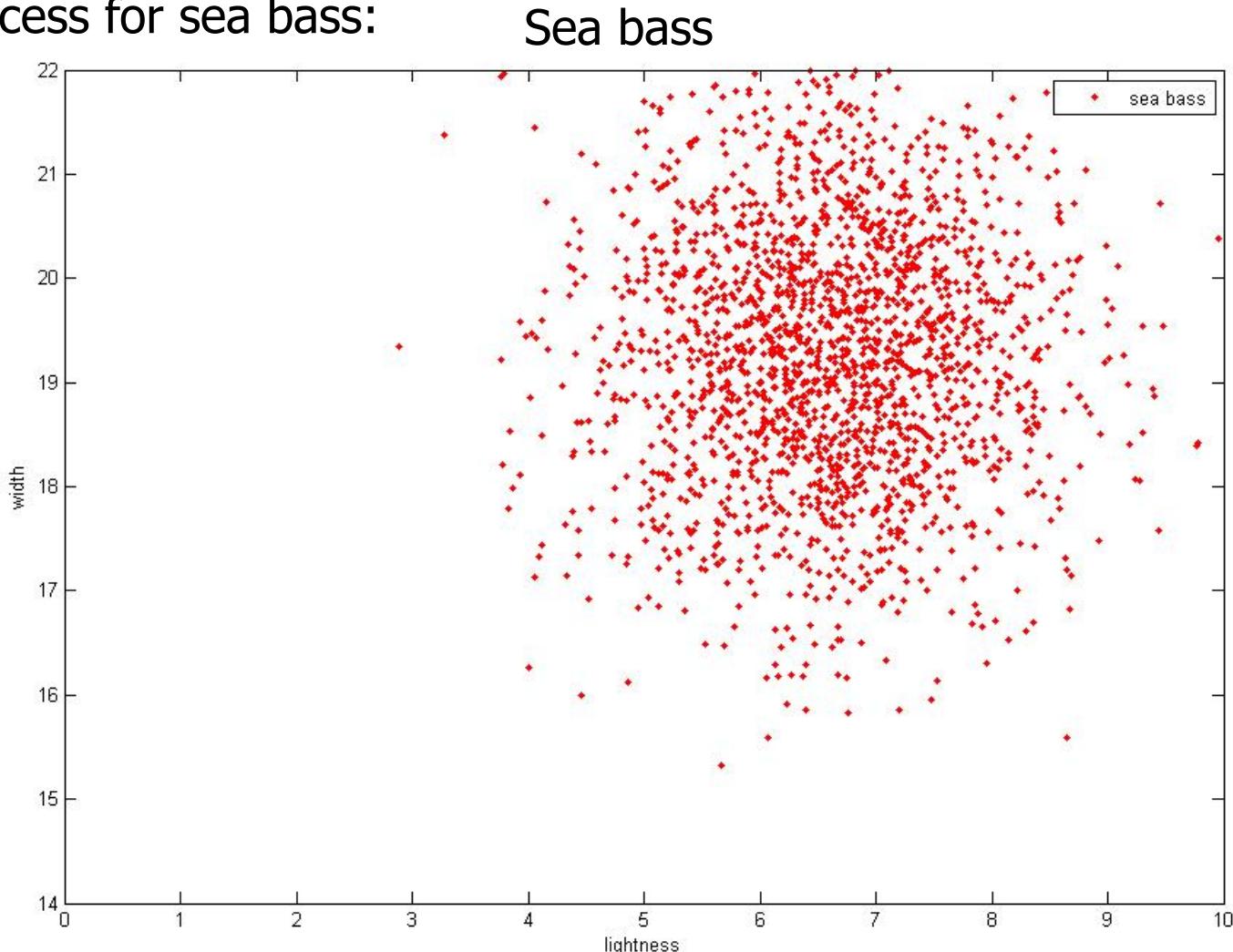


Example: Representation of fish samples using Gaussian Densities

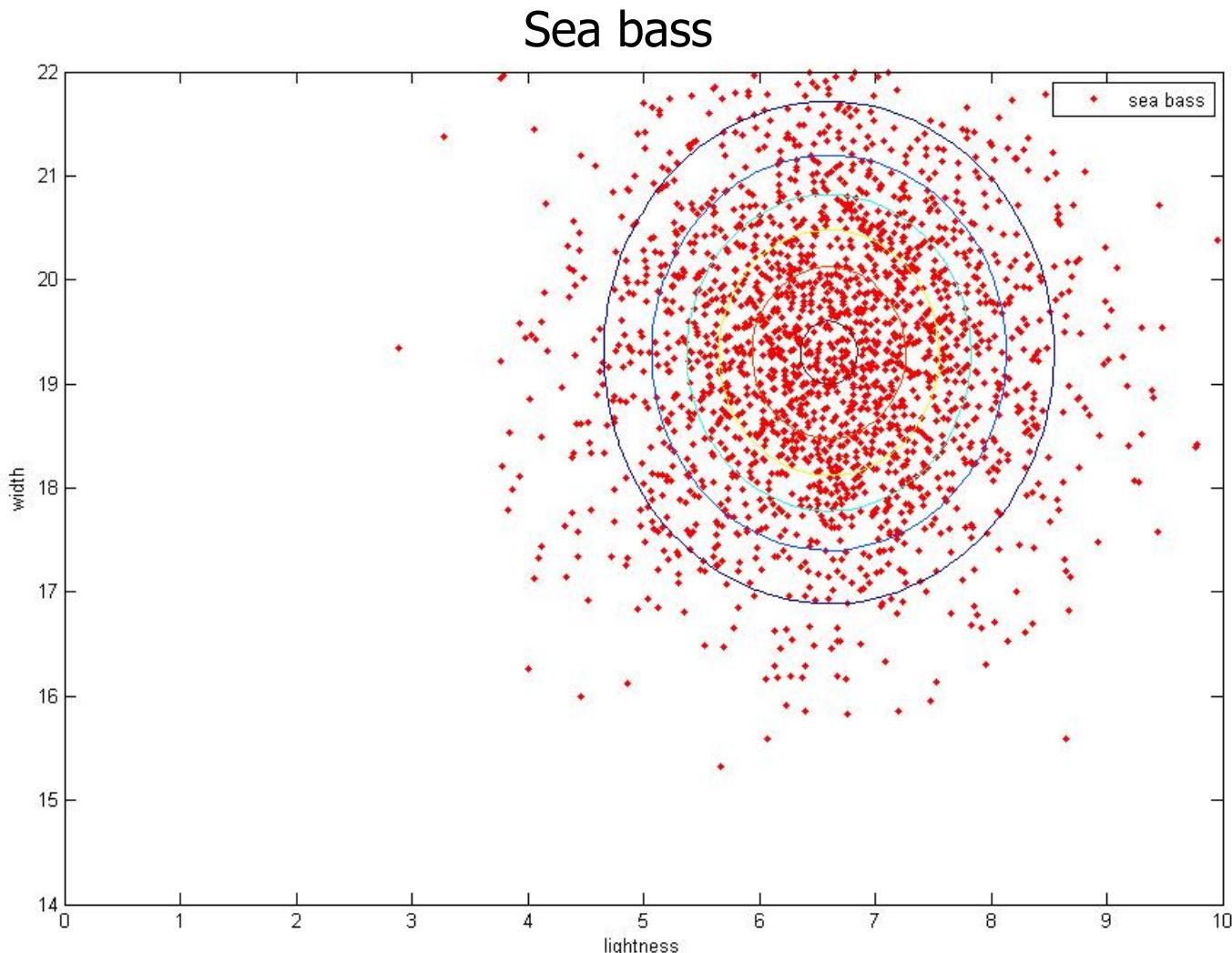


Example: Representation of fish samples using Gaussian Densities

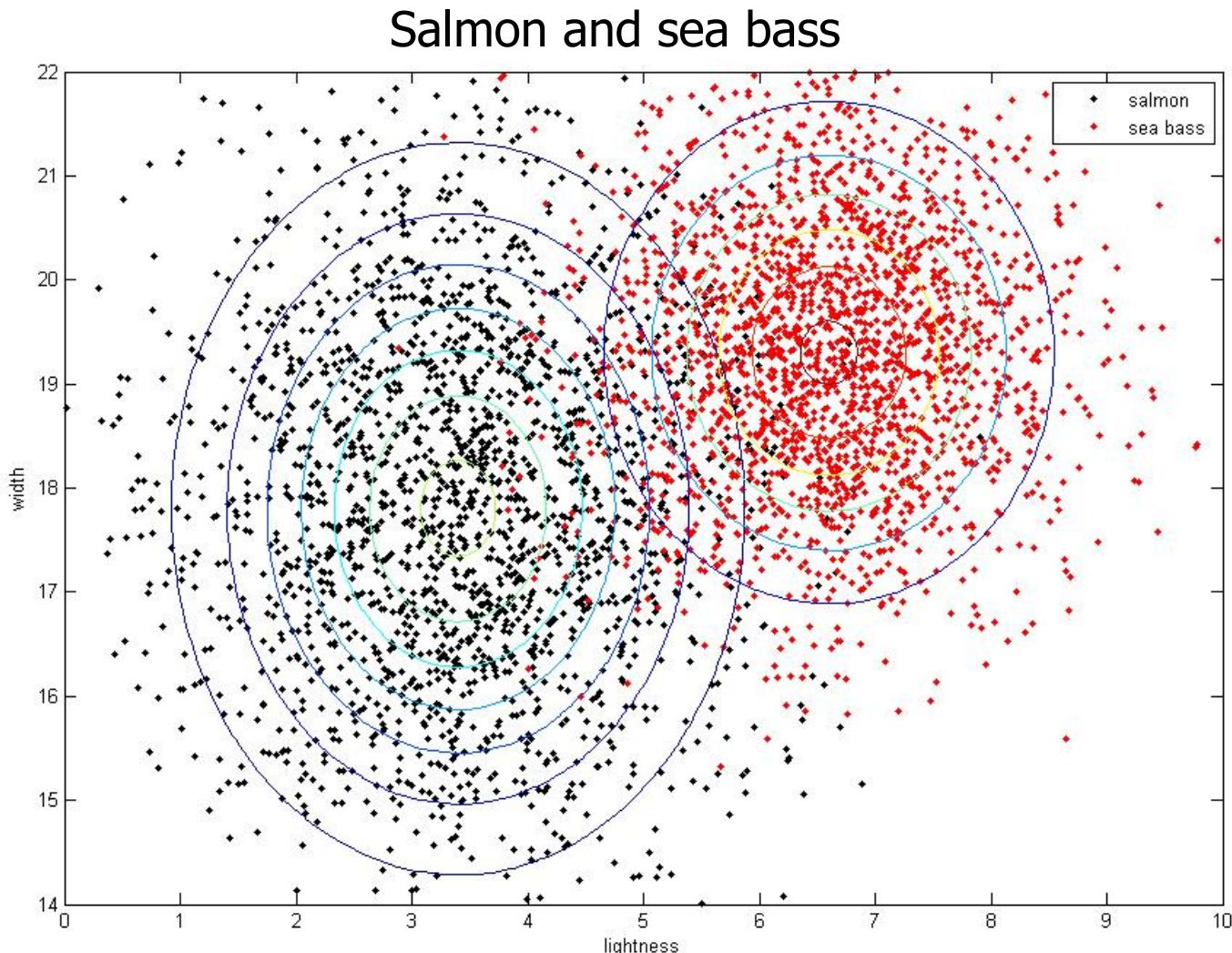
Repeat process for sea bass:



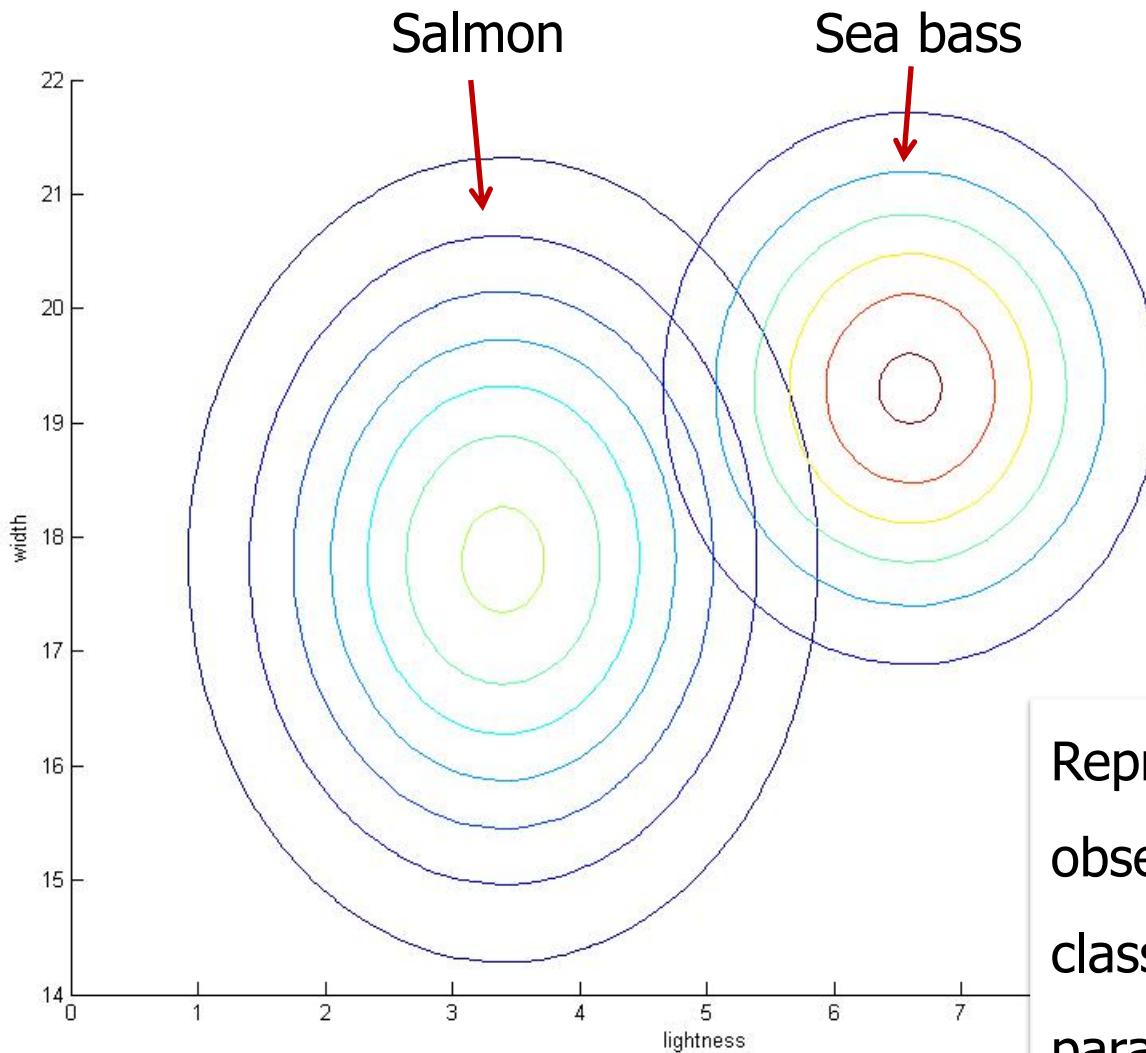
Example: Representation of fish samples using Gaussian Densities



Example: Representation of fish samples using Gaussian Densities



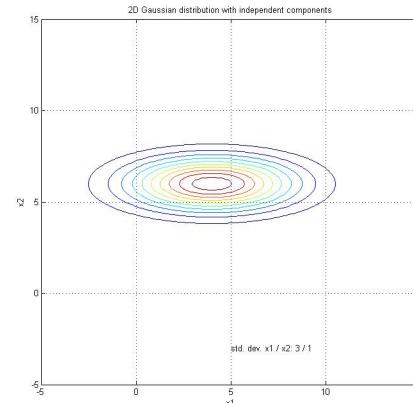
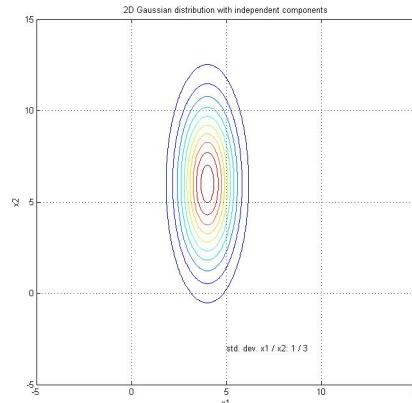
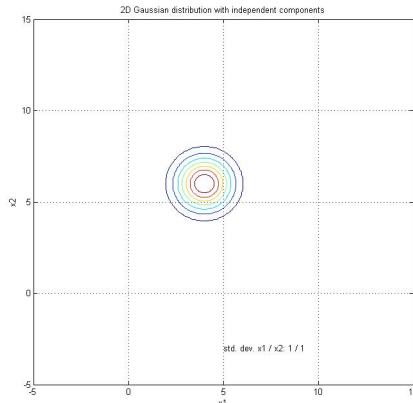
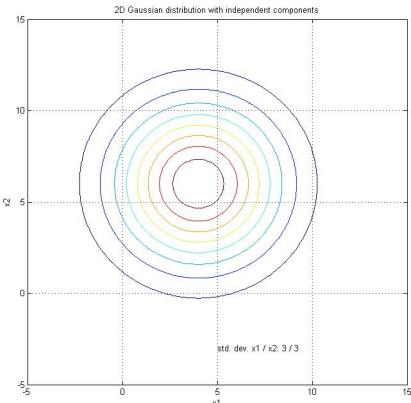
Example: Representation of fish samples using Gaussian Densities



Representation of all
observations of both
classes with only 8
parameters!

