

Statistical Methods in AI (CSE 471)

The Generative Way

Vineet Gandhi
Centre for Visual Information Technology (CVIT)



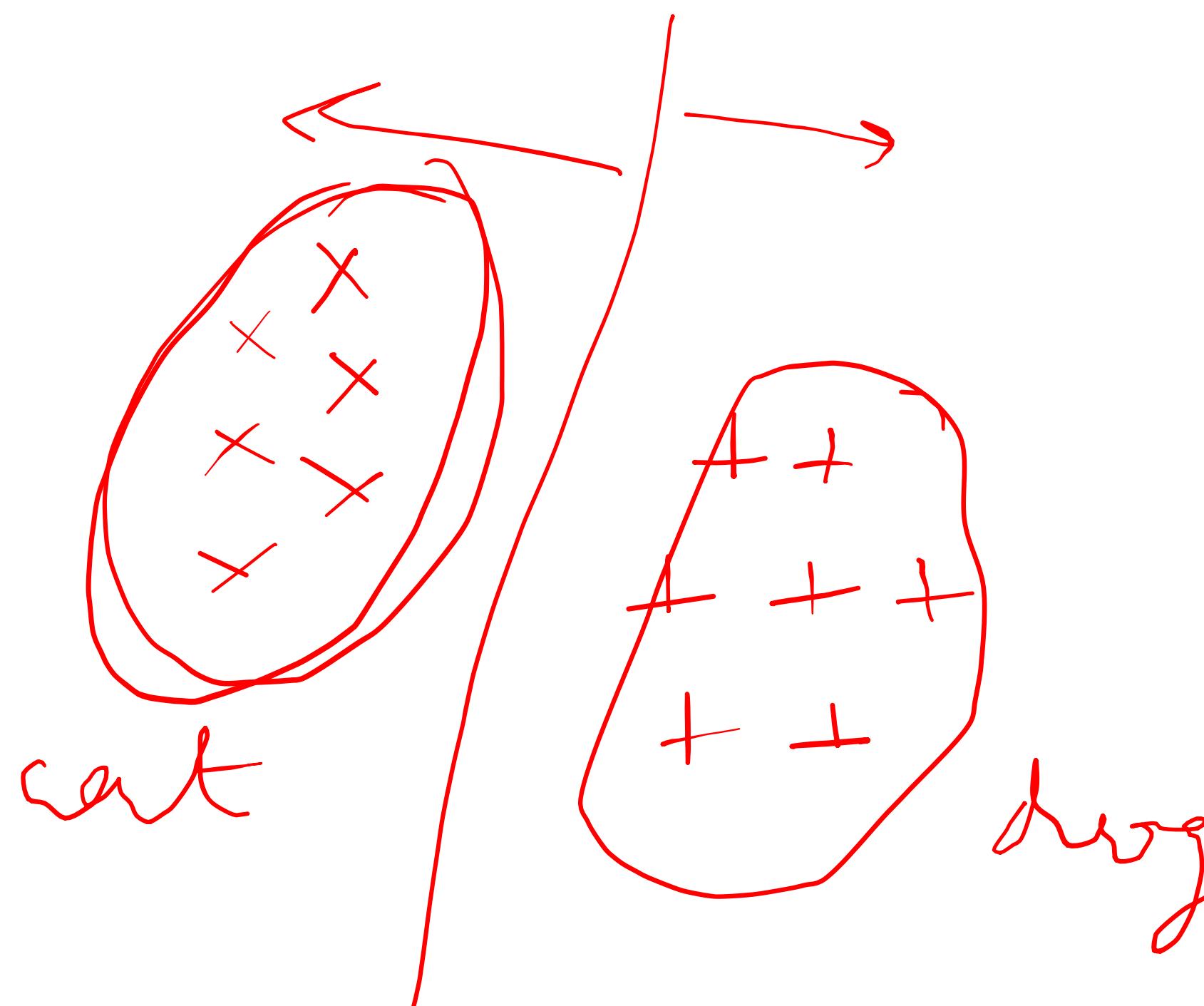
Many slides, figures, examples from: Ravi Kiran, Victor Lavrenko, Sho Nakagome, Andrew Ng,
Ben Tasker, Amos Storkey

Overview

- Basics
- Gaussian Distribution - univariate + multivariate
- How to generate gaussian distribution
- Density estimation + Anomaly detection
- Bayesian Classification + Naive Bayes

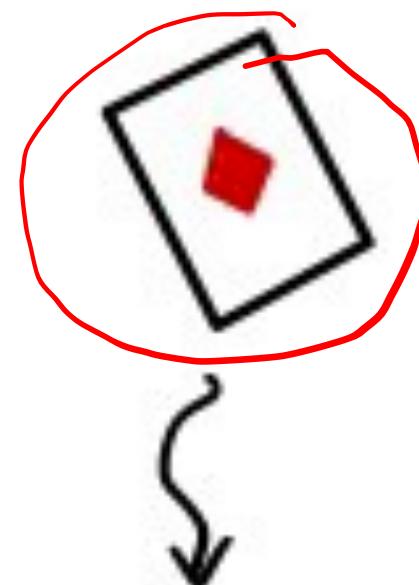
Generative vs Discriminative way

- When we estimate $P(X, Y) = P(X|Y)P(Y)$, then we call it *generative learning*.
- When we only estimate $P(Y|X)$ directly, then we call it *discriminative learning*.



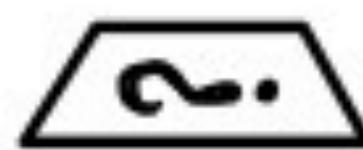
Probability and Random Variables

Probability



$$\begin{cases} x = 1 & (\text{Front}) \\ x = 0 & (\text{Back}) \end{cases}$$

$$\Pr(x=1) = 0.4$$



$$\Pr(x=0) = 0.6$$

$$\Pr(x=1) + \Pr(x=0) = 1$$

- Random Variable (RV) is a variable representing the outcome of our interest
- RV of x is used to represent the state of a card. It only has 2 states, either facing up ($x = 1$) or facing down ($x = 0$)
- RV could change it's state every time the event occurs. So if you dropped a card 5 times, x could be $[0, 0, 1, 0, 1]$
- $\Pr(x=0) = 3/5 = 0.6$

Joint probability

$$p(x, y)$$

Joint Probability



$\begin{cases} x = 1 & (\text{Rains}) \\ x = 0 & (\text{Doesn't rain}) \end{cases}$



$\begin{cases} y = 1 & (\text{Have umbrella}) \\ y = 0 & (\text{Don't have umbrella}) \end{cases}$

$$\begin{array}{l|l} \Pr(x=1) = 0.6 & \\ \Pr(x=0) = 0.4 & \\ \Pr(y=1) = 0.3 & \\ \Pr(y=0) = 0.7 & \checkmark \end{array}$$

Case 1 : Rains but you have an umbrella

$$\begin{aligned} \Pr(x=1, y=1) &= \Pr(x=1) \times \Pr(y=1) \\ &= 0.6 \times 0.3 \\ &= 0.18 \end{aligned}$$

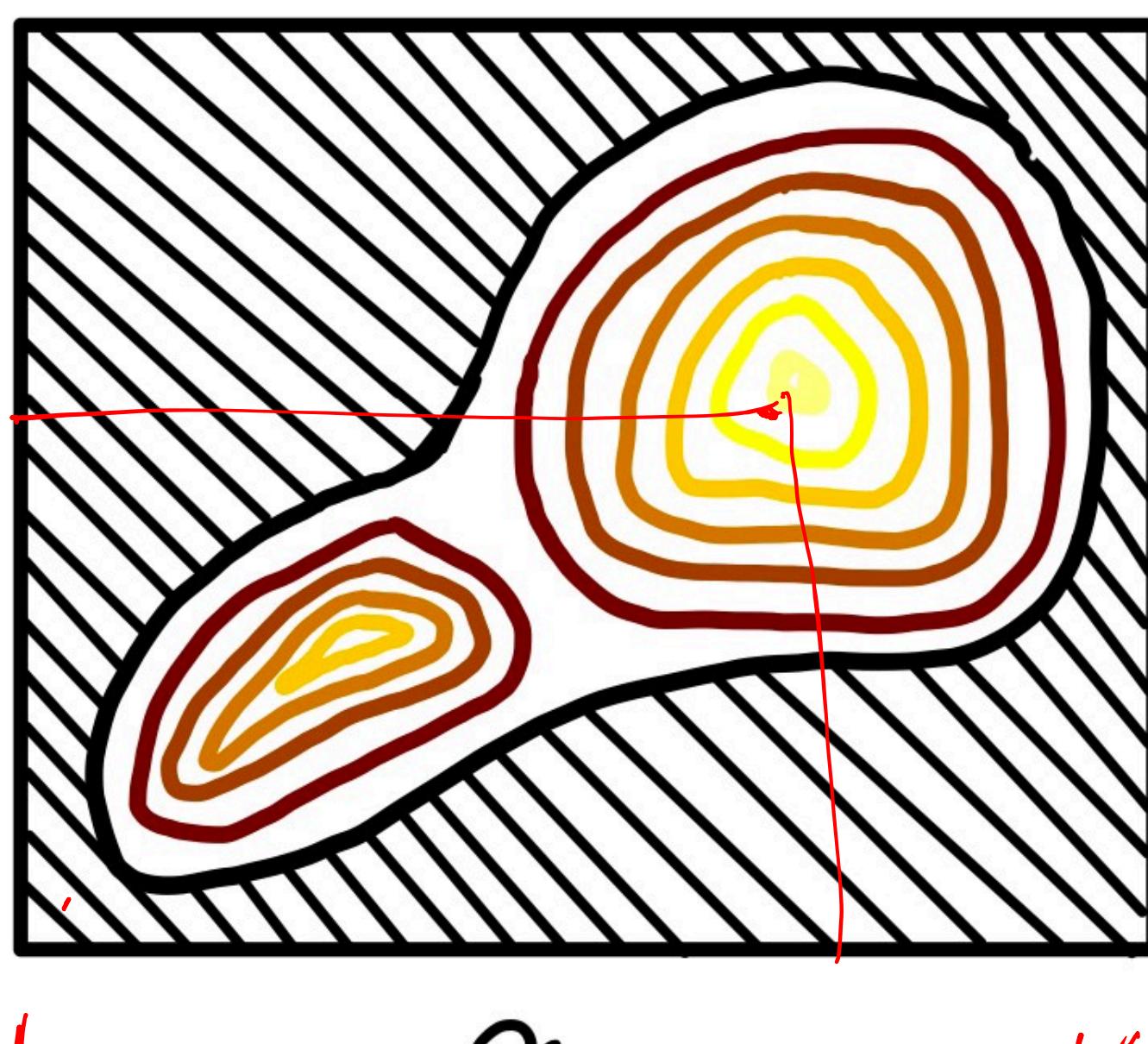
Case 2 : Rains but you DON'T have an umbrella

$$\begin{aligned} \Pr(x=1, y=0) &= \Pr(x=1) \times \Pr(y=0) \\ &= 0.6 \times 0.7 \\ &= 0.42 \end{aligned}$$

Joint probability

$P(x=1, y=1)$
= $\#(x=1 \& y=1)$
 $\#_{\text{obs}}$

Joint Probability in 2D



Marginalization

Marginalization

$$\left\{ \begin{array}{l} \Pr(x=0, y=0) = 0.28 \\ \Pr(x=1, y=0) = 0.42 \\ \Pr(x=0, y=1) = 0.12 \\ \Pr(x=1, y=1) = 0.18 \end{array} \right.$$

How to get $\Pr(x=0)$?

$$\Pr(x) = \int \Pr(x, y) dy$$

(Continuous)

$$\Pr(x) = \sum_y \Pr(x, y)$$

(Discrete)