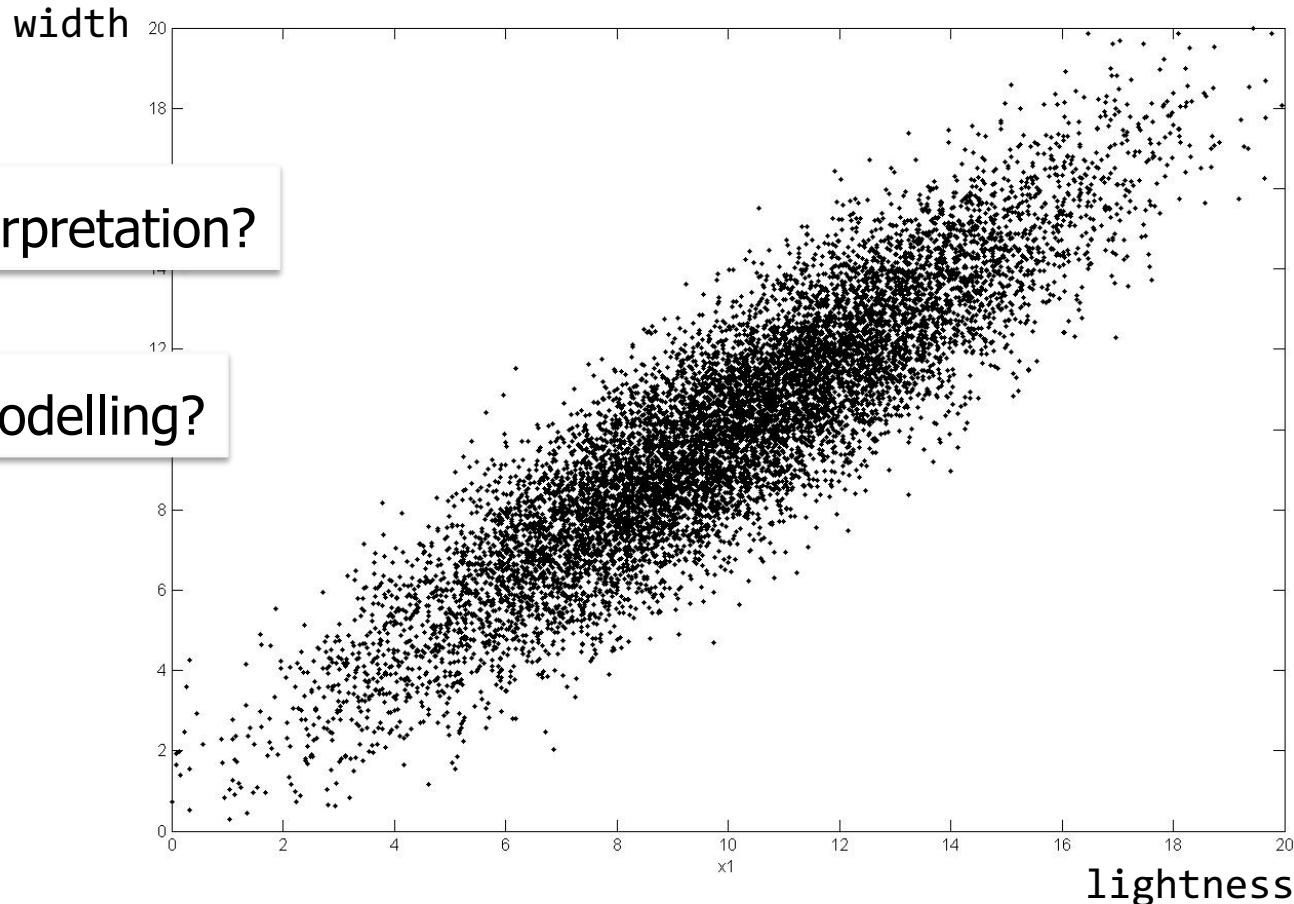
1.13 Multivariate Normal Densities

Multivariate Gaussian densities with **independent components** cover distributions of the following form:

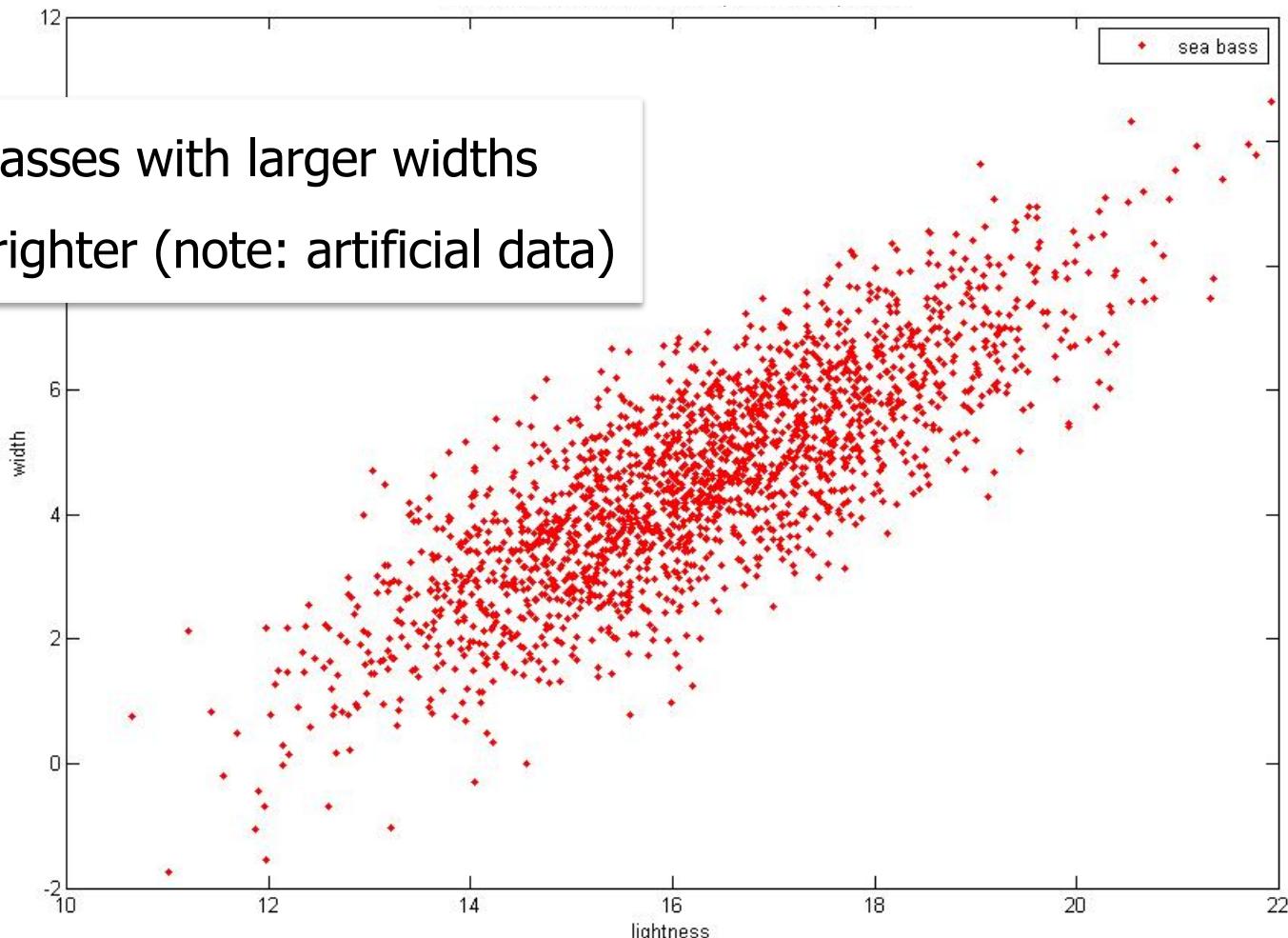


1.13 Multivariate Normal Densities

However, what about distributions like this?

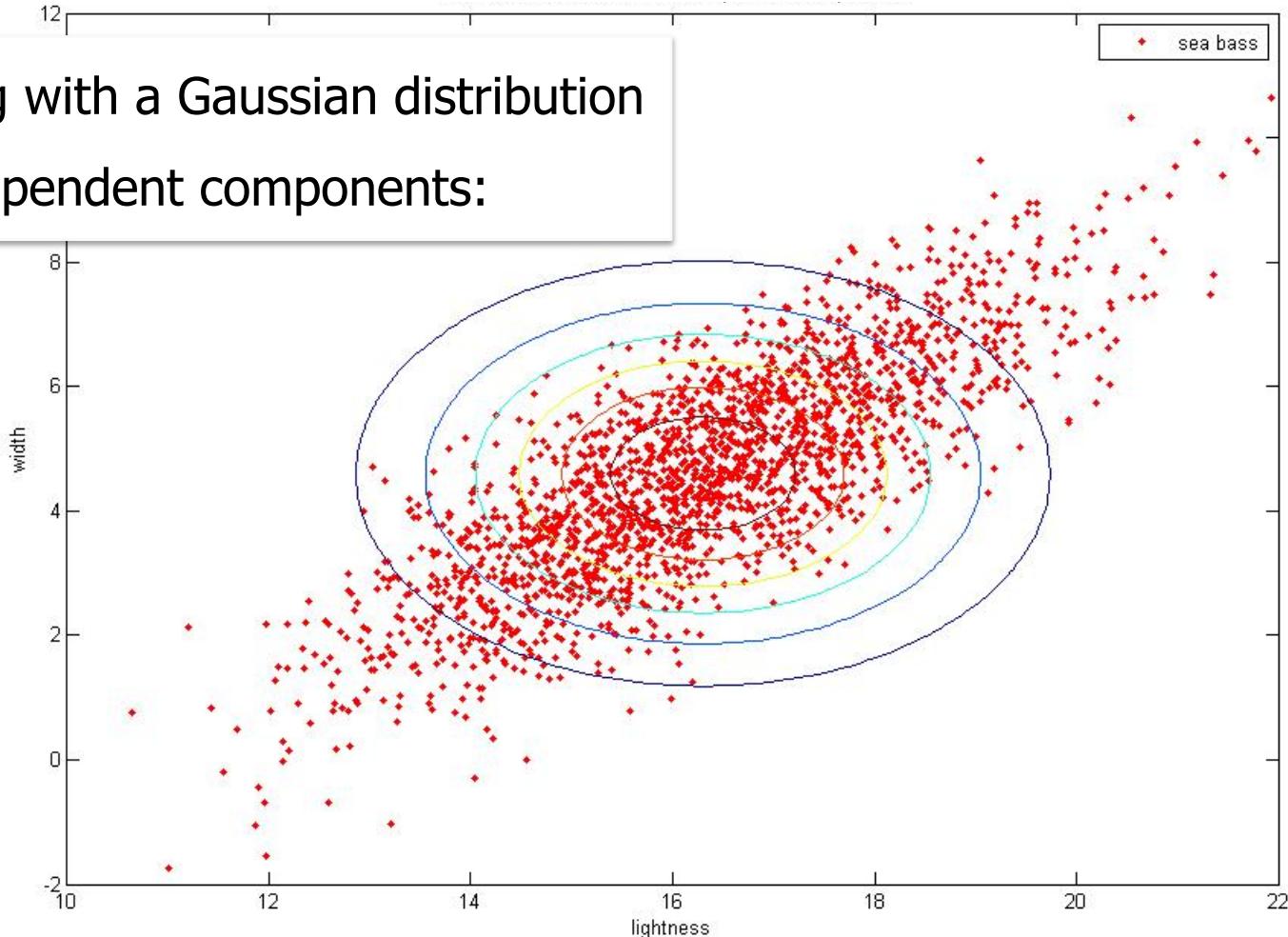


1.13 Multivariate Normal Densities



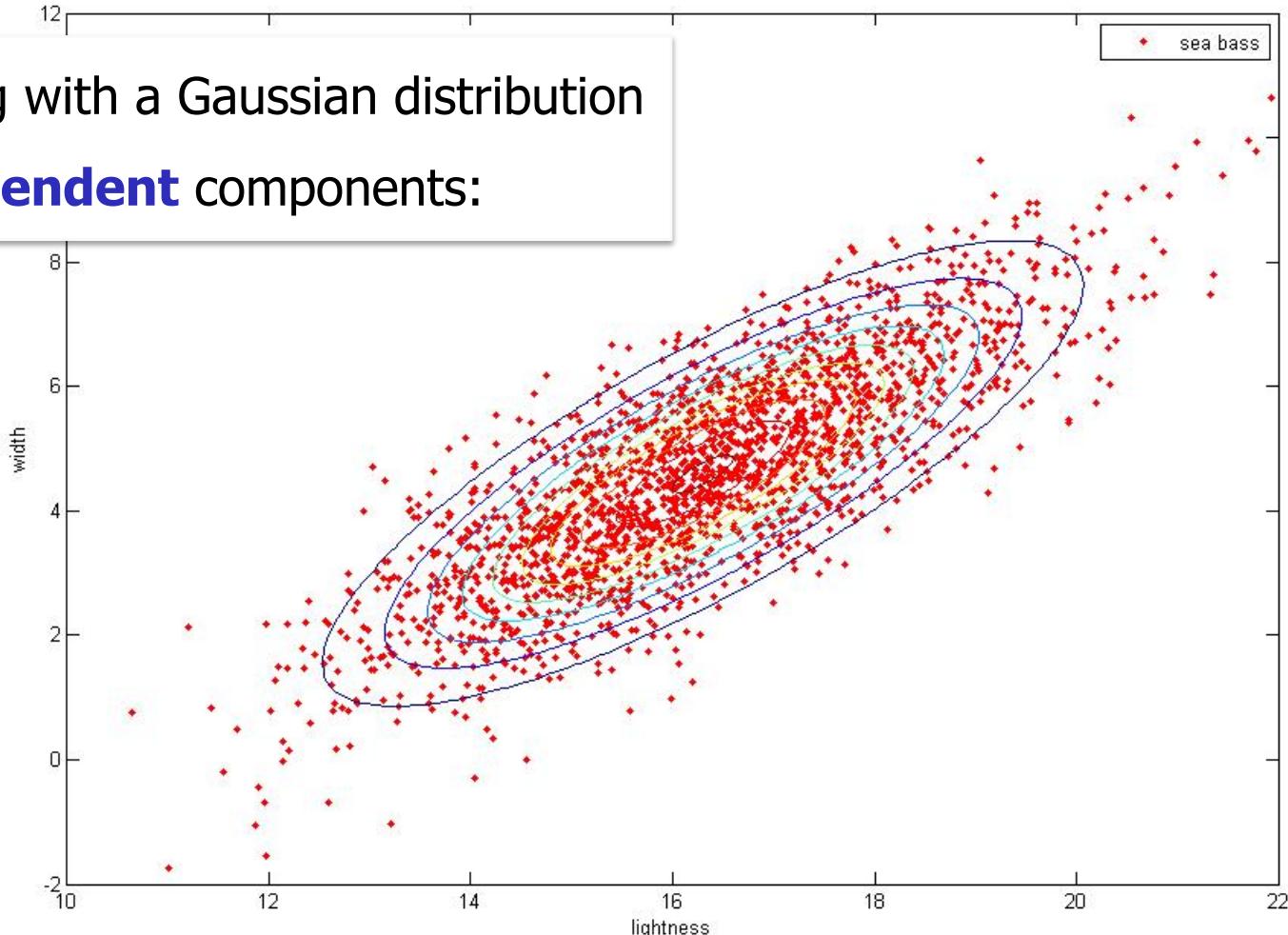
1.13 Multivariate Normal Densities

Modelling with a Gaussian distribution
with independent components:



1.13 Multivariate Normal Densities

Modelling with a Gaussian distribution
with **dependent** components:



1.13 Multivariate Normal Densities

Case 2: Statistically dependent variables

$$p(x_1, x_2, \dots, x_D) = \frac{1}{\sqrt{(2\pi)^D \cdot \det(\Sigma)}} \cdot e^{-\frac{1}{2} \cdot (x - \mu)^t \cdot \Sigma^{-1} \cdot (x - \mu)}$$

mean vector

determinant

covariance matrix

inverse matrix

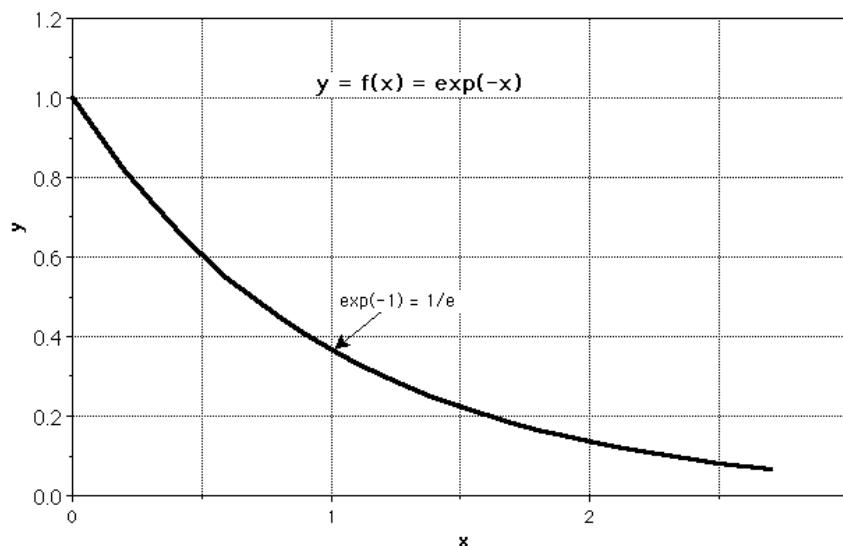
transpose
→ row vector

Multivariate Gaussian Density With Dependent Components

1.13 Multivariate Normal Densities

Case 2: Statistically dependent variables

$$p(x_1, x_2, \dots, x_D) = \frac{1}{\sqrt{(2\pi)^D \cdot \det(\Sigma)}} \cdot e^{-\frac{1}{2} \cdot (x - \mu)^t \cdot \Sigma^{-1} \cdot (x - \mu)}$$

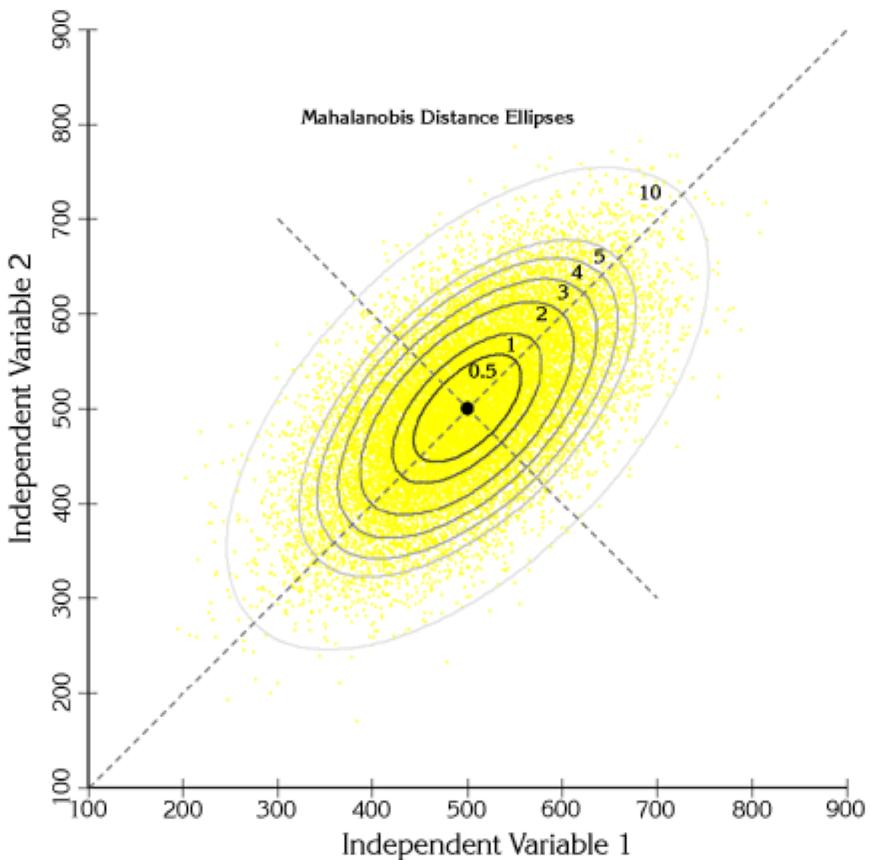


1.13 Multivariate Normal Densities

Mahalanobis Distance

$$d(\mathbf{x}, \boldsymbol{\mu}) = \sqrt{(\mathbf{x} - \boldsymbol{\mu})^t \cdot \boldsymbol{\Sigma}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu})}$$

Mahalanobis Distance between
observation vector \mathbf{x} and distribution
with $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$



1.13 Multivariate Normal Densities

Practical realization with Octave:

```
mx=[7 9]'; % mean vector
Cx=[5 3; 3 4]; % covariance matrix
x1=0:0.1:14;
x2=0:0.1:14;
for i=1:length(x1)
    for j=1:length(x2)
        p(j,i)=(1/(2*pi*det(Cx)^1/2))*  

            exp((-1/2)*([x1(i) x2(j)]-mx')*inv(Cx)*([x1(i);x2(j)]-mx));
    end
end
figure(1); mesh(x1,x2,p);
xlabel("x1");
ylabel("x2");
figure(2); contour(x1,x2,p);
```

1.13 Multivariate Normal Densities

Practical realization with Octave:

```
for i=1:length(x1)
    for j=1:length(x2)
        p(j,i)=1/(2*pi*det(Cx)^(1/2))*  

        exp((-1/2)*([x1(i) x2(j)]-mx')*inv(Cx)*([x1(i);x2(j)]-mx));
    end
end
```

Remark:

$D = 2$ in our case (2D observation vector)

$$\rightarrow \sqrt{2\pi^D} = 2\pi$$

$$p(x_1, x_2, \dots, x_D) = \frac{1}{\sqrt{(2\pi)^D \cdot \det(\Sigma)}} \cdot e^{-\frac{1}{2} \cdot (x - \mu)^t \cdot \Sigma^{-1} \cdot (x - \mu)}$$

1.13 Multivariate Normal Densities

Practical realization with Octave:

```
for i=1:length(x1)
    for j=1:length(x2)
        p(j,i)=(1/(2*pi*det(Cx)^(1/2)))*
            exp((-1/2)*([x1(i) x2(j)]-mx')*inv(Cx)*([x1(i);x2(j)]-mx));
    end
end
```

$$p(x_1, x_2, \dots, x_D) = \frac{1}{\sqrt{(2\pi)^D \cdot \det(\Sigma)}} \cdot e^{-\frac{1}{2} \cdot (x - \mu)^t \cdot \Sigma^{-1} \cdot (x - \mu)}$$

1.13 Multivariate Normal Densities

Practical realization with Octave:

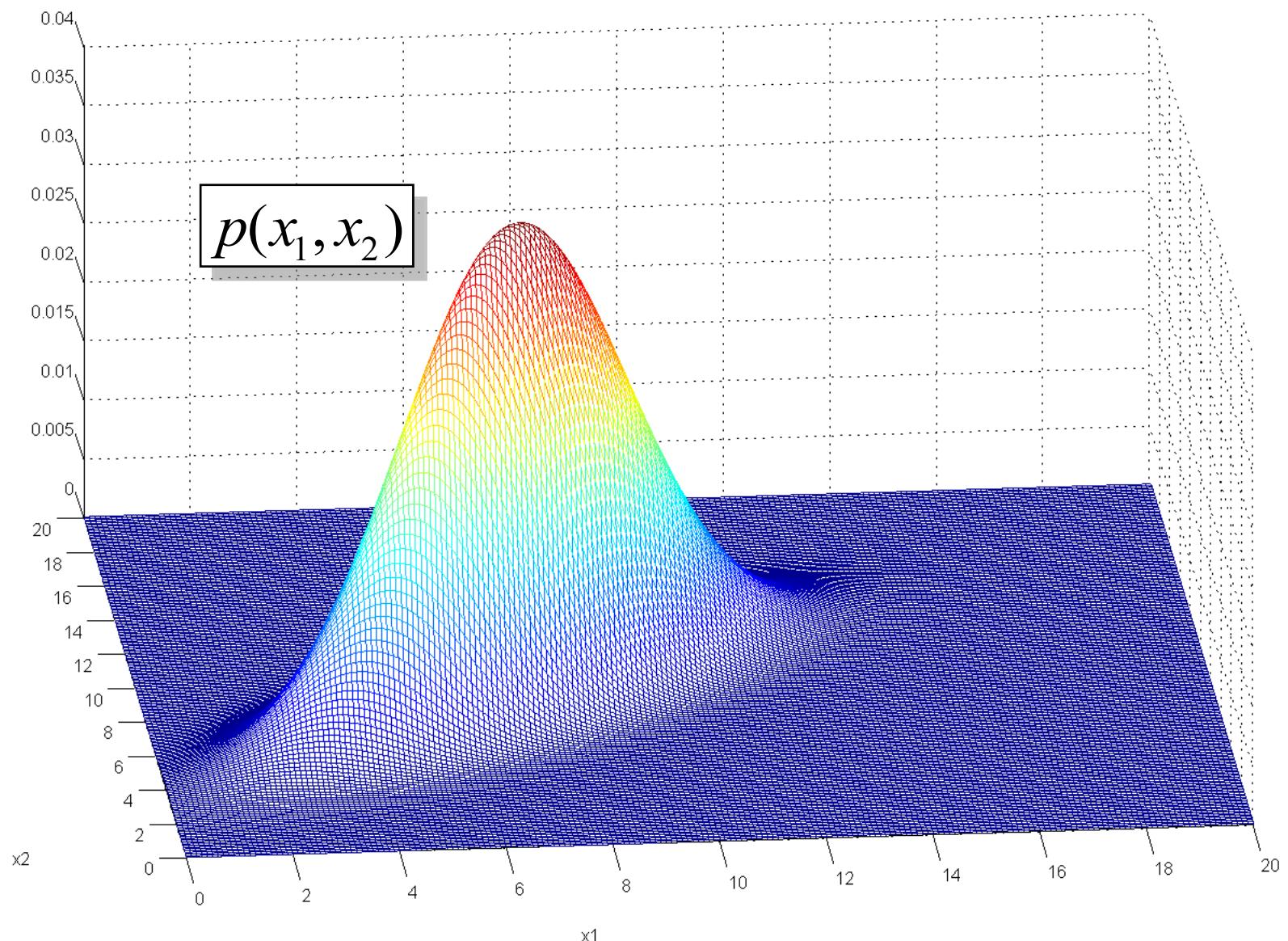
Example:

Covariance matrix:

$$\begin{pmatrix} 5 & 3 \\ 3 & 4 \end{pmatrix}$$

Mean vector:

$$\begin{pmatrix} 7 \\ 9 \end{pmatrix}$$



Multivariate Gaussian Distribution with dependent variables

1.13 Multivariate Normal Densities

$$p(x_1, x_2, \dots, x_D) = \frac{1}{\sqrt{(2\pi)^D \cdot \det(\Sigma)}} \cdot e^{-\frac{1}{2} \cdot (x - \mu)^T \cdot \Sigma^{-1} \cdot (x - \mu)}$$

Realization with function `mvnpdf` (statistics package)

```

mx=[1 -1]; % mean vector
Cx=[.9 .4; .4 .3]; % covariance matrix

% generate a grid to consider combinations of x1-x2 values
[X1,X2] = meshgrid(linspace(-1,3,25)', linspace(-3,1,25)');

% mvnpdf computes probability for each point i in X (i.e. X(i,:))
% therefore: collect all x1-x2-value combinations in X
X = [X1(:) X2(:)];

% result: probability vector, each element contains prob. of one grid point
p = mvnpdf(X, mx, Cx);

% for display we must reshape p into a 25 x 25 plane
surf(X1,X2,reshape(p,25,25));

```

1.13 Multivariate Normal Densities

Simplified example for usage of meshgrid

```
[X1, X2] = meshgrid(linspace(-1,3,5)',linspace(-3,1,5)')
```

X1 =

-1	0	1	2	3
-1	0	1	2	3
-1	0	1	2	3
-1	0	1	2	3
-1	0	1	2	3

X2 =

-3	-3	-3	-3	-3
-2	-2	-2	-2	-2
-1	-1	-1	-1	-1
0	0	0	0	0
1	1	1	1	1

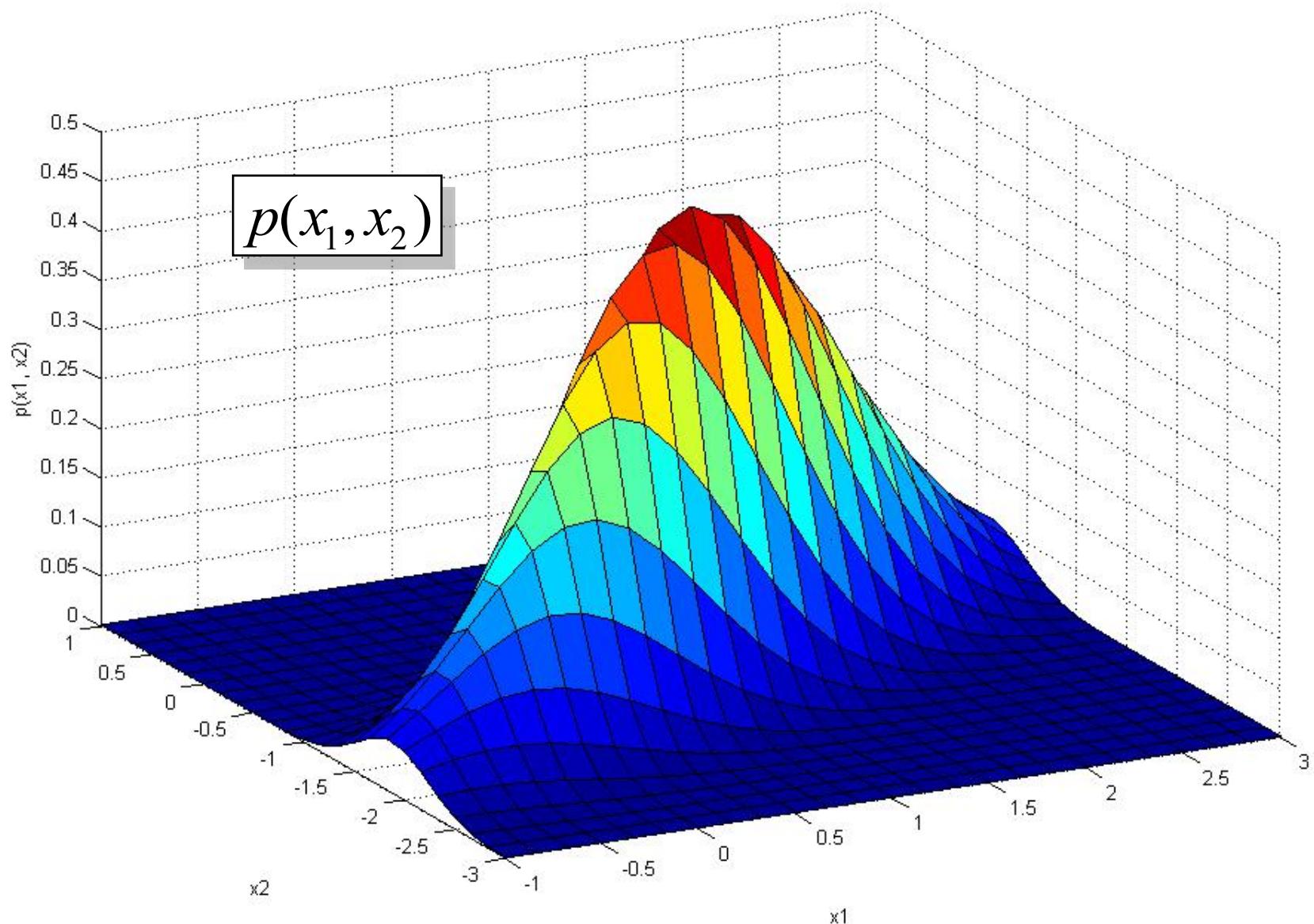
1.13 Multivariate Normal Densities

Simplified example for usage of meshgrid

```
[X1, X2] = meshgrid(linspace(-1,3,5)', linspace(-3,1,5)');  
X = [X1(:) X2(:)];
```

X =

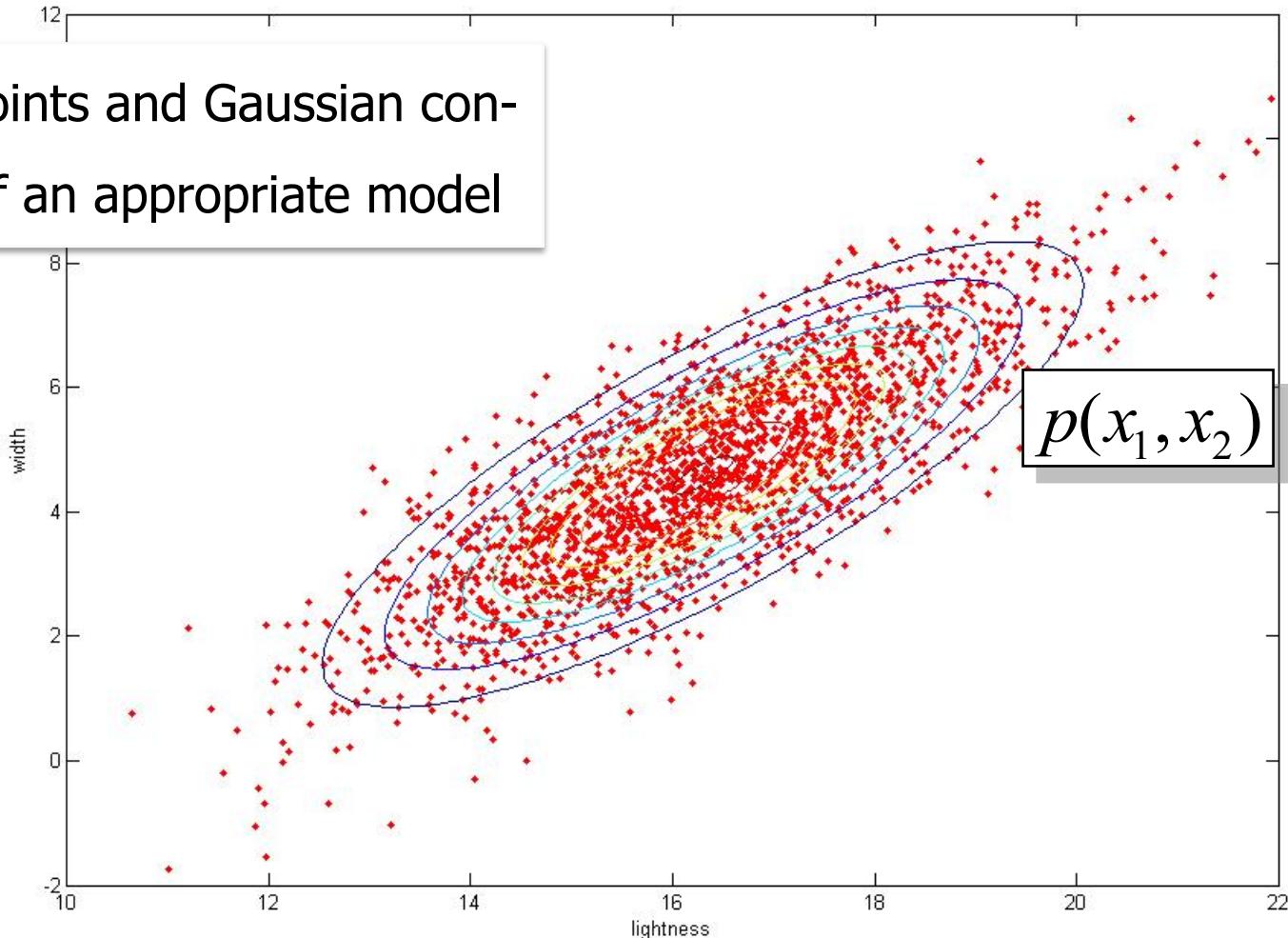
```
-1      -3  
-1      -2  
-1      -1  
-1      0  
-1      1  
0      -3  
0      -2  
0      -1  
0      0  
...  
...
```



Multivariate Gaussian Distribution with dependent variables

1.13 Multivariate Normal Densities

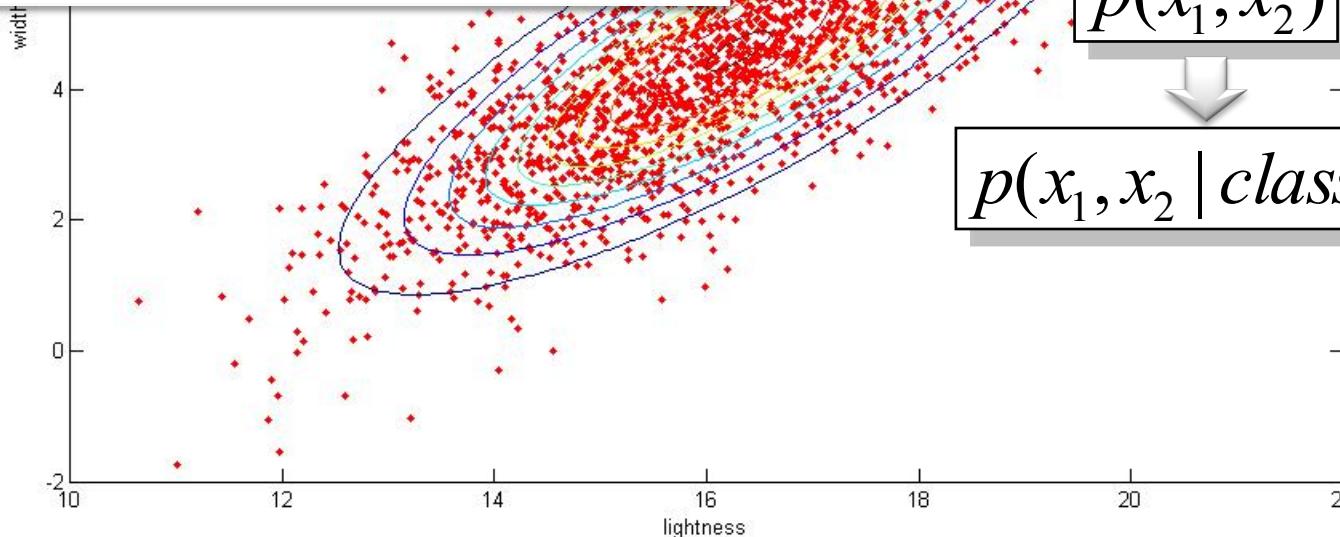
Data points and Gaussian contours of an appropriate model



1.13 Multivariate Normal Densities

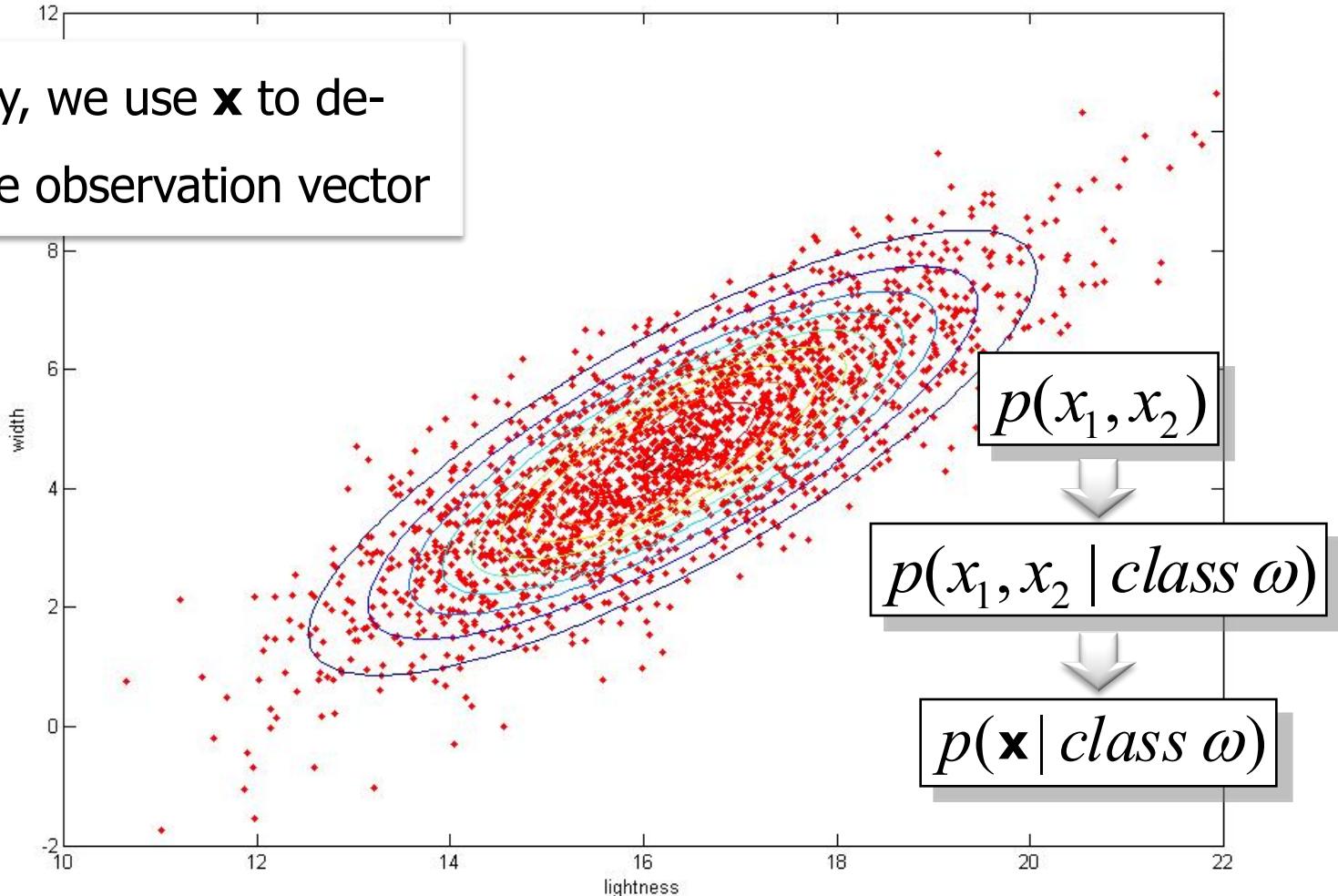
12

Important remarks: If these data points originate from the same class ω , the joint probability becomes a **class conditional probability**.