$$\Pr(x=0) = \sum_{y=0}^1 \Pr(x=0, y)$$

$$= \underbrace{\Pr(x=0, y=0)} + \underbrace{\Pr(x=0, y=1)}$$

$$= \underline{0.28} + \underline{0.12} = 0.4$$

Conditional probability

$$p(x|y) p(y) = p(y|x) p(x)$$

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



Given

$$x=1 \text{ (Rains)}$$



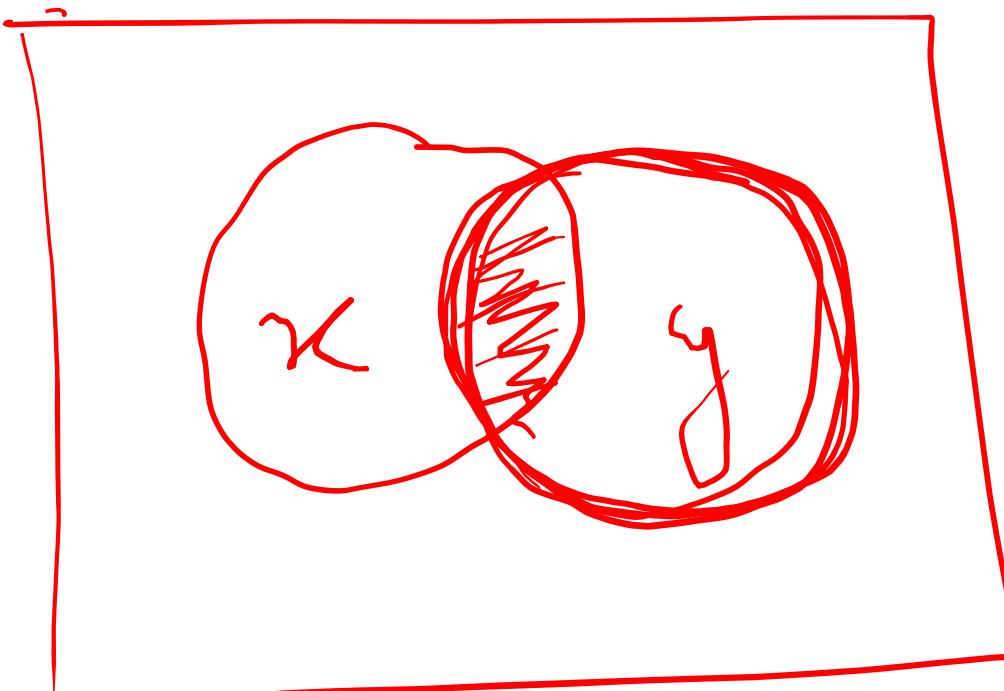
What's the Probability of
 $y=1$ (Bring umbrella)

Given y ,
probability of x

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

Joint probability
of x & y

Probability of y
(Normalization Term)



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

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

Marginalization

$$p(x=0) = p(x=0|y=0) + p(x=0|y=1)$$

$$p(x) = p(x,a) + p(x,b)$$

$$p(x) = p(x|a)p(a) + p(x|b)p(b)$$

$$p(x=0) = p(x=0|y=0) \cdot p(y=0)$$

$$+ p(x=0|y=1) p(y=1)$$

Bayes' Rule

A photograph of a large-scale neon sign for Bayes' Rule. The sign is illuminated in blue neon against a dark background. It displays the formula:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

The sign is mounted on a structure with visible cables and is set against a dark, grid-like ceiling.

- A disease occurs in 0.5% of population
- A diagnostic test gives a positive result
 - in 99% of people that have the disease
 - in 5% of people that do not have the disease (false positive)

A random person from the street is found to be positive on this test.
What is the probability that they have the disease?

A: 0-30%

B: 30-60%

C: 60-90%

- A disease occurs in 0.5% of population
- A diagnostic test gives a positive result
 - in 99% of people that have the disease
 - in 5% of people that do not have the disease (false positive)

A = disease

B = positive test result

$$\underline{P(A)} = 0.005 \quad \text{probability of having disease}$$

$$\underline{P(\sim A)} = 1 - 0.005 = 0.995 \quad \text{probability of not having disease}$$

$$\underline{P(B)} = 0.005 * 0.99 \text{ (people with disease)} + 0.995 * 0.05 \text{ (people without disease)} = 0.0547 \text{ (slightly more than 5\% of all tests are positive)}$$
✓

conditional probabilities

$$P(B|A) = 0.99 \quad \text{probability of pos result **given** you have disease}$$

$$P(\sim B|A) = 1 - 0.99 = 0.01 \quad \text{probability of neg result **given** you have disease}$$

$$P(B|\sim A) = 0.05 \quad \text{probability of pos result **given** you do not have disease}$$

$$P(\sim B|\sim A) = 1 - 0.05 = 0.95 \quad \text{probability of neg result **given** you do not have disease}$$

$P(A|B)$ is probability of disease **given** the test is positive (which is what we're interested in)

~~Very~~ different from $P(B|A)$: probability of positive test results given you have the disease.



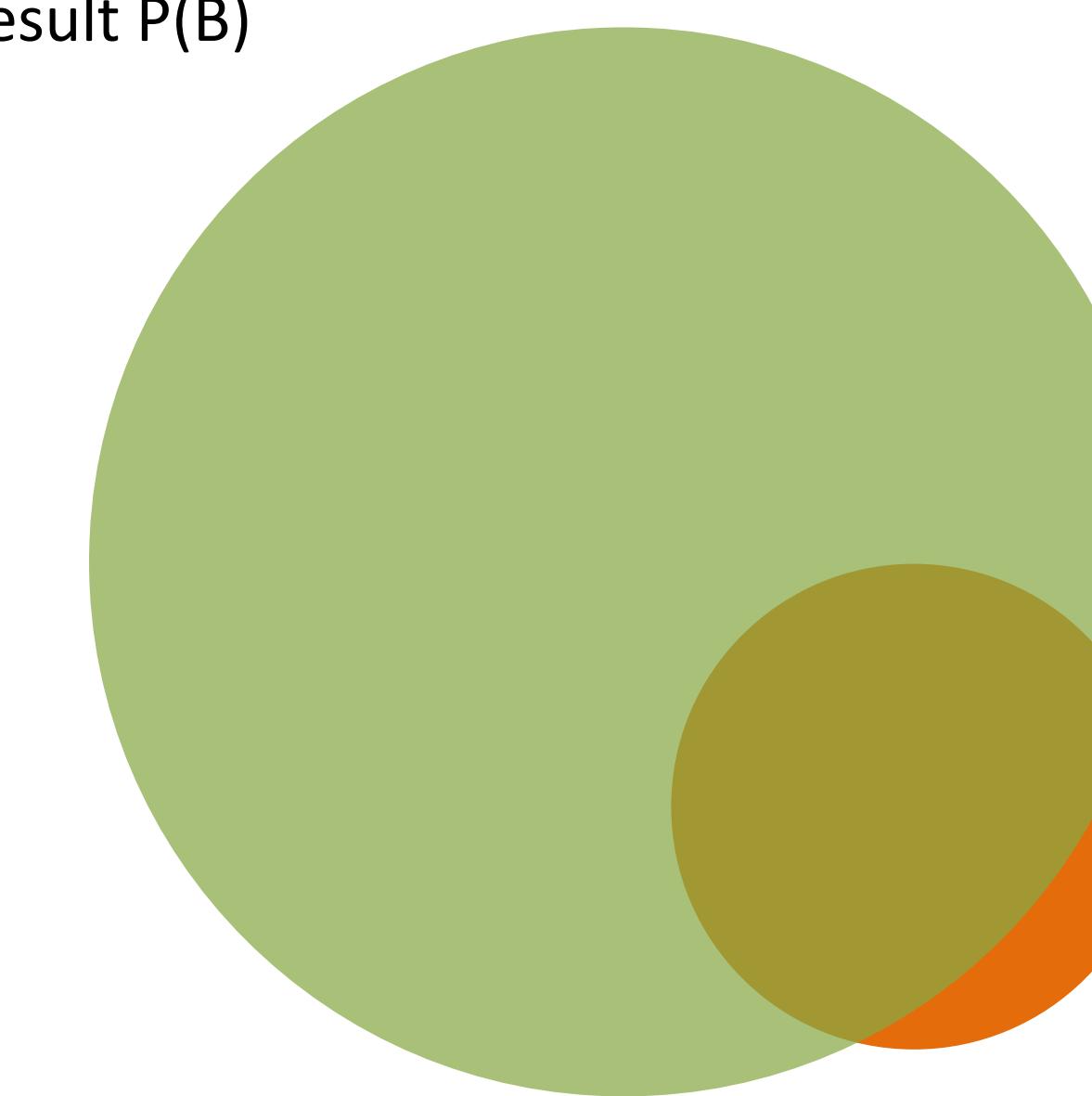
A = disease

B = positive test result

population = 1000

positive test result $P(B)$
0.0547

disease $P(A)$
0.005



A = disease

B = positive test result

$P(A, B)$ is the *joint probability*, or the probability that both events occur.

$P(A, B)$ is the same as $P(B, A)$.

But we already *know* that the test was positive, so we have to take that into account.

Of all the people already in the green circle, how many fall into the $P(A, B)$ part? That's the probability we want to know!

That is:

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

You can write down same thing for the inverse:

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

The *joint probability* can be expressed in two ways by rewriting the equations

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

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

Equating the two gives

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

$$P(A|B) = P(B|A) * P(A) / P(B)$$

population = 100

positive test result $P(B)$
0.0547

$P(B, \sim A)$

$P(A, B)$

disease $P(A)$
0.005

$P(A, \sim B)$

A = disease

B = positive test result

$$P(A) = 0.005$$

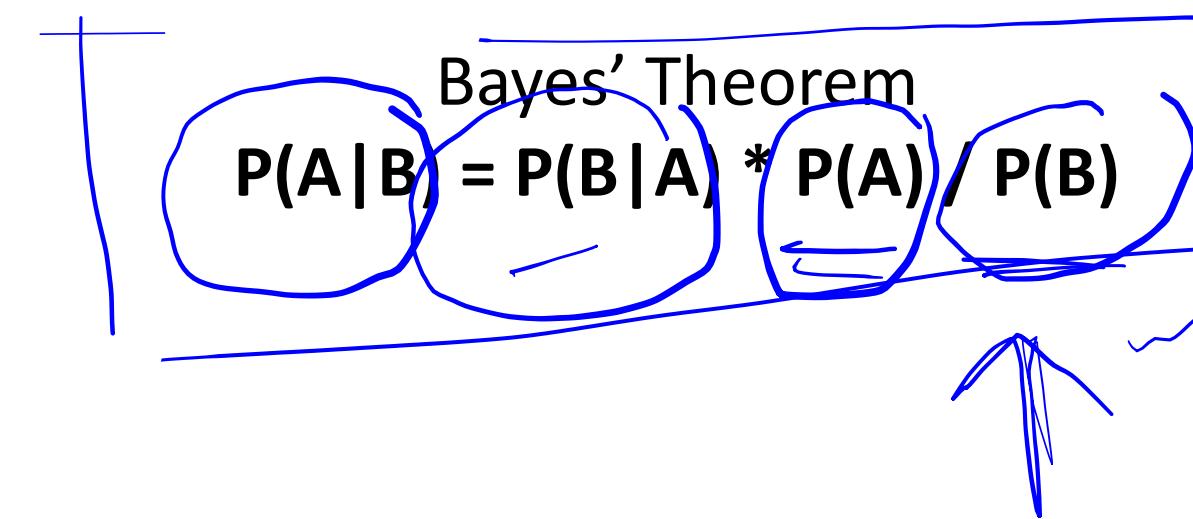
probability of having disease

$$P(B|A) = 0.99$$

probability of pos result **given** you have disease

$$P(B) = 0.005 * 0.99 \text{ (people with disease)} + 0.995 * 0.05 \text{ (people without disease)} = 0.0547$$

$$P(A|B) = 0.99 * 0.005 / 0.0547 \\ = 0.09$$



So a positive test result increases your probability of having the disease to 'only' 9%, simply because the disease is very rare (relative to the false positive rate).

$P(A)$ is called the **prior**: before we have any information, we estimate the chance of having the disease 0.5%

$P(B|A)$ is called the **likelihood**: probability of the data (pos test result) given an underlying cause (disease)

$P(B)$ is the **marginal probability of the data**: the probability of observing this particular outcome, taken over all possible values of A (disease and no disease)

$P(A|B)$ is the **posterior probability**: it is a combination of what you thought before obtaining the data, and the new information the data provided (combination of **prior** and **likelihood**)

Let's do another one...

It rains on 20% of days.

When it rains, it was forecasted 80% of the time

When it doesn't rain, it was erroneously forecasted 10% of the time.

The weatherman forecasts rain. What's the probability of it actually raining?

A = forecast rain

B = it rains

What information is given in the story?