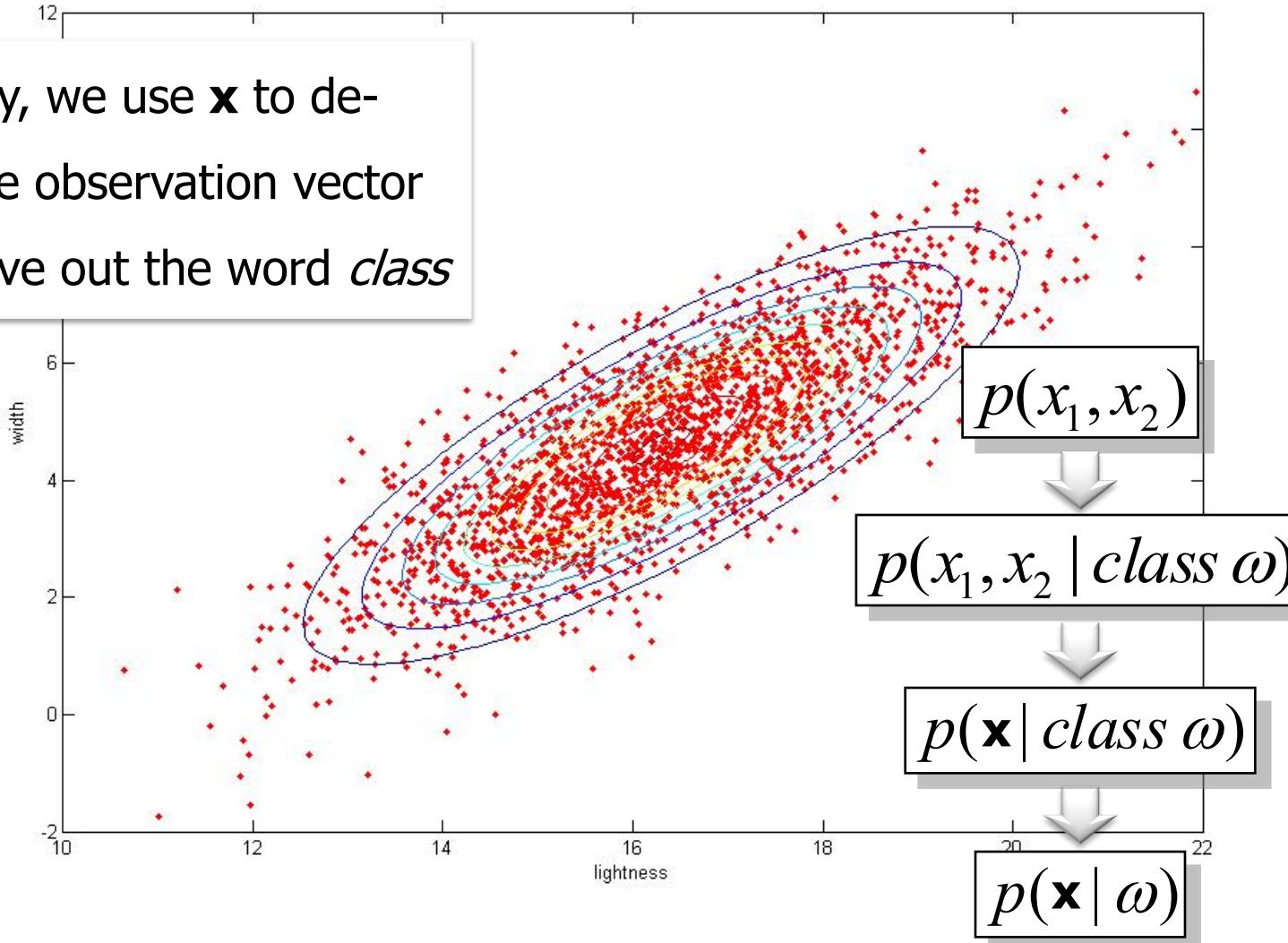
1.13 Multivariate Normal Densities

Typically, we use \mathbf{x} to denote the observation vector

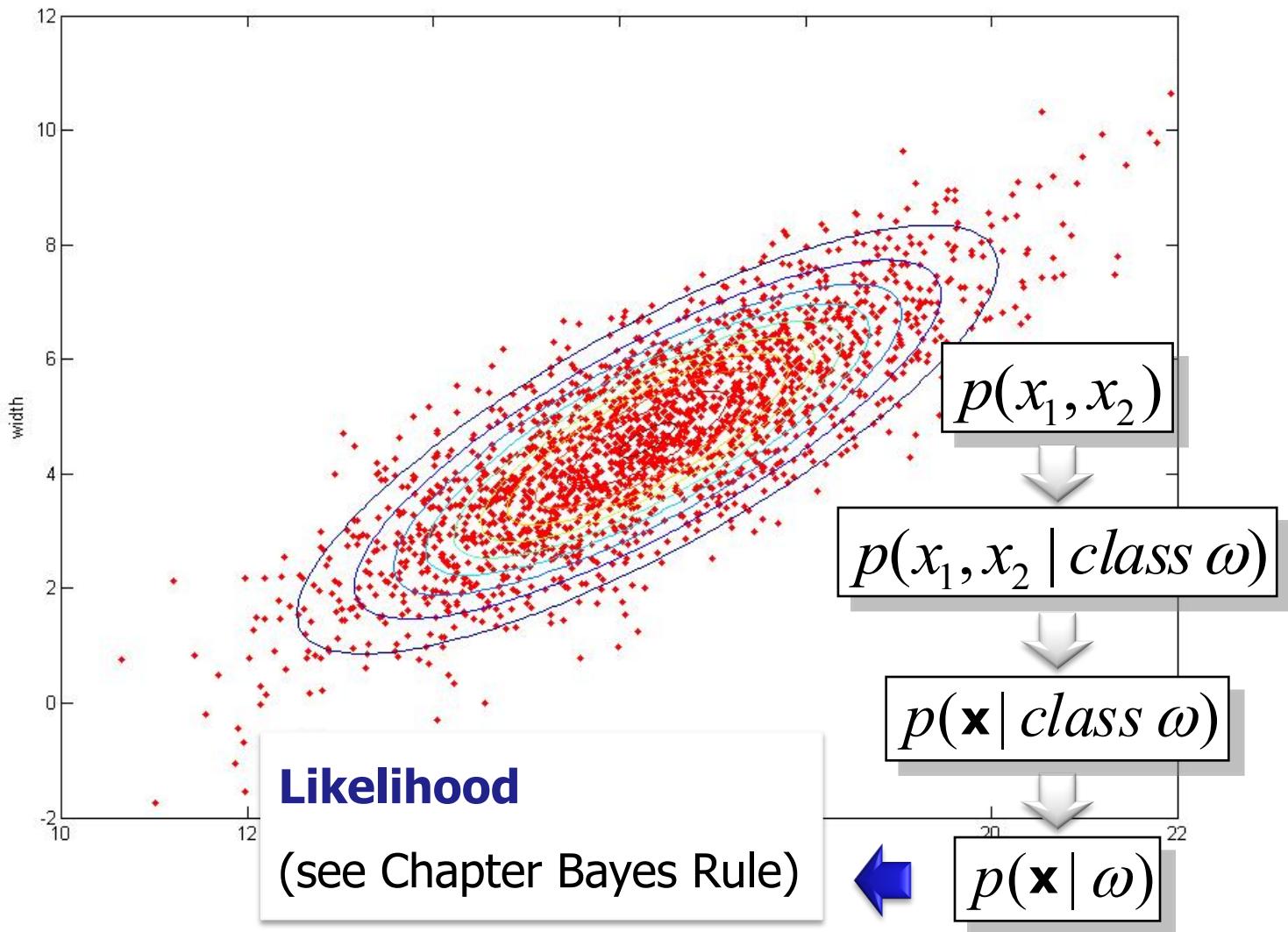


1.13 Multivariate Normal Densities

Typically, we use \mathbf{x} to denote the observation vector and leave out the word *class*



1.13 Multivariate Normal Densities



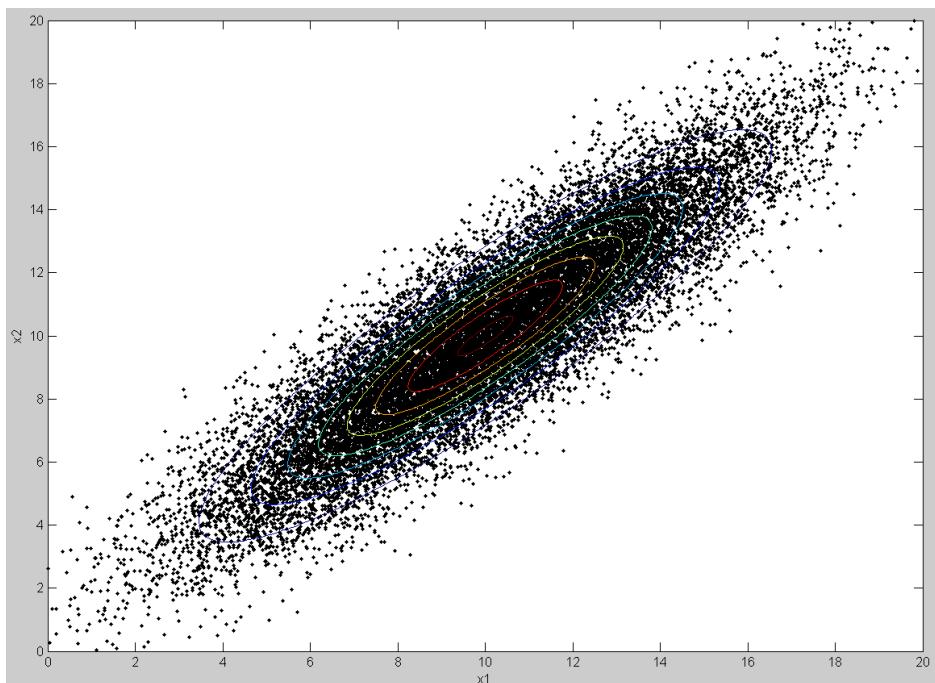
1.14 Mixture Densities

Explanation: next slides

Up to now:

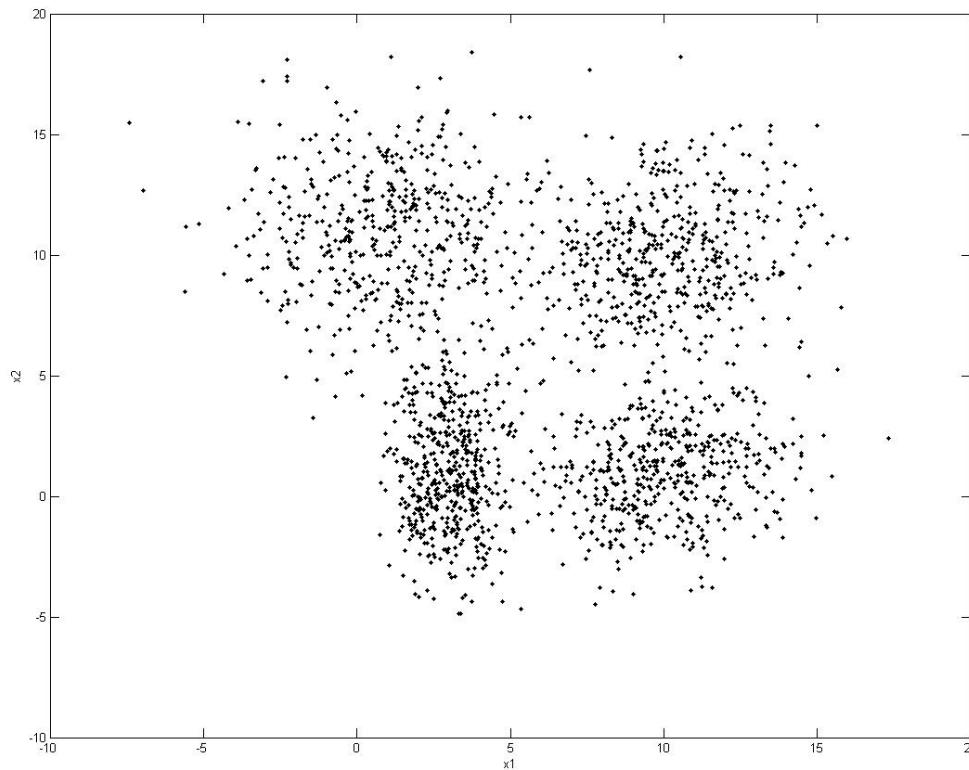
Arbitrary distributions which could be modeled by a **unimodal** but **multivariate** normal distribution

depending on a vector
of random variables



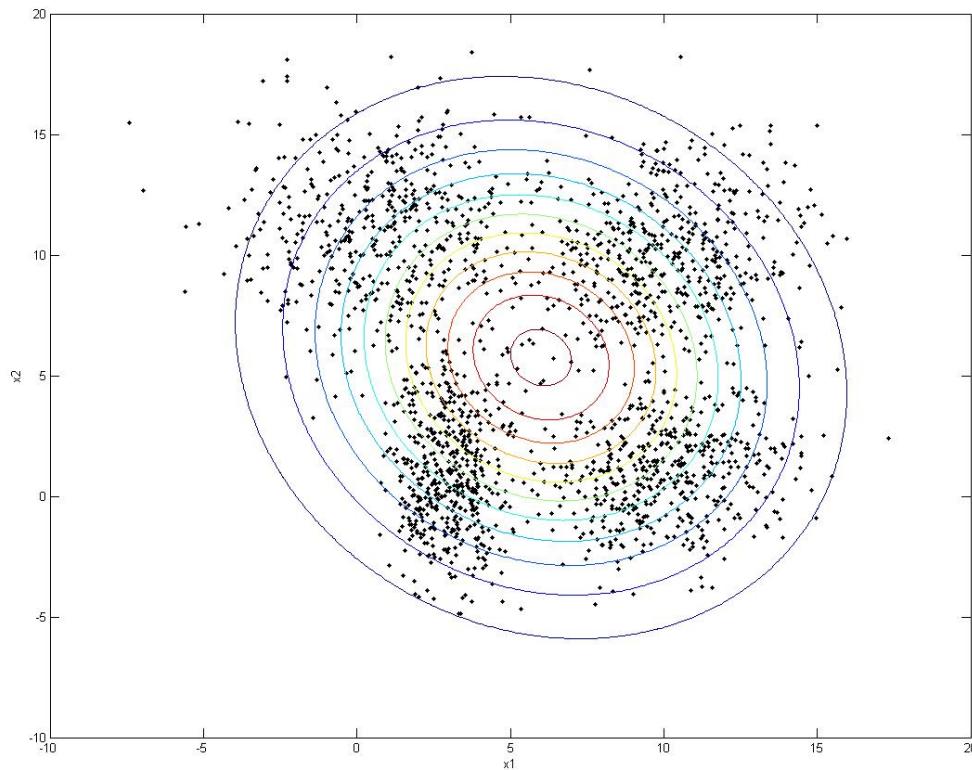
1.14 Mixture Densities

What can we do with the following type of distribution?



1.14 Mixture Densities

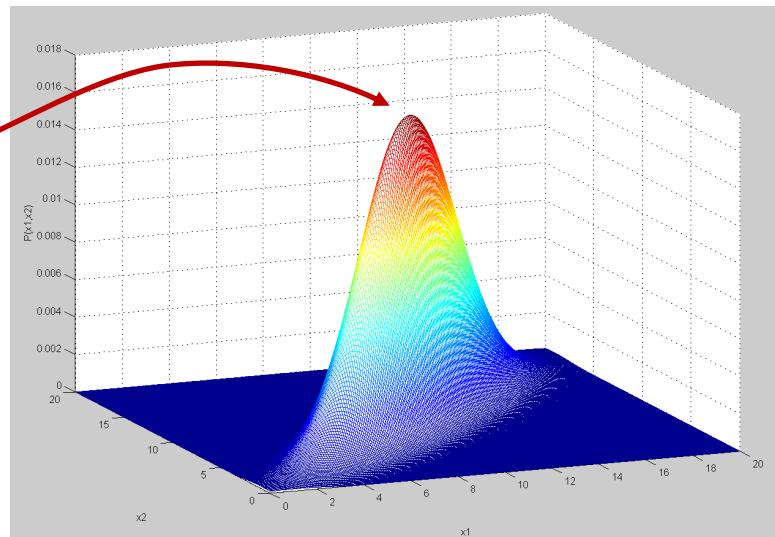
What can we do with the following type of distribution?



Definition

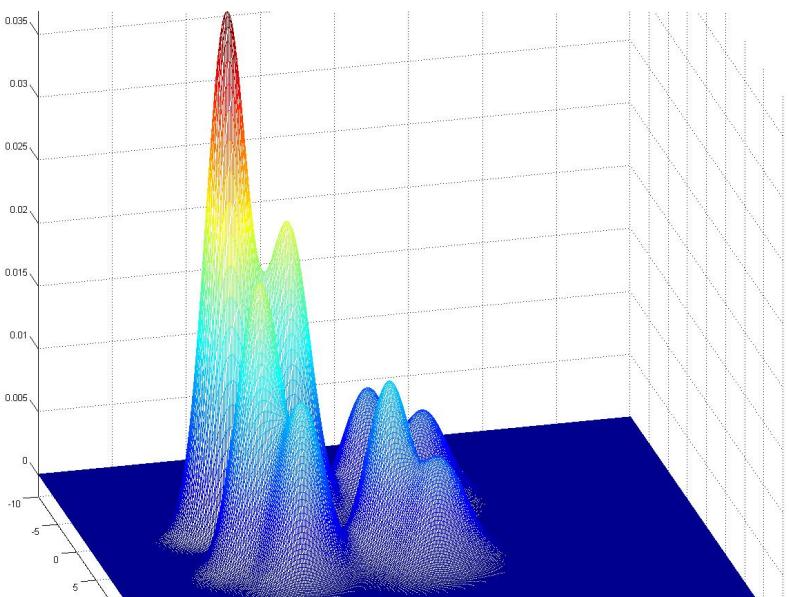
Unimodal distributions:

- exhibit only a **single** local maximum



Multimodal distributions
or **Mixture Densities**

- exhibit **several** local maxima



1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model: Number of modes

$$p(\mathbf{x}) = \sum_{i=1}^I p(\mathbf{x}, i)$$

Law of total probability

1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model:

$$p(\mathbf{x}) = \sum_{i=1}^I p(\mathbf{x}, i) = \sum_{i=1}^I p(i) \cdot p(\mathbf{x} | i)$$

Prior probability of mode i

Conditional probability of observation x given mode i

1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model:

$$p(\mathbf{x}) = \sum_{i=1}^I p(\mathbf{x}, i) = \sum_{i=1}^I p(i) \cdot p(\mathbf{x} | i) = \sum_{i=1}^I p(i) \cdot N(\mathbf{x} | \mu_i, \Sigma_i)$$

Modeling of conditional probability by unimodal normal distribution $N(\cdot)$.

1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model:

$$p(\mathbf{x}) = \sum_{i=1}^I p(c_i) \cdot N(\mathbf{x} | \mu_i, \Sigma_i)$$

Unimodal normal distribution

Interpretation:

A Mixture Density is a **combination of weighted unimodal normal distributions**.

1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model:

$$p(\mathbf{x}) = \sum_{i=1}^I p(c_i) \cdot N(\mathbf{x} | \mu_i, \Sigma_i)$$

Interpretation:

A Mixture Density is a **combination of weighted unimodal normal distributions**.

Remark: mean vector and covariance matrix can be individually adjusted for each of the I modes.

1.14 Mixture Densities

Multimodal distributions can be generated by a combination of distributions.