$$P(B) = 0.2 \text{ (prior)}$$

$$P(A|B) = 0.8 \text{ (likelihood)}$$

$$P(A|\sim B) = 0.1$$

$$P(B|A) = P(A|B) * P(B) / P(A)$$

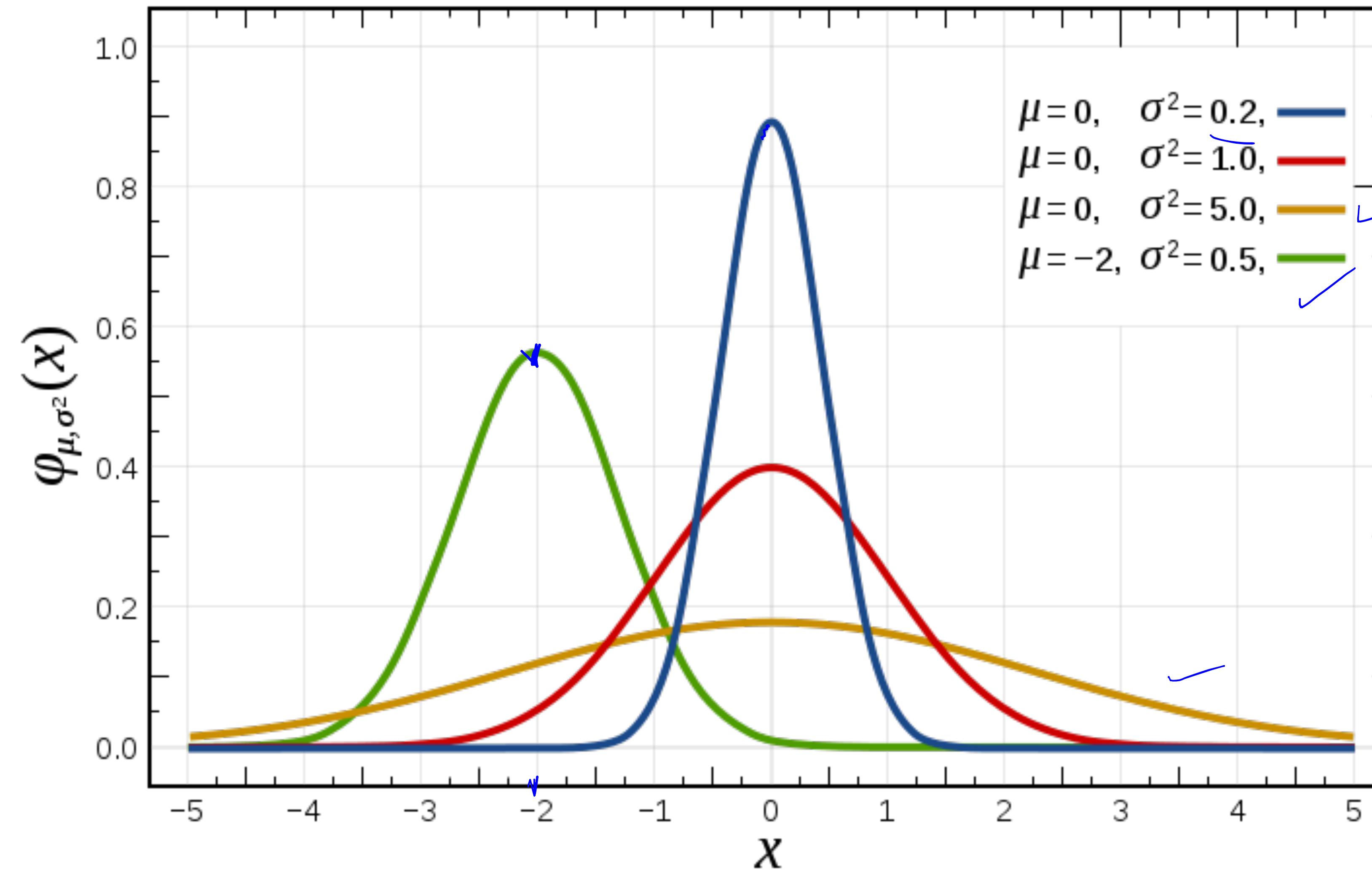
What is $P(A)$, probability of rain forecast? Calculate over all possible values of B (**marginal probability**)

$$P(A|B) * P(B) + P(A|\sim B) * P(\sim B) = 0.8 * 0.2 + 0.1 * 0.8 = 0.24$$

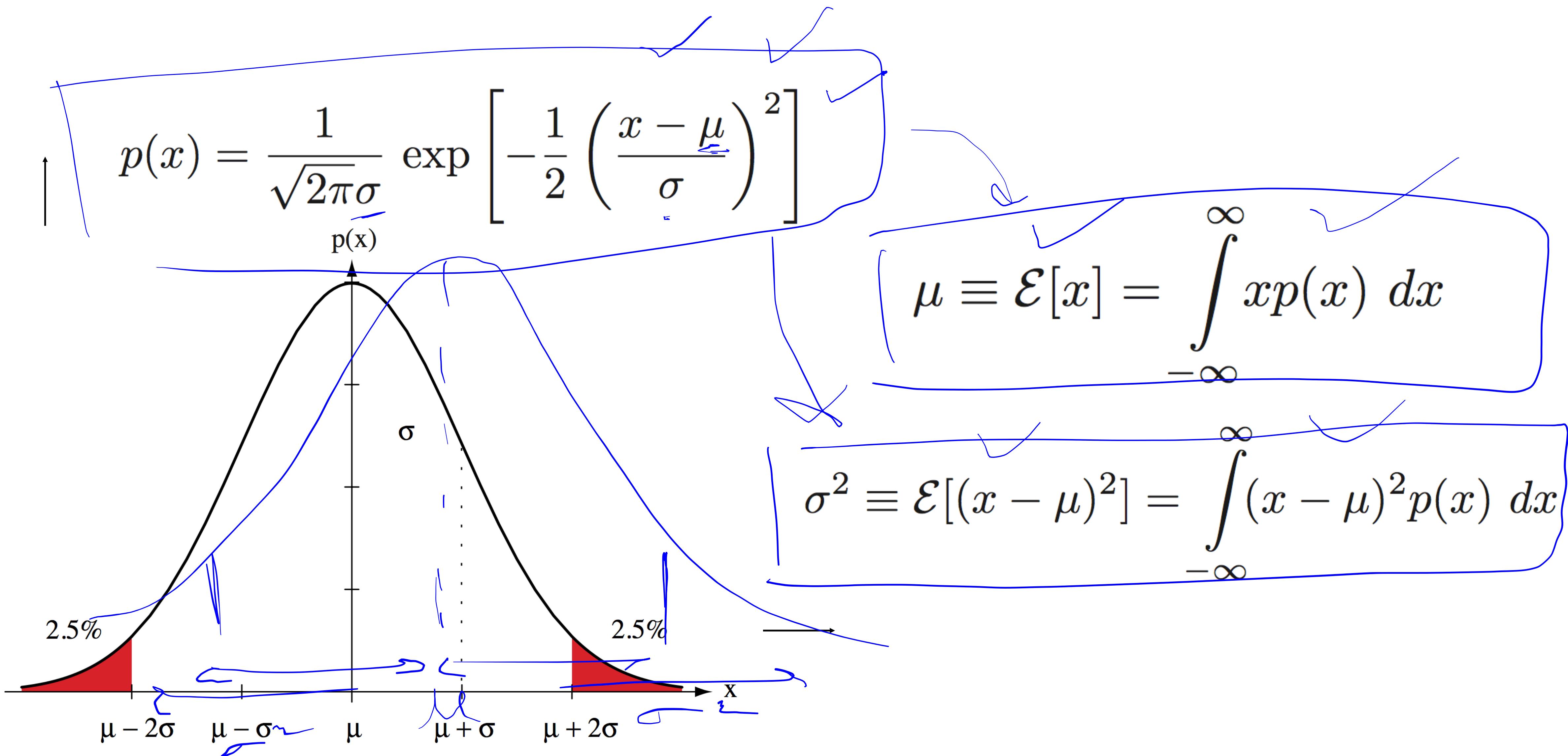
$$\begin{aligned} P(B|A) &= 0.8 * 0.2 / 0.24 \\ &= 0.67 \end{aligned}$$

So before you knew anything you thought $P(\text{rain})$ was 0.2. Now that you heard the weather forecast, you adjust your expectation upwards $P(\text{rain}|\text{forecast}) = 0.67$

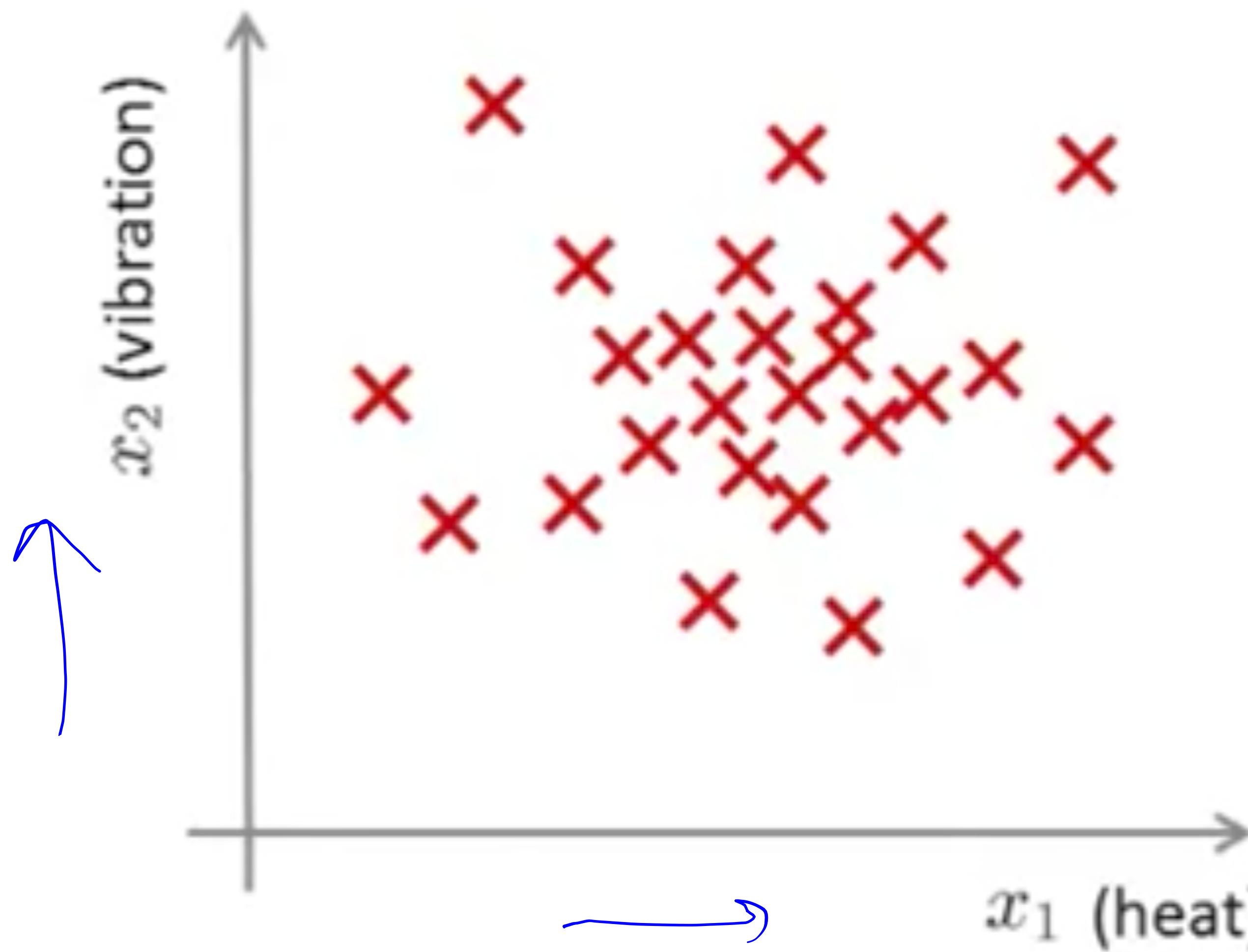
Univariate Gaussian Distribution



Parameter estimation



Anomaly Detection



$$p(x)$$