Standard model:

$$p(\mathbf{x}) = \sum_{i=1}^I p(c_i) \cdot N(\mathbf{x} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$$

Class-conditional probability (likelihood):

$$p(\mathbf{x} | k) = \sum_{i=1}^I p(\mathbf{x}, i | k) = \sum_{i=1}^I p(i | k) \cdot p(\mathbf{x} | i, k)$$

$$p(\mathbf{x} | k) = \sum_{i=1}^I p(i | k) \cdot N(\mathbf{x} | \boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik})$$

1.14 Mixture Densities

Class-conditional distribution:

$$p(\mathbf{x} | k) = \sum_{i=1}^I p(i | k) \cdot N(\mathbf{x} | \mu_{ik}, \Sigma_{ik})$$

class dependency

1.14 Mixture Densities

In praxis usually sufficient: **Maximum Approximation**

$$p(\mathbf{x} | k) = \sum_{i=1}^I p(i | k) \cdot N(\mathbf{x} | \mu_{ik}, \Sigma_{ik})$$

$$p(\mathbf{x} | k) \approx \max_i \{ p(i | k) \cdot N(\mathbf{x} | \mu_{ik}, \Sigma_{ik}) \}$$

Remarks:

- Sum is replaced by maximum operation
 - distribution with largest weighted probab. represents conditional prob.
- Motivation: exponential drop of normal distribution
 - sum is dominated by largest summand

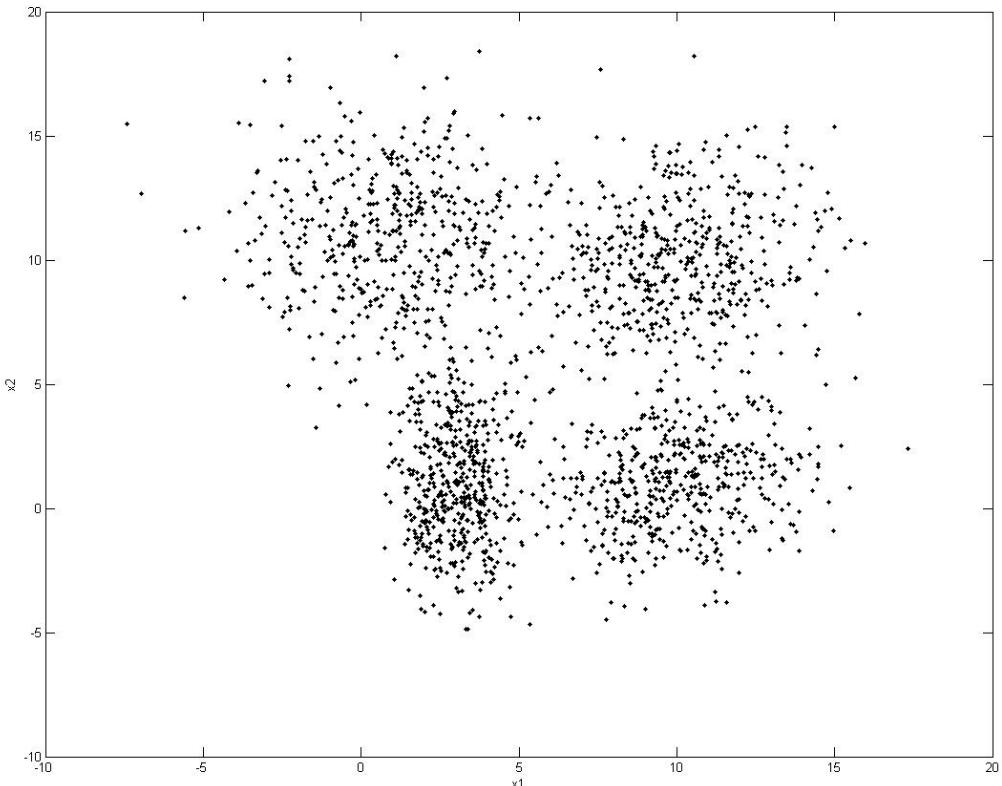
1.14 Mixture Densities

Question: How to estimate

- number of modes I ,
- priors $p(i | k)$,
- mean vectors,
- covariance matrices?

$$p(\mathbf{x} | k) = \sum_{i=1}^I p(i | k) \cdot N(\mathbf{x} | \boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik})$$

$$p(\mathbf{x} | k) \approx \max_i \{ p(i | k) \cdot N(\mathbf{x} | \boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik}) \}$$



K-Means Clustering

Goal:

Split dataset into I partitions to minimize square error criterion

$$G = \sum_{i=1}^I \sum_{x \in S_i} \|x - \mu_i\|^2$$

Diagram illustrating the square error criterion:

- All vectors x associated to cluster S_i** : Circled in red, representing the data points assigned to cluster S_i .
- Squared distance**: Circled in red, representing the squared Euclidean distance between a data vector x and the cluster center μ_i .
- All clusters S_i** : Circled in red, representing the sum of squared distances for all data points in cluster S_i .

K-Means Clustering

Iterative Technique:

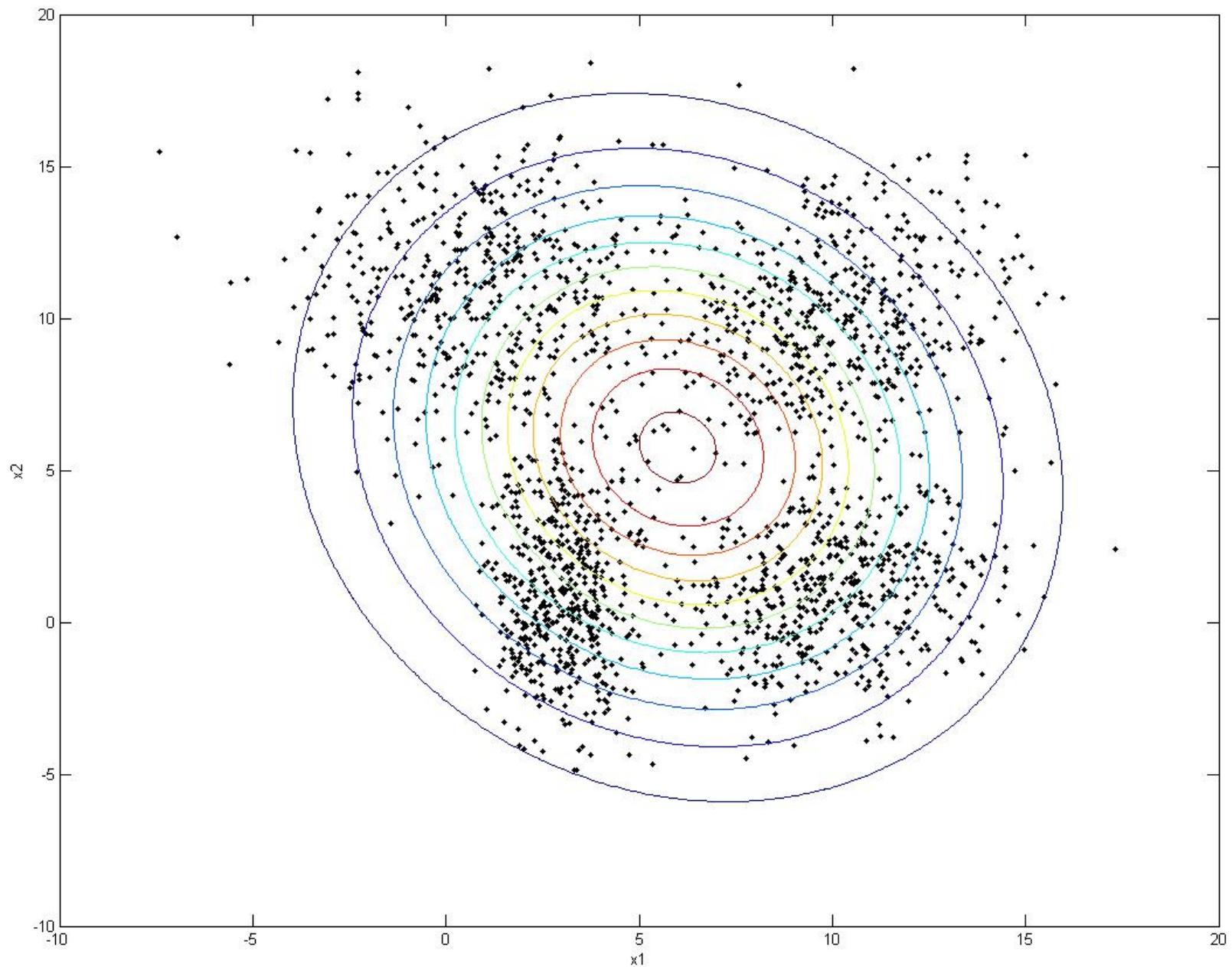
Given: Training observations $x_1, \dots, x_j, \dots, x_n$

(with

$$x_j = \begin{bmatrix} x_{j1} \\ \dots \\ x_{jD} \end{bmatrix}$$

1. Determine parameters of unimodal multivariate distribution

Compute empirical mean vector and covariance matrix from all training observations



K-Means Clustering

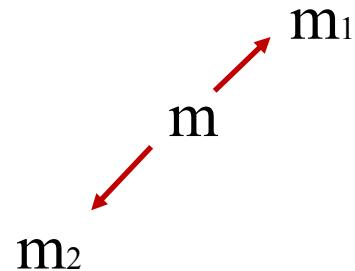
2. Split distribution

a. Split mean

$$\text{eps} = 0.001;$$

$$m1 = m + \text{eps};$$

$$m2 = m - \text{eps};$$



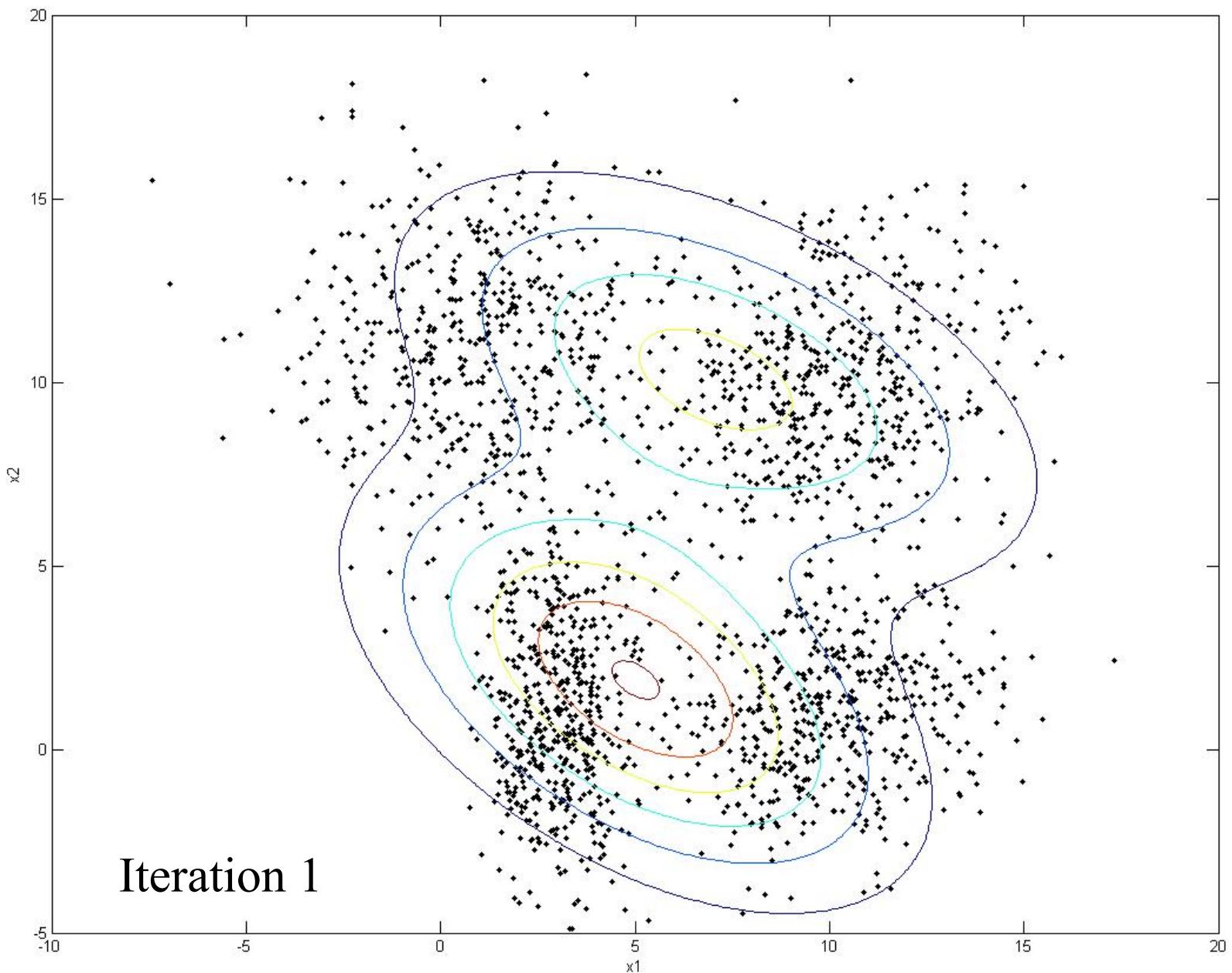
b. Assign each training sample x_j to optimal cluster optimal \rightarrow minimum increase of cluster variance

$$S_k^{(t)} = \left\{ x_j : \|x_j - \mu_k^{(t)}\|^2 \leq \|x_j - \mu_i^{(t)}\|^2 \text{ for all clusters } i = 1, \dots, I, i \neq k \right\}$$

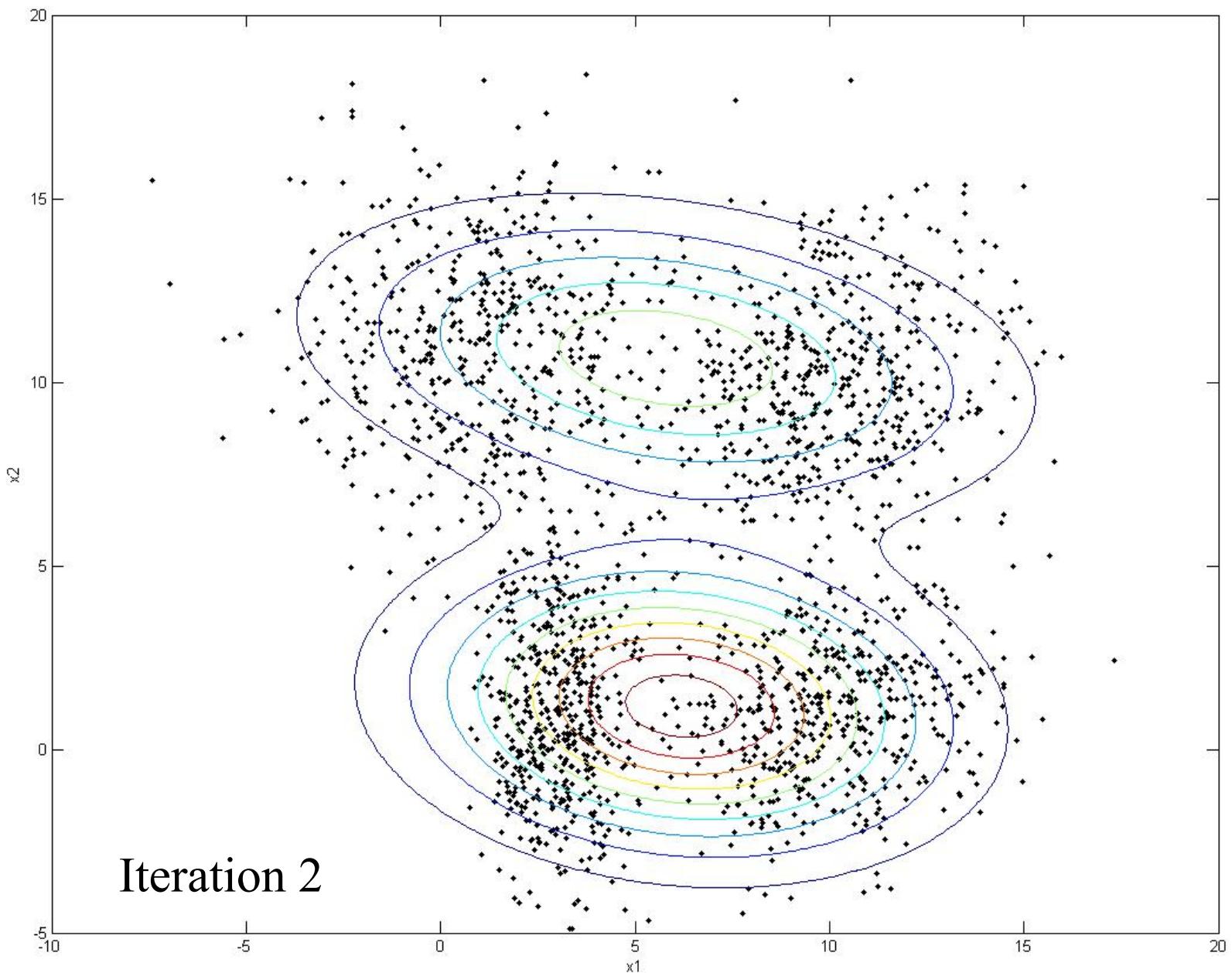
c. Determine the new means for all clusters with new associated x_j

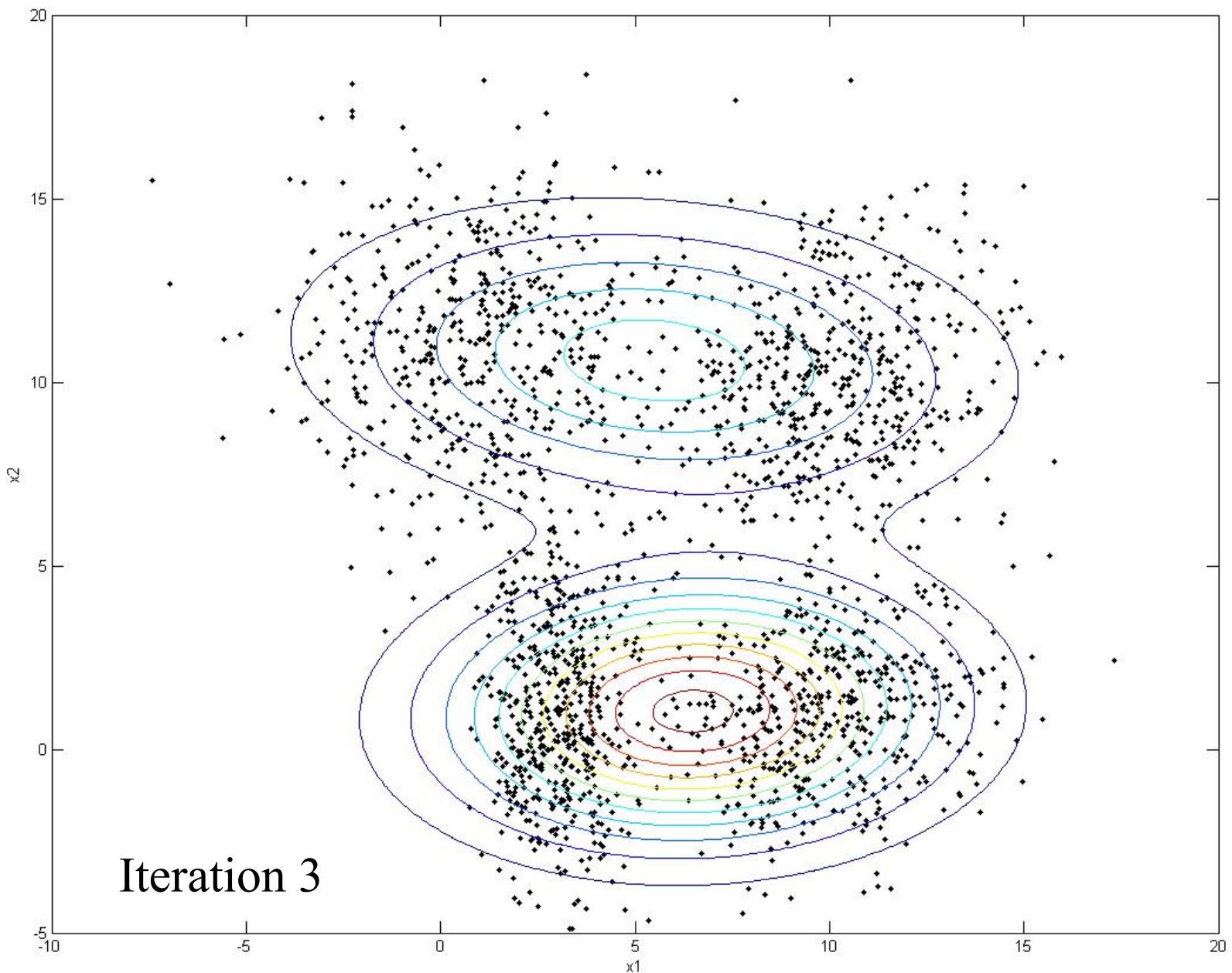
$$\mu_k^{(t+1)} = \frac{1}{|S_k^{(t)}|} \sum_{x_j \in S_k^{(t)}} x_j$$