$$b(n_1 n_2 \dots n_c)$$

$$f(x)$$

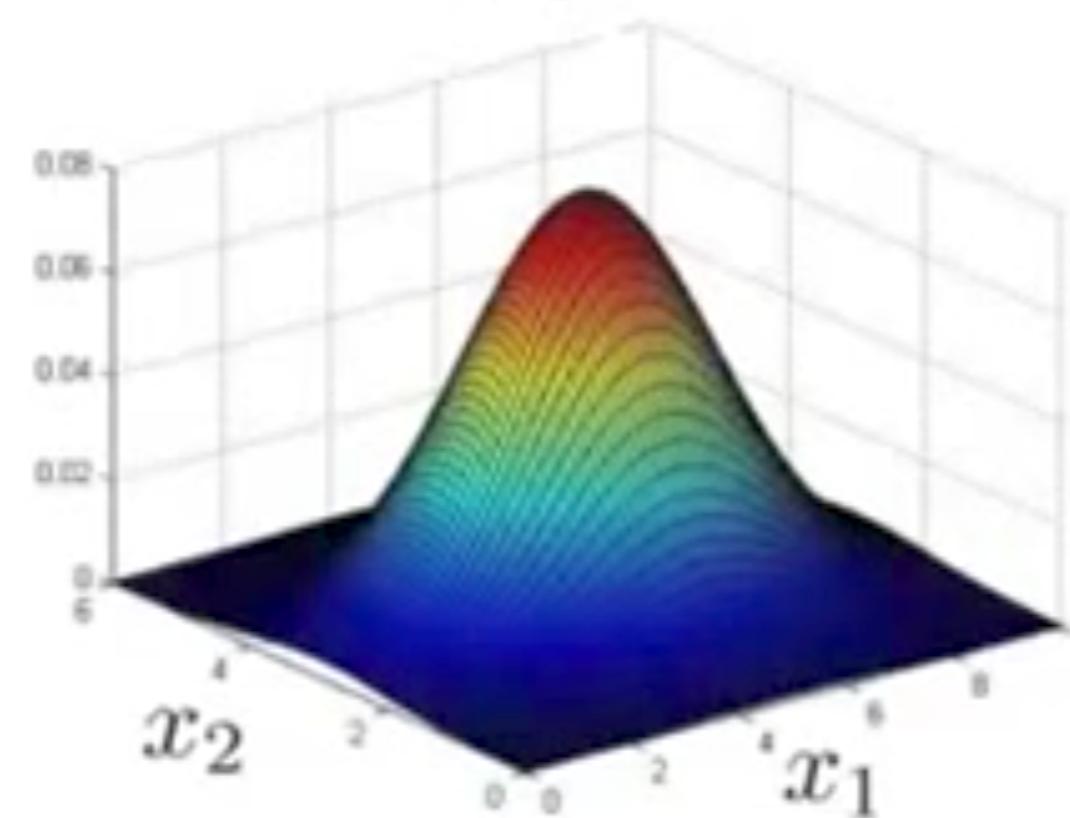
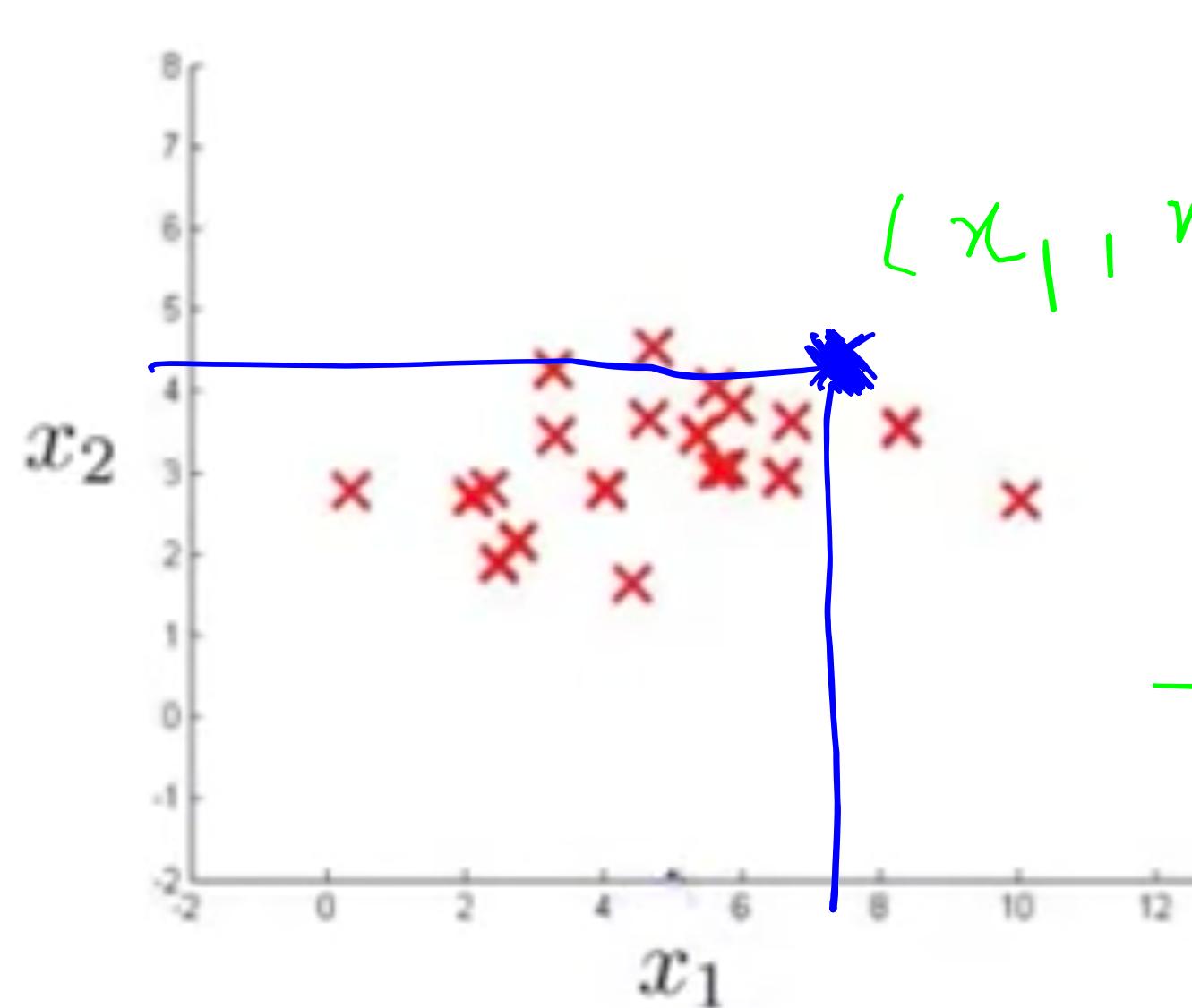
figure credit: Andrew Ng

Density Estimation

$$p(x) = p(x_1, x_2, x_3, \dots, x_d) = \dots$$

$$= p(x_1 | u_1, \sigma_1) \cdot p(x_2 | u_2, \sigma_2) \cdot p(x_3 | u_3, \sigma_3) \cdot \dots$$

Example: Anomaly Detection



$$= 0.2 \times 0.1$$

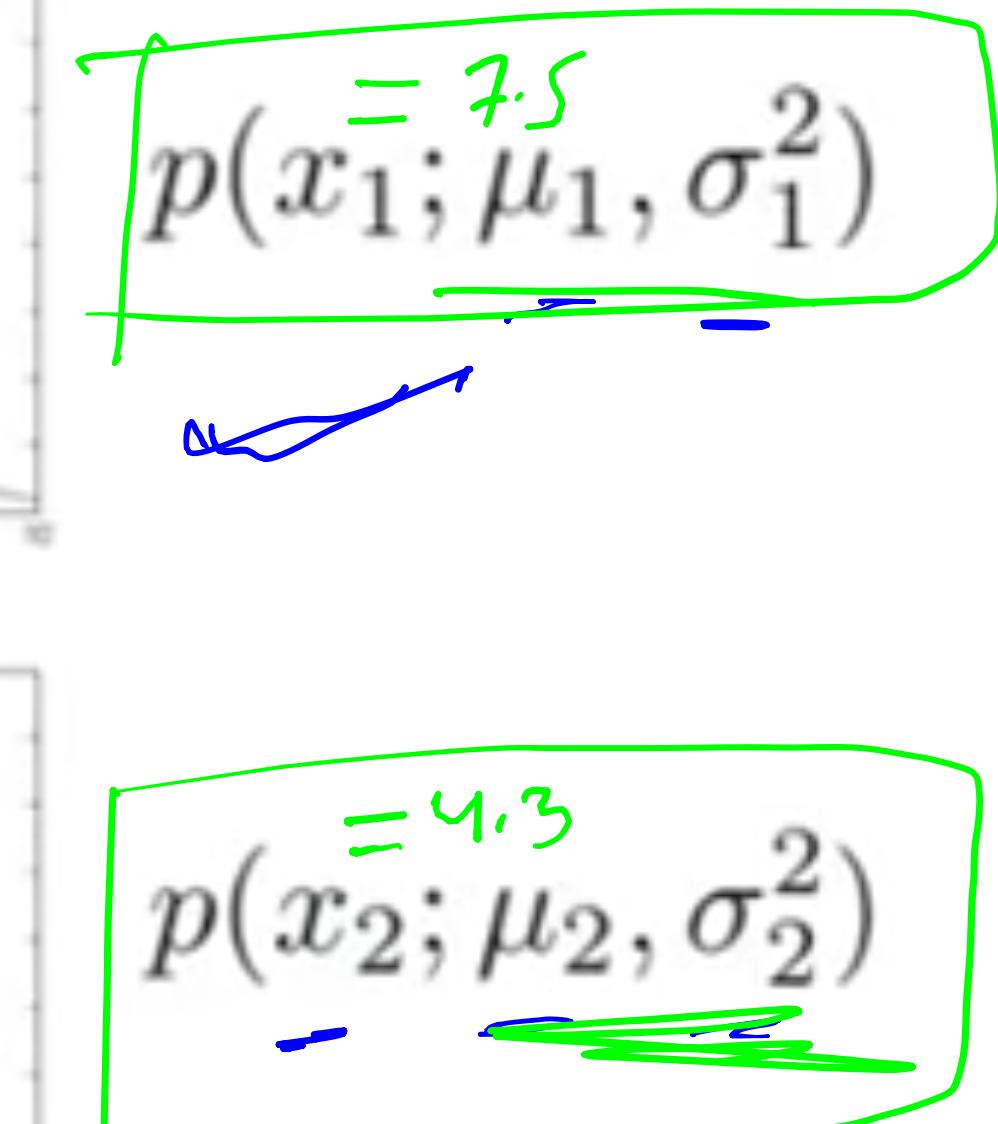
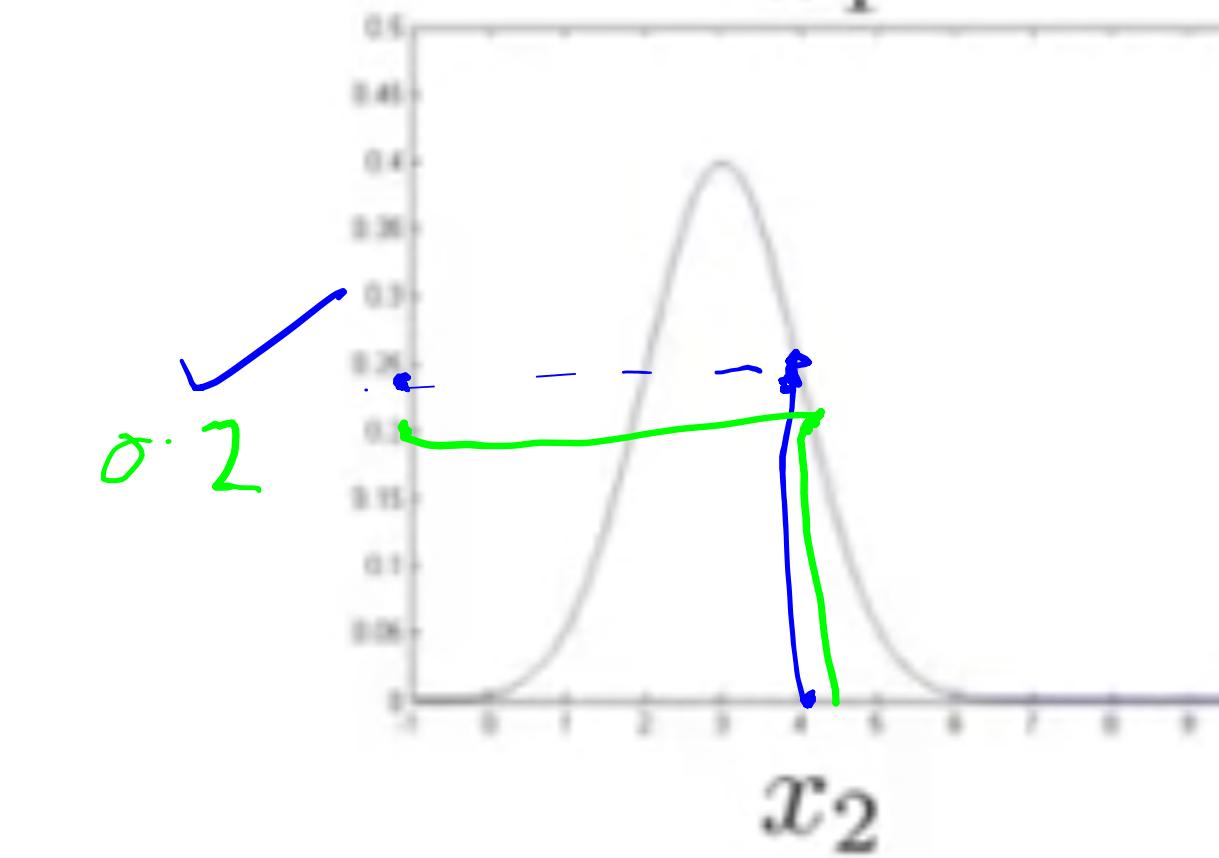