iterate

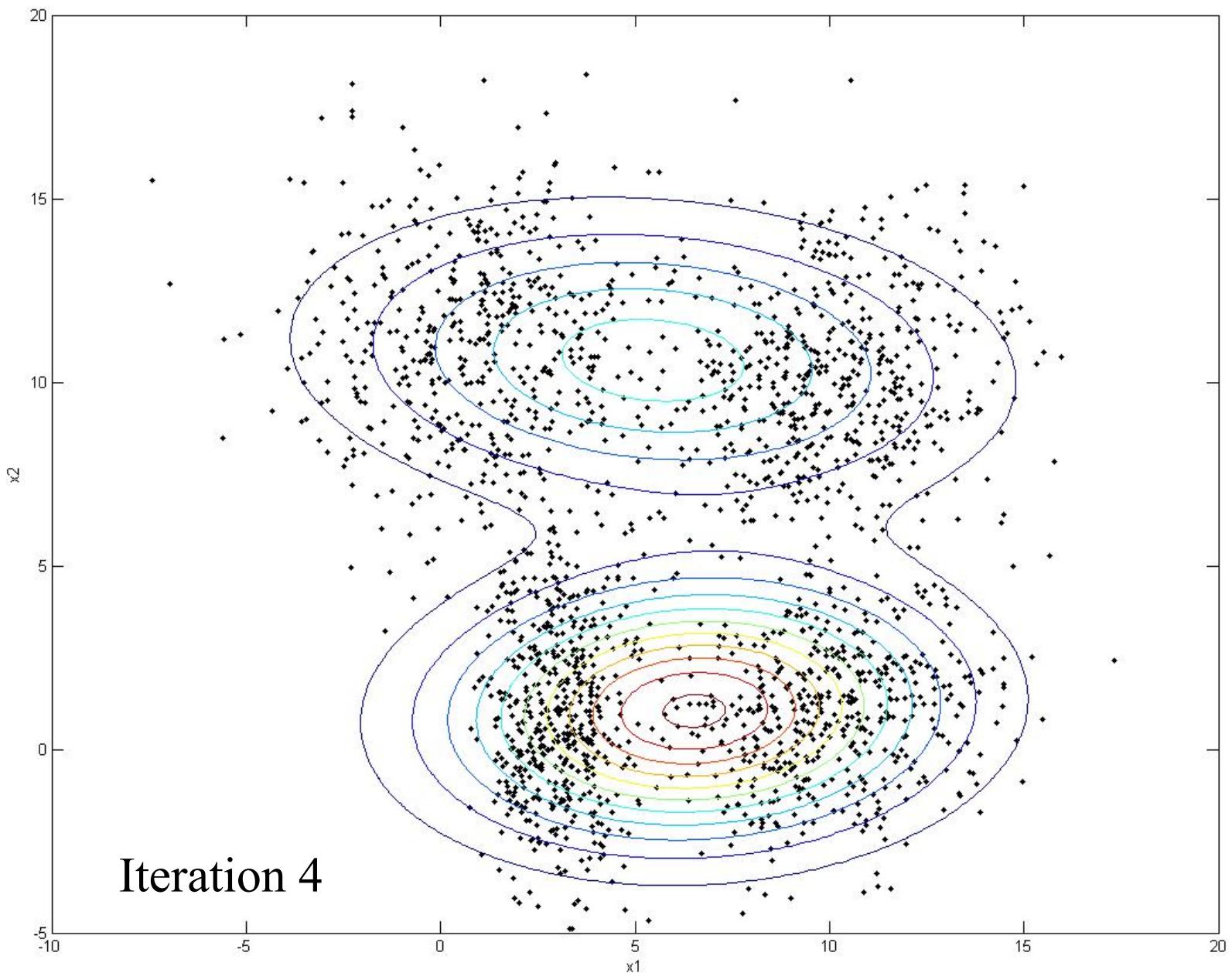


Iteration 1





Iteration 3



Iteration 4

K-Means Clustering

2. Split distribution

a. Split means

$$\text{eps} = 0.001;$$

$$m11 = m1 + \text{eps}; \quad m12 = m1 - \text{eps};$$

$$m21 = m2 + \text{eps}; \quad m22 = m2 - \text{eps};$$

b. Assign each training sample x_j to optimal cluster

optimal \rightarrow minimum increase of cluster variance

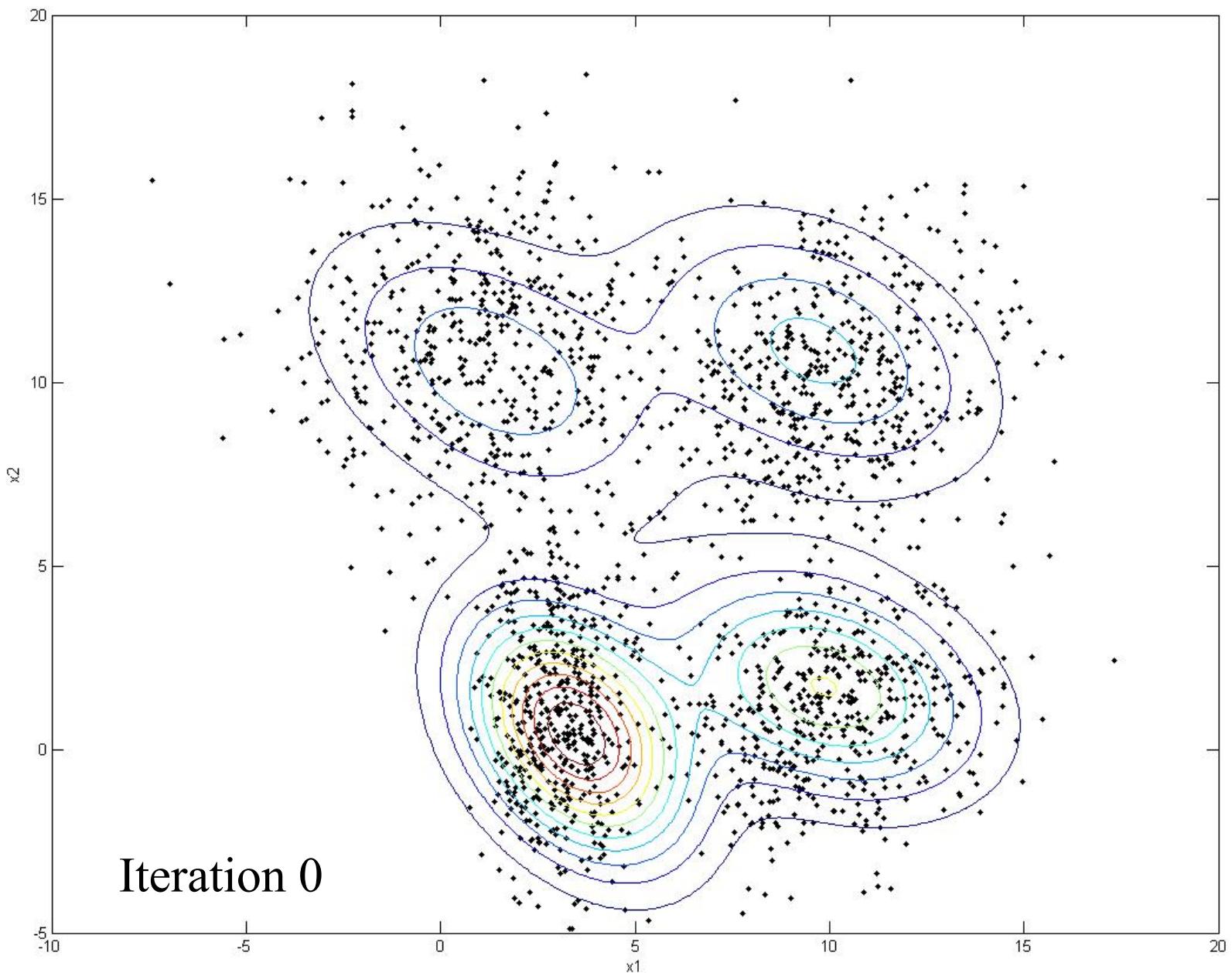
iterate

$$S_k^{(t)} = \left\{ x_j : \|x_j - \mu_k^{(t)}\|^2 \leq \|x_j - \mu_i^{(t)}\|^2 \text{ for all clusters } i = 1, \dots, I \right\}$$

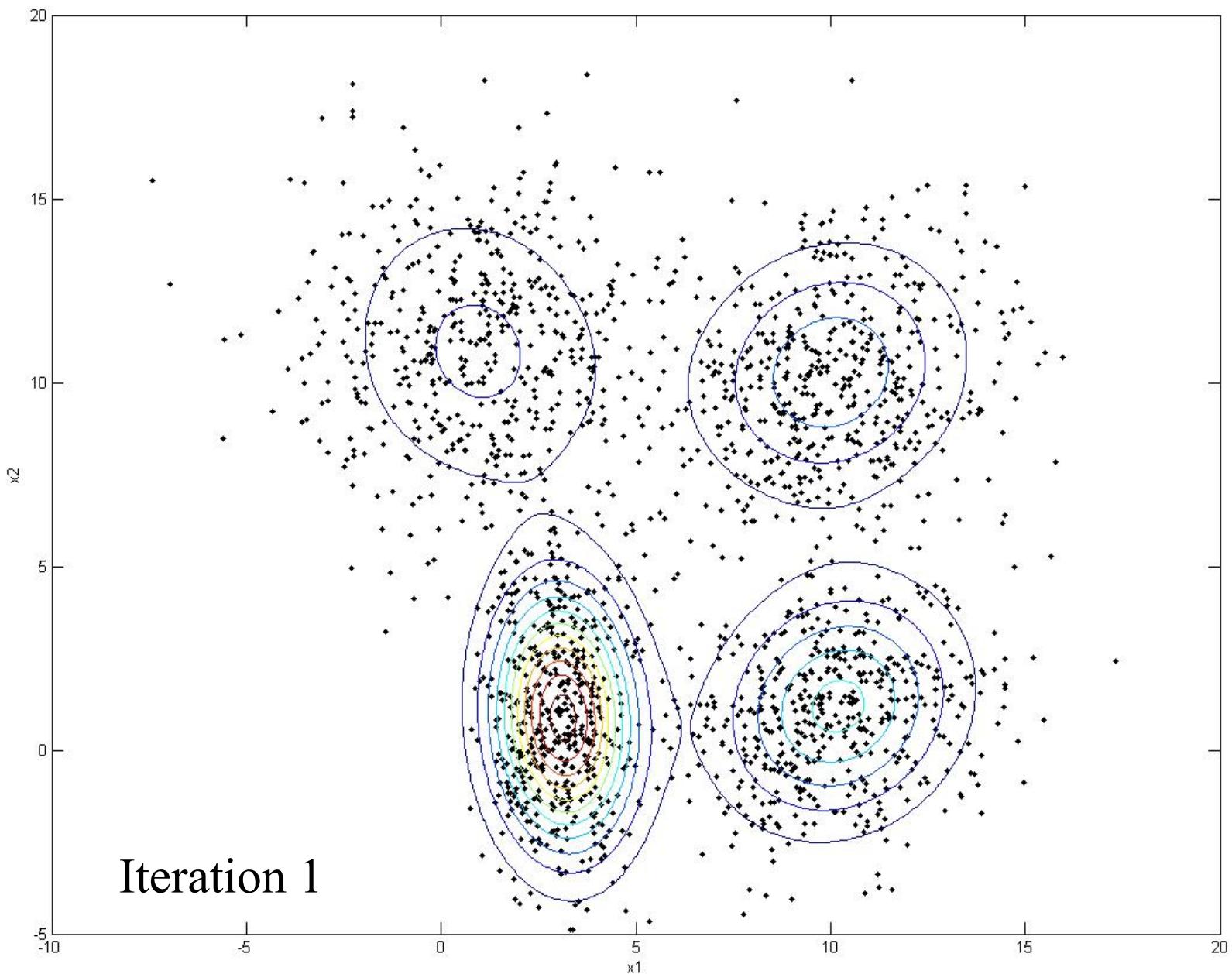
iterate

c. Determine the new means for all clusters with new associated x_j

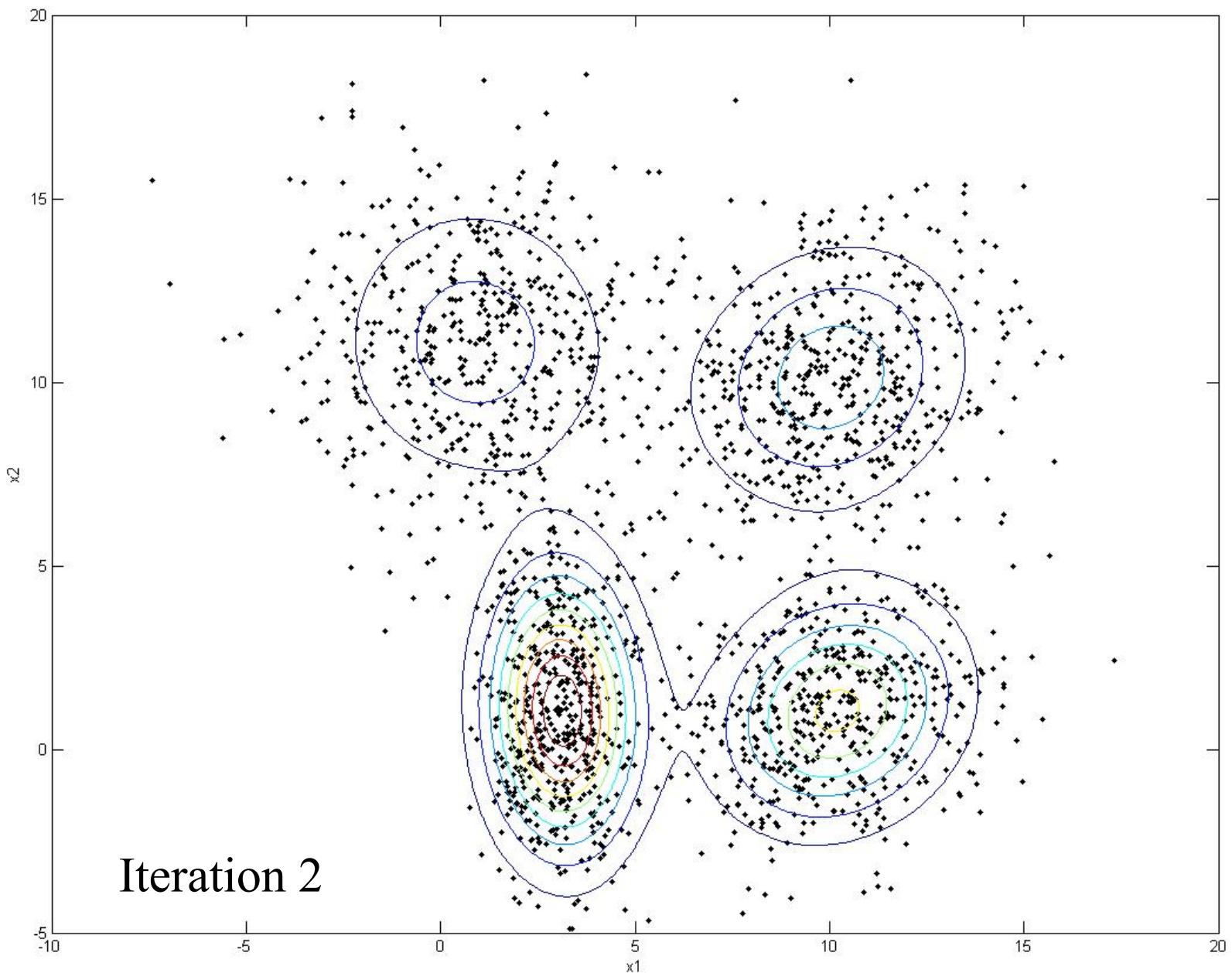
$$\mu_k^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$



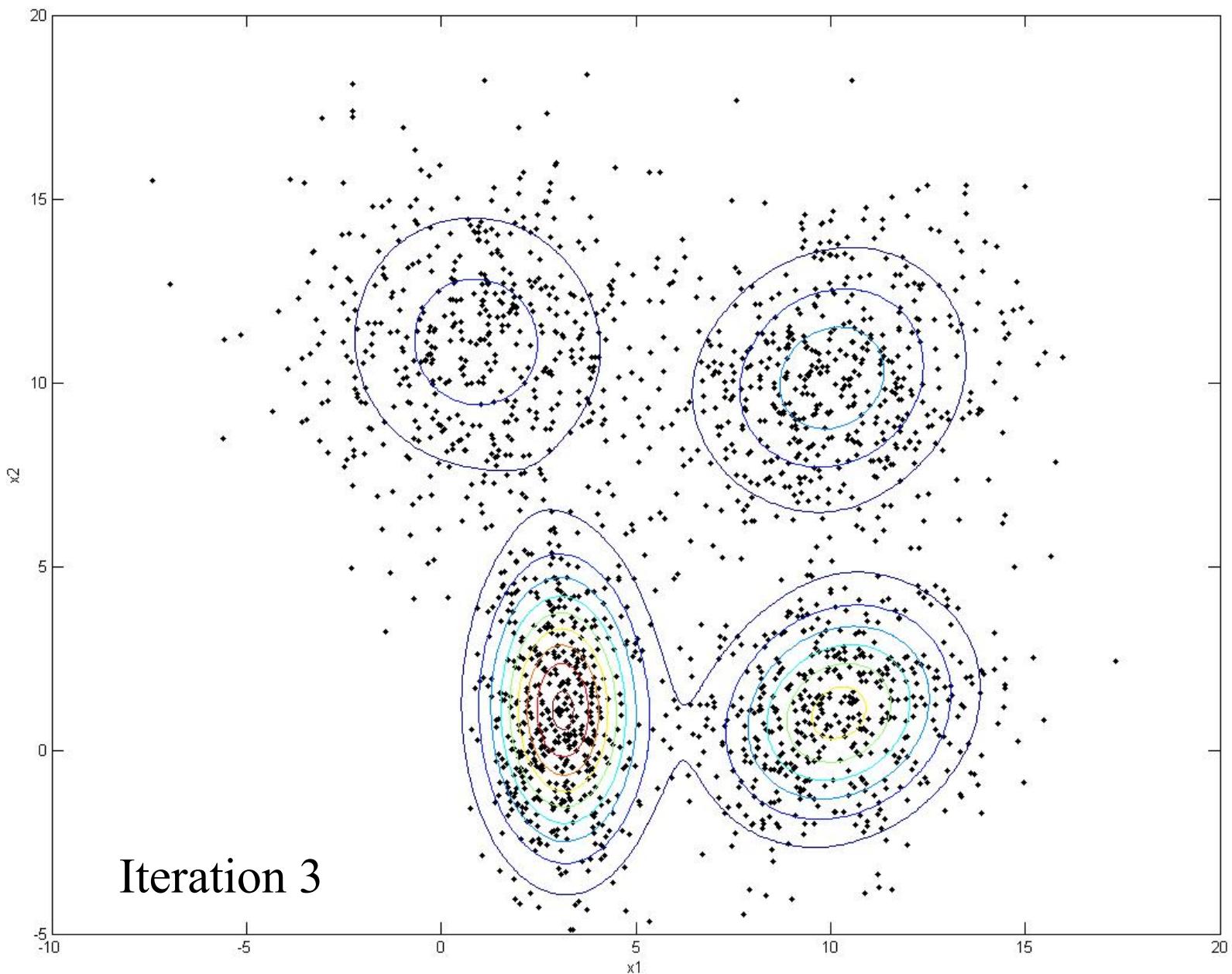
Iteration 0



Iteration 1



Iteration 2



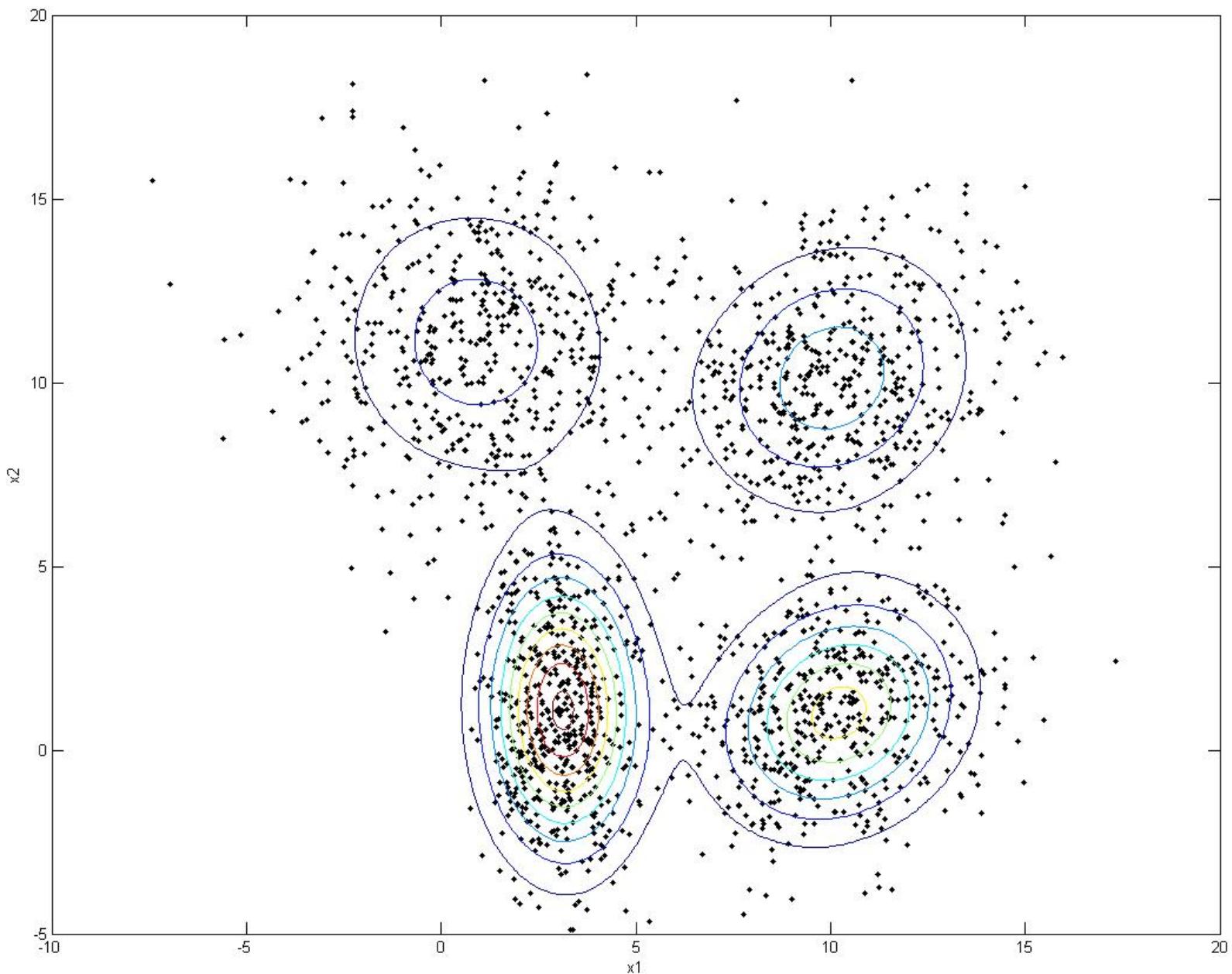
Iteration 3

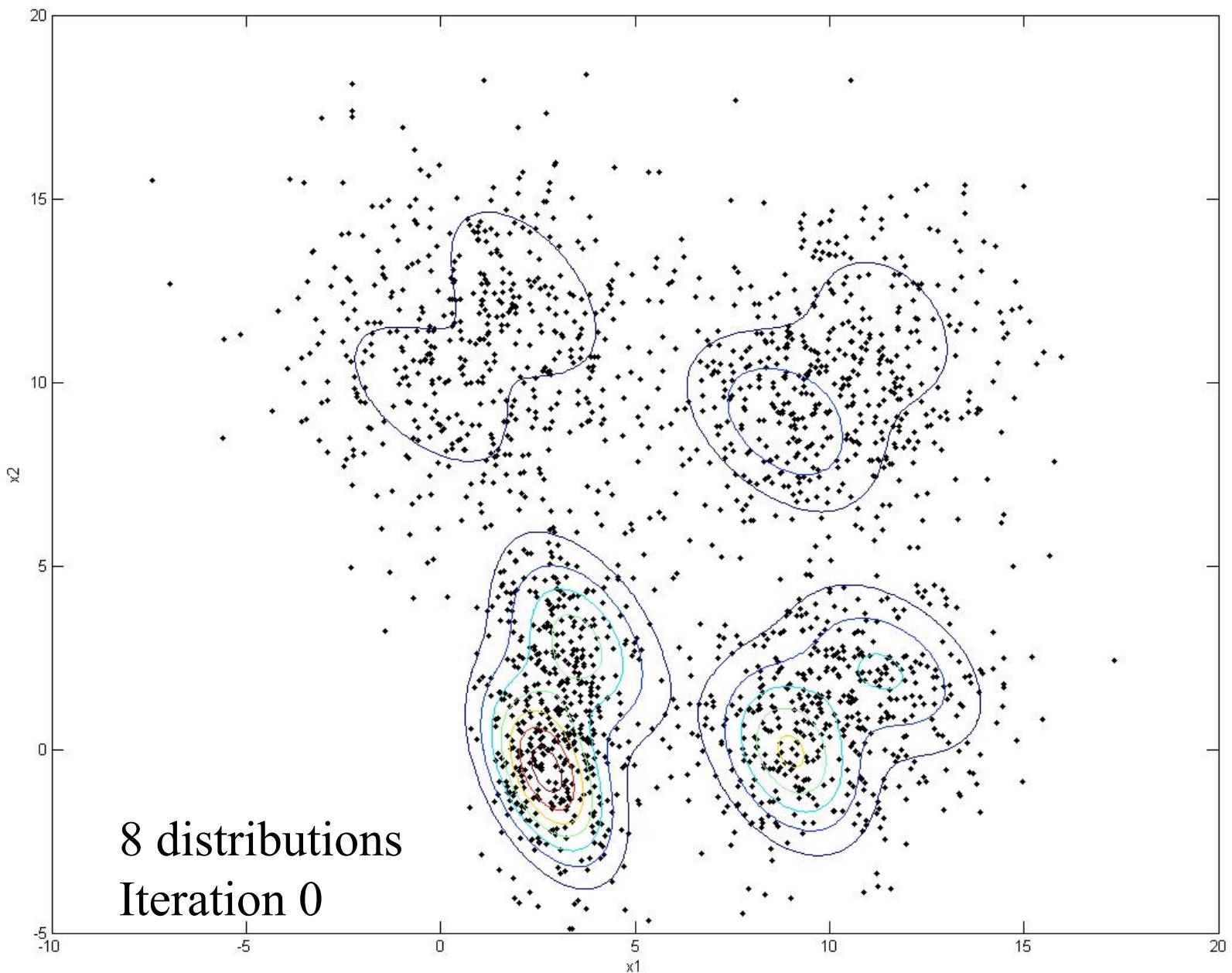
K-Means Clustering

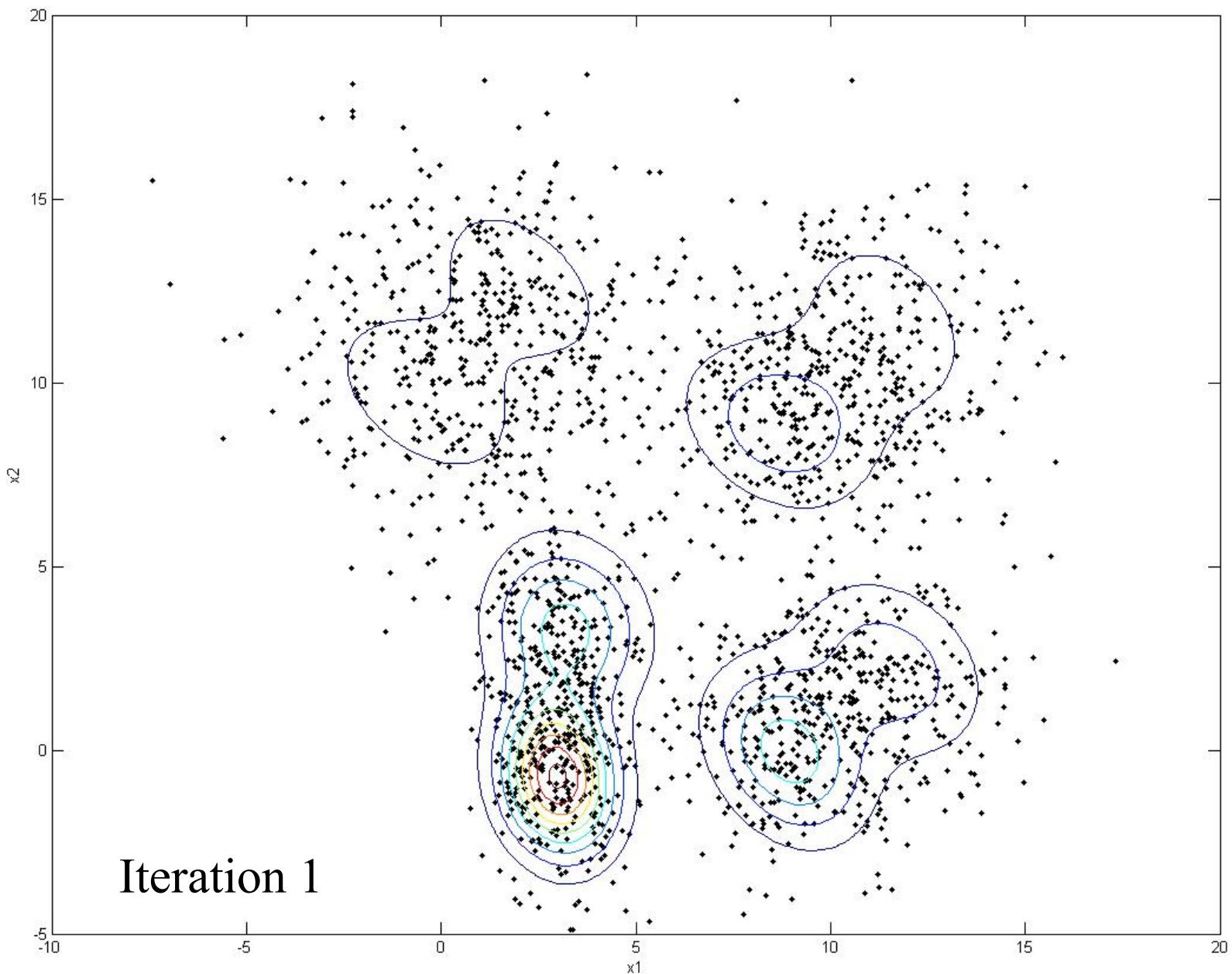
Repeat of split and reestimation until minimum square error achieved.

$$G = \sum_{i=1}^I \sum_{x \in S_i} \|x - \mu_i\|^2$$

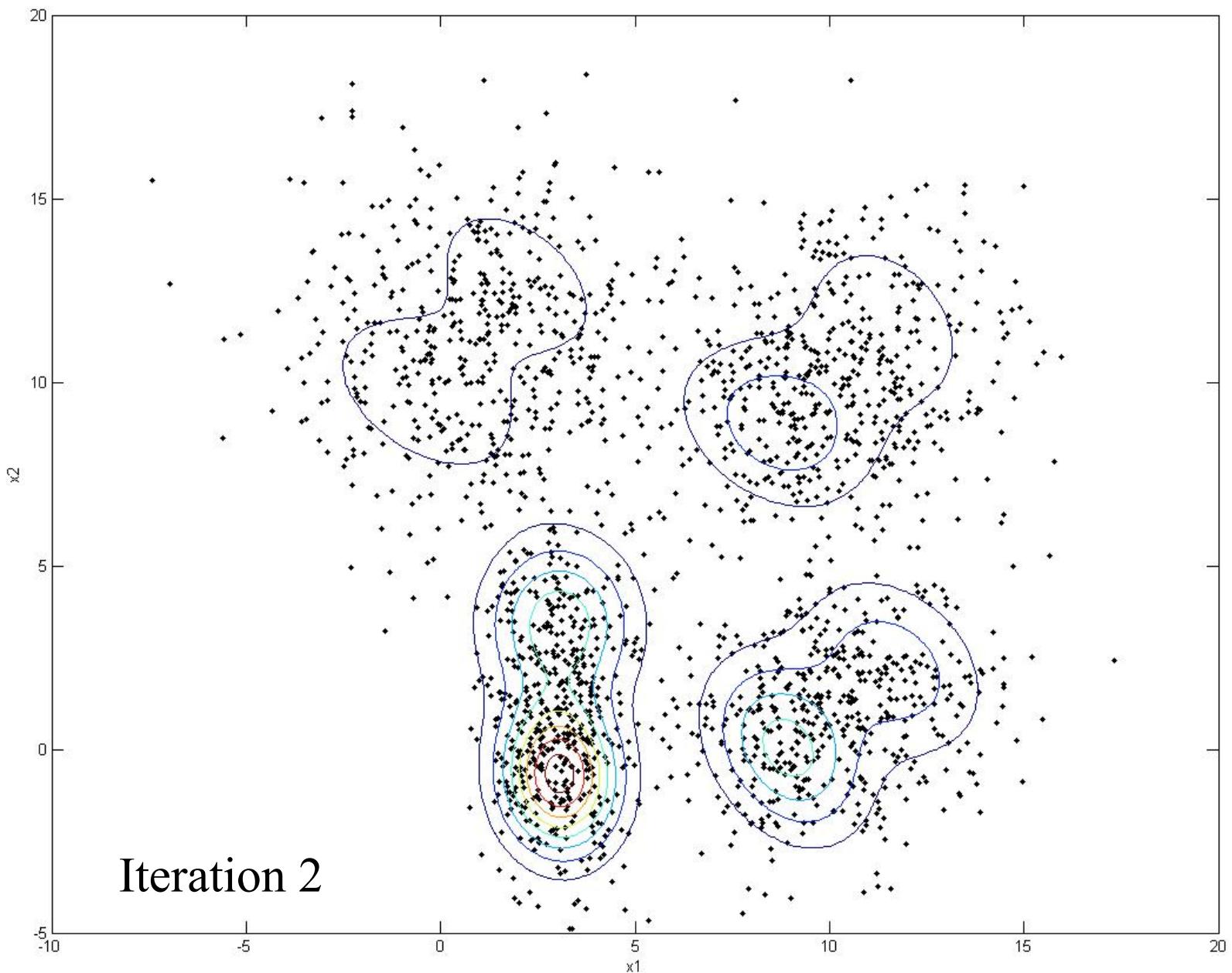
- Determination of optimal number of splits/iterations requires experience
- Possible other criteria:
 - Classification error rate on test data
 - Likelihood increase on training data below threshold
- Possible problem: **Overfitting**
 - Too strong adaptation to training data
 - No generalization capability of the model to unknown data.



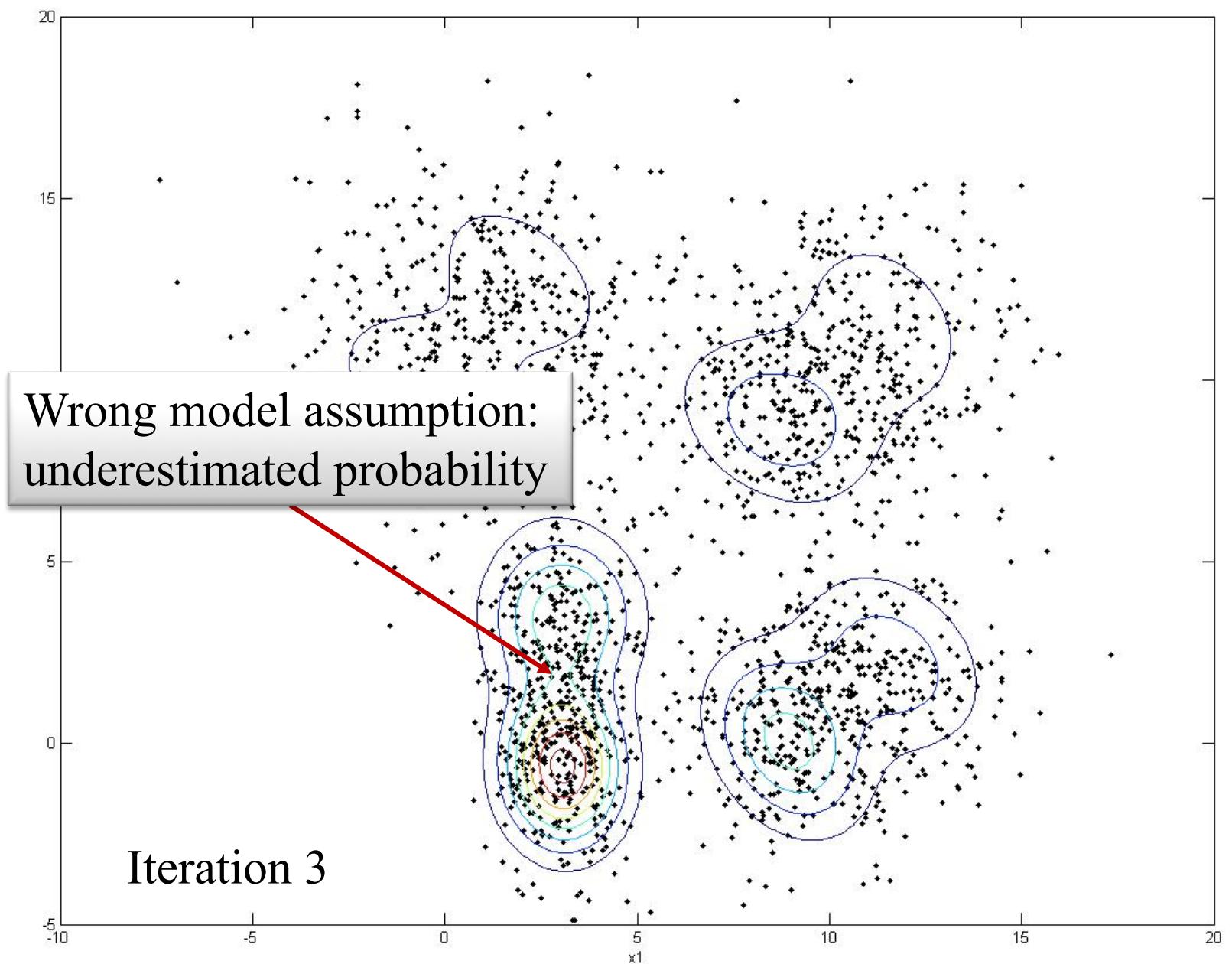




Iteration 1



Iteration 2



1.15 Bayes rule

We already know: $P(x | y) = \frac{P(x, y)}{P(y)}$ and $P(y | x) = \frac{P(x, y)}{P(x)}$

1.15 Bayes rule

We already know: $P(x | y) = \frac{P(x, y)}{P(y)}$ and $P(y | x) = \frac{P(x, y)}{P(x)}$

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

insert

1.15 Bayes rule

We already know: $P(x | y) = \frac{P(x, y)}{P(y)}$ and $P(y | x) = \frac{P(x, y)}{P(x)}$

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

insert

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

Bayes Rule

1.15 Bayes rule

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

With

$$P(y) = \sum_{x \in X} P(x, y) = \sum_{x \in X} P(y, x) = \sum_{x \in X} P(y | x) \cdot P(x)$$

conditional probability (see above)

Law of Total Probability
(Marginal Distribution)

1.15 Bayes rule

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

With

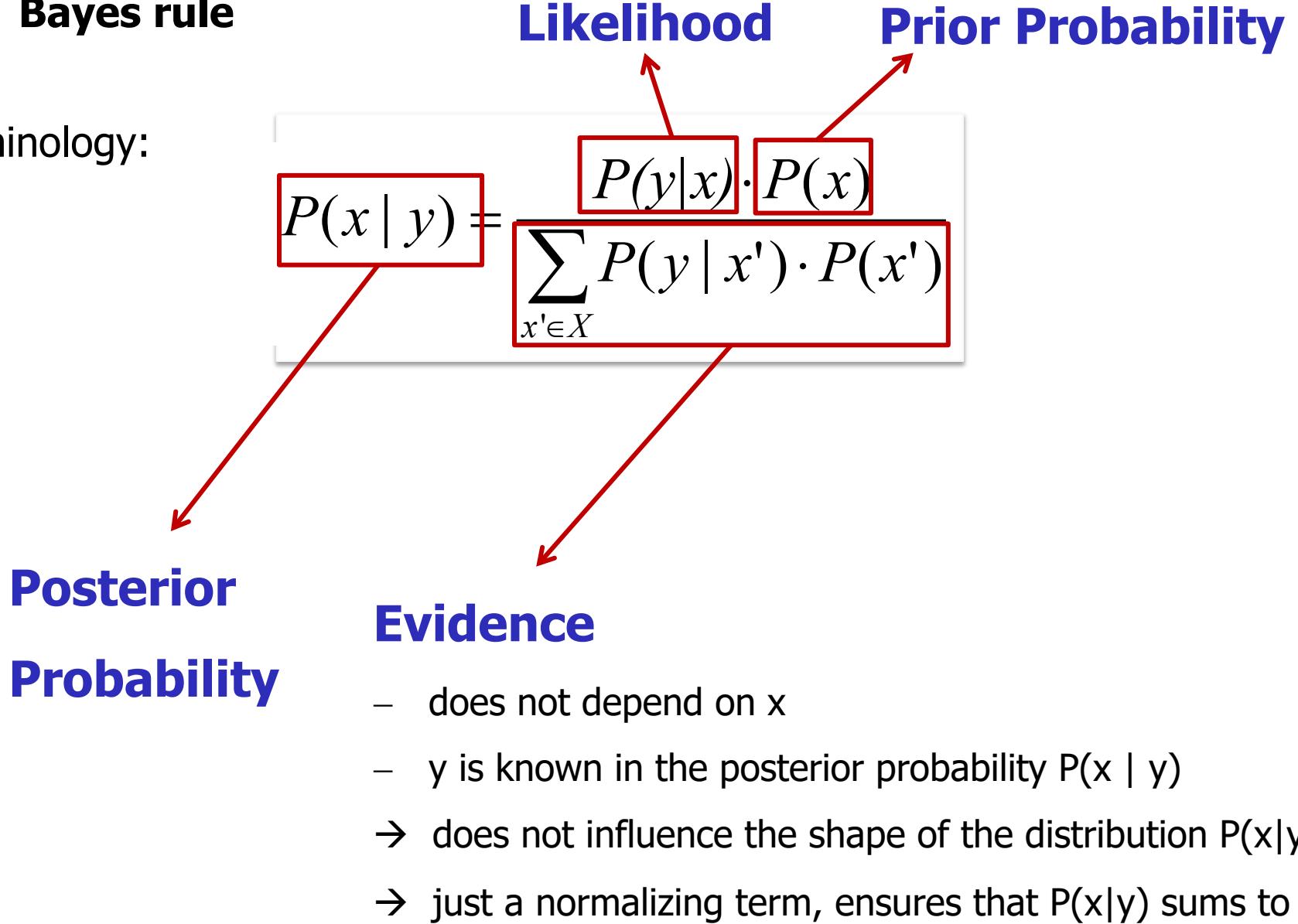
$$P(y) = \sum_{x \in X} P(x, y) = \sum_{x \in X} P(y, x) = \sum_{x \in X} P(y | x) \cdot P(x)$$

we get another important representation of Bayes rule:

$$P(x | y) = \frac{P(y|x) \cdot P(x)}{P(y)} = \frac{P(y|x) \cdot P(x)}{\sum_{x' \in X} P(y | x') \cdot P(x')}$$

1.15 Bayes rule

Terminology:



1.15 Bayes rule

Relevance and interpretation:

- "Inverts" statistical connections
- If $P(y|x)$ is known but $P(x|y)$ is required

1.15 Bayes rule

Example: Spam-Filter

Given: 10.000 different mails, 5.000 spam, 5.000 ham (not spam)

It is easy to learn the probability of words, given that we know it is spam, that is:

$P(\text{Words} \mid \text{Spam})$ → We know the **class** (spam) and estimate the probability of **words**.

However: A spam filter must decide on Spam / Not-Spam based on **a given mail**

- We know the **words** (from mail) and want to get the probability of the **class**
- We need a model for: $P(\text{Spam} \mid \text{Words})$
- How can we invert the conditional probability?

1.15 Bayes rule

$$P(x | y) = \frac{P(y|x) \cdot P(x)}{P(y)} = \frac{P(y|x) \cdot P(x)}{\sum_{x \in X} P(y | x) \cdot P(x)}$$

Example: Spam-Filter

Answer: with Bayes Rule

$$P(\text{Spam} | \text{Words}) = \frac{P(\text{Words} | \text{Spam}) \cdot P(\text{Spam})}{\sum_{\substack{\text{all classes} \\ C \in [\text{Spam, Ham}]}} P(\text{Words} | C) \cdot P(C)}$$

1.15 Bayes rule

Example: Spam-Filter

Answer: with Bayes Rule

$$P(x | y) = \frac{P(y|x) \cdot P(x)}{P(y)} = \frac{P(y|x) \cdot P(x)}{\sum_{x \in X} P(y | x) \cdot P(x)}$$

$$P(\text{Spam} | \text{Words}) = \frac{P(\text{Words} | \text{Spam}) \cdot P(\text{Spam})}{\sum_{\substack{\text{all classes} \\ C \in [\text{Spam, Ham}]}} P(\text{Words} | C) \cdot P(C)}$$

1.15 Bayes rule

Example: Spam-Filter

Answer: with Bayes Rule

$$P(\text{Spam} \mid \text{Words}) = \frac{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam})}{\sum_{\substack{\text{all classes} \\ C \in [\text{Spam, Ham}]}} P(\text{Words} \mid C) \cdot P(C)}$$

$$= \frac{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam})}{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam}) + P(\text{Words} \mid \text{Ham}) \cdot P(\text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

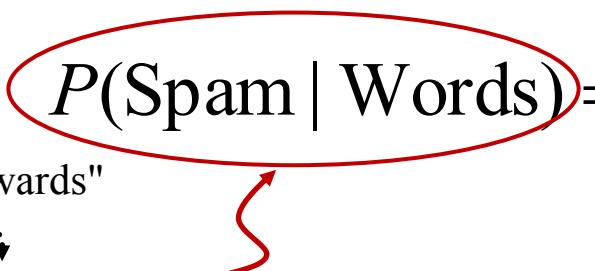
Answer: with Bayes Rule

$$P(\text{Spam} \mid \text{Words}) = \frac{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam})}{\sum_{\text{all classes}} P(\text{Words} \mid C) \cdot P(C)}$$

C ∈ [Spam, Ham]

"afterwards"

Posterior Probability



Probability of class spam

after we observed the words.

1.15 Bayes rule

Example: Spam-Filter

Answer: with Bayes Rule

$$P(\text{Spam} | \text{Words}) = \frac{P(\text{Words} | \text{Spam}) \cdot P(\text{Spam})}{\sum_{\substack{\text{all classes} \\ C \in [\text{Spam, Ham}]}} P(\text{Words} | C) \cdot P(C)}$$

Likelihood or class dependent probability

Probability of words if we know that it is spam.

1.15 Bayes rule

Example: Spam-Filter

Answer: with Bayes Rule

$$P(\text{Spam} | \text{Words}) = \frac{P(\text{Words} | \text{Spam}) \cdot P(\text{Spam})}{\sum_{\text{all classes } C \in [\text{Spam, Ham}]} P(\text{Words} | C) \cdot P(C)}$$

Prior Spam Probability

Probability of class spam

before we made any observation
(i.e. without observation)

1.15 Bayes rule

Example: Spam-Filter

$$P(\text{Spam} \mid \text{Words})$$

$$= \frac{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam})}{P(\text{Words} \mid \text{Spam}) \cdot P(\text{Spam}) + P(\text{Words} \mid \text{Ham}) \cdot P(\text{Ham})}$$

How do we get the required probabilities?

→ from training data!

1.15 Bayes rule

Example: Spam-Filter

P(Spam), P(Ham) ?

- Realistic value: $P(\text{Spam}) = 0.8, P(\text{Ham}) = 0.2$
- Most spam filters, however, use $P(\text{Spam})=0.5, P(\text{Ham})=0.5$
(give each incoming mail a chance)

1.15 Bayes rule

Example: Spam-Filter

Simplification: Same prior probability of spam and ham

$$P(\text{Spam} \mid \text{Words}) = \frac{P(\text{Words} \mid \text{Spam})}{P(\text{Words} \mid \text{Spam}) + P(\text{Words} \mid \text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

P(Words | Spam), P(Words | Ham) ?

→ Probability of a word in spam mails can be estimated from training spams

$P(\text{Word} | \text{Spam}) = \text{Frequency of word in spam} / \text{Total number of spam-words}$

But this is the estimation of $P(\text{Word} | \text{Spam})$ instead of $P(\text{Words} | \text{Spam})$

Assumption: word probabilities are **statistically independent**

→ we can factorize as follows:

$$P(W_1, W_2, \dots, W_N | \text{Spam}) = P(W_1 | \text{Spam}) \cdot P(W_2 | \text{Spam}) \cdot \dots \cdot P(W_N | \text{Spam})$$

1.15 Bayes rule

Example: Spam-Filter

Details:

$$\begin{aligned} P(\text{Words} \mid \text{Spam}) &= P(\text{Word1}, \text{Word2}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2}, \text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= \dots \\ &= P(\text{Word1} \mid \text{Word2}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot \dots \cdot P(\text{WordN-1} \mid \text{WordN}, \text{Spam}) \cdot P(\text{WordN} \mid \text{Spam}) \end{aligned}$$

1.15 Bayes rule

Example: Spam-Filter

Details:

$$\begin{aligned} P(\text{Words} \mid \text{Spam}) &= P(\text{Word1}, \text{Word2}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2}, \text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= \dots \\ &= P(\text{Word1} \mid \text{Word2}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot \dots \cdot P(\text{WordN-1} \mid \text{WordN}, \text{Spam}) \cdot P(\text{WordN} \mid \text{Spam}) \end{aligned}$$

Naive assumption which often works well in practice:

Probability of one word is not influenced by other words in the mail.

1.15 Bayes rule

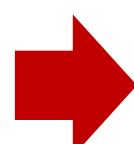
Example: Spam-Filter

Details:

$$\begin{aligned} P(\text{Words} \mid \text{Spam}) &= P(\text{Word1}, \text{Word2}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2}, \text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= P(\text{Word1} \mid \text{Word2}, \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word3}, \dots, \text{WordN} \mid \text{Spam}) \\ &= \dots \\ &= P(\text{Word1} \mid \text{Word2}, \dots, \text{WordN}, \text{Spam}) \cdot P(\text{Word2} \mid \text{Word3}, \dots, \text{WordN}, \text{Spam}) \cdot \dots \cdot P(\text{WordN-1} \mid \text{WordN}, \text{Spam}) \cdot P(\text{WordN} \mid \text{Spam}) \end{aligned}$$

Naive assumption which often works well in practice:

Probability of one word is not influenced by other words in the mail.


$$P(\text{Word1}, \text{Word2}, \dots, \text{WordN} \mid \text{Spam})$$

$$= P(\text{Word1} \mid \text{Spam}) \cdot P(\text{Word2} \mid \text{Spam}) \cdot \dots \cdot P(\text{WordN} \mid \text{Spam})$$

1.15 Bayes rule

Example: Spam-Filter

Using this so-called **Naive Bayes Rule** leads to

$$P(\text{Spam} \mid w_1, w_2, \dots, w_N)$$

$$= \frac{P(w_1, w_2, \dots, w_N \mid \text{Spam})}{P(w_1, w_2, \dots, w_N \mid \text{Spam}) + P(w_1, w_2, \dots, w_N \mid \text{Ham})}$$

$$= \frac{P(w_1 \mid \text{Spam}) \cdot P(w_2 \mid \text{Spam}) \cdot \dots \cdot P(w_N \mid \text{Spam})}{P(w_1 \mid \text{Spam}) \cdot \dots \cdot P(w_N \mid \text{Spam}) + P(w_1 \mid \text{Ham}) \cdot \dots \cdot P(w_N \mid \text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

Similar we obtain

$$P(\text{Ham} \mid w_1, w_2, \dots, w_N)$$

$$= \frac{P(w_1 \mid \text{Ham}) \cdot P(w_2 \mid \text{Ham}) \cdot \dots \cdot P(w_N \mid \text{Ham})}{P(w_1 \mid \text{Spam}) \cdot \dots \cdot P(w_N \mid \text{Spam}) + P(w_1 \mid \text{Ham}) \cdot \dots \cdot P(w_N \mid \text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

Decision (Spam/Ham) may be based on **Likelihood-Ratio**:

$$\frac{P(\text{Spam} \mid w_1, \dots, w_N)}{P(\text{Ham} \mid w_1, \dots, w_N)} = \frac{P(w_1 \mid \text{Spam}) \cdot \dots \cdot P(w_N \mid \text{Spam})}{P(w_1 \mid \text{Ham}) \cdot \dots \cdot P(w_N \mid \text{Ham})}$$

$$\frac{P(\text{Spam} \mid w_1, \dots, w_N)}{P(\text{Ham} \mid w_1, \dots, w_N)} = \frac{\prod_i P(w_i \mid \text{Spam})}{\prod_i P(w_i \mid \text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

Likelihood-Ratio:

$$\frac{P(\text{Spam} \mid w_1, \dots, w_N)}{P(\text{Ham} \mid w_1, \dots, w_N)} = \prod_i \frac{P(w_i \mid \text{Spam})}{P(w_i \mid \text{Ham})}$$

Log-Likelihood-Ratio:

$$\log \frac{P(\text{Spam} \mid w_1, \dots, w_N)}{P(\text{Ham} \mid w_1, \dots, w_N)} = \sum_i \log \frac{P(w_i \mid \text{Spam})}{P(w_i \mid \text{Ham})}$$

1.15 Bayes rule

Example: Spam-Filter

Decision:

$$\sum_i \log \frac{P(w_i | \text{Spam})}{P(w_i | \text{Ham})} > 0 \quad \rightarrow \quad \text{Spam}$$

$$\sum_i \log \frac{P(w_i | \text{Spam})}{P(w_i | \text{Ham})} < 0 \quad \rightarrow \quad \text{Ham}$$

1.15 Bayes rule

Example: Spam-Filter

Example Spam Mail

Let us just take the first
4 words for the decision

Order Original Viagra directly from Pfizer

Here: <http://www.dummynameshop.com/>

All prices are tax/vat free and same-day free worldwide
shipping also included.

1.15 Bayes rule

Example: Spam-Filter

$$\begin{aligned} \text{Log-Lik-Ratio} &= \log \frac{P(\text{order} \mid \text{Spam})}{P(\text{order} \mid \text{Ham})} + \log \frac{P(\text{original} \mid \text{Spam})}{P(\text{original} \mid \text{Ham})} \\ &\quad + \log \frac{P(\text{viagra} \mid \text{Spam})}{P(\text{viagra} \mid \text{Ham})} + \log \frac{P(\text{directly} \mid \text{Spam})}{P(\text{directly} \mid \text{Ham})} \end{aligned}$$

1.15 Bayes rule

Example: Spam-Filter

Estimate $P(\text{word}|\text{Spam})$ and $P(\text{word}|\text{Ham})$ from training data:

Spam corpus: 467748 words

word	frequency in spam corpus	probability $P(\text{word} \text{Spam})$
order	355	$7.59 \cdot 10^{-4}$
original	179	$3.83 \cdot 10^{-4}$
viagra	422	$9.02 \cdot 10^{-4}$
directly	114	$2.44 \cdot 10^{-4}$

1.15 Bayes rule

Example: Spam-Filter

Estimate $P(\text{word}|\text{Spam})$ and $P(\text{word}|\text{Ham})$ from training data:

We use Hamlet as Ham corpus: 32630 words

(in praxis: corpora should have similar sizes)

word	frequency in ham corpus	probability $P(\text{word} \text{Ham})$
order	1	$3.06 \cdot 10^{-5}$
original	0	0
viagra	0	0
directly	1	$3.06 \cdot 10^{-5}$

Problem: Division by zero

1.15 Bayes rule

Example: Spam-Filter

Practical solution to the problem of unseen words

Take probability mass from 'singletons' (words with only one observation) and distribute it equally to unseen words.

word	frequency in ham corpus	probability $P(\text{word} \text{Spam})$
order	1	$2.36 \cdot 10^{-5}$
original	0	$0.13 \cdot 10^{-5}$
viagra	0	$0.13 \cdot 10^{-5}$
directly	1	$2.36 \cdot 10^{-5}$

reduced probability
(hypothetical numbers!)

1.15 Bayes rule

Example: Spam-Filter

Practical solution to the problem of unseen words

Take probability mass from 'singletons' (words with only one observation) and distribute it equally to unseen words.

word	frequency in ham corpus	probability $P(\text{word} \text{Spam})$
order	1	$2.36 \cdot 10^{-5}$
original	0	$0.13 \cdot 10^{-5}$ 
viagra	0	$0.13 \cdot 10^{-5}$ 
directly	1	$2.36 \cdot 10^{-5}$

increased probability (hypothetical numbers!)

1.15 Bayes rule

Example: Spam-Filter

$$\text{Log-Lik-Ratio} = \log \frac{P(\text{order} \mid \text{Spam})}{P(\text{order} \mid \text{Ham})} + \log \frac{P(\text{original} \mid \text{Spam})}{P(\text{original} \mid \text{Ham})} \\ + \log \frac{P(\text{viagra} \mid \text{Spam})}{P(\text{viagra} \mid \text{Ham})} + \log \frac{P(\text{directly} \mid \text{Spam})}{P(\text{directly} \mid \text{Ham})}$$

$$\text{Log-Lik-Ratio} = \log \frac{7.59 \cdot 10^{-4}}{2.36 \cdot 10^{-5}} + \log \frac{3.83 \cdot 10^{-4}}{0.13 \cdot 10^{-5}} \\ + \log \frac{9.02 \cdot 10^{-4}}{0.13 \cdot 10^{-5}} + \log \frac{2.44 \cdot 10^{-4}}{2.36 \cdot 10^{-5}} \\ = 18.04$$

word	frequency in spam corpus	probability $P(\text{word} \mid \text{Spam})$
order	355	$7.59 \cdot 10^{-4}$
original	179	$3.83 \cdot 10^{-4}$
viagra	422	$9.02 \cdot 10^{-4}$
directly	114	$2.44 \cdot 10^{-4}$

word	frequency in ham corpus	probability $P(\text{word} \mid \text{Spam})$
order	1	$2.36 \cdot 10^{-5}$
original	0	$0.13 \cdot 10^{-5}$
viagra	0	$0.13 \cdot 10^{-5}$
directly	1	$2.36 \cdot 10^{-5}$

1.15 Bayes rule

Example: Spam-Filter

$$\text{Log-Lik-Ratio} = \log \frac{P(\text{order} | \text{Spam})}{P(\text{order} | \text{Ham})} + \log \frac{P(\text{original} | \text{Spam})}{P(\text{original} | \text{Ham})} \\ + \log \frac{P(\text{viagra} | \text{Spam})}{P(\text{viagra} | \text{Ham})} + \log \frac{P(\text{directly} | \text{Spam})}{P(\text{directly} | \text{Ham})}$$

$$\text{Log-Lik-Ratio} = \log \frac{7.59 \cdot 10^{-4}}{2.36 \cdot 10^{-5}} + \log \frac{3.83 \cdot 10^{-4}}{0.13 \cdot 10^{-5}} \\ + \log \frac{9.02 \cdot 10^{-4}}{0.13 \cdot 10^{-5}} + \log \frac{2.44 \cdot 10^{-4}}{2.36 \cdot 10^{-5}} \\ = 18.04 \quad \text{→} \quad \text{Spam}$$

word	frequency in spam corpus	probability $P(\text{word} \text{Spam})$
order	355	$7.59 \cdot 10^{-4}$
original	179	$3.83 \cdot 10^{-4}$
viagra	422	$9.02 \cdot 10^{-4}$
directly	114	$2.44 \cdot 10^{-4}$

word	frequency in ham corpus	probability $P(\text{word} \text{Spam})$
order	1	$2.36 \cdot 10^{-5}$
original	0	$0.13 \cdot 10^{-5}$
viagra	0	$0.13 \cdot 10^{-5}$
directly	1	$2.36 \cdot 10^{-5}$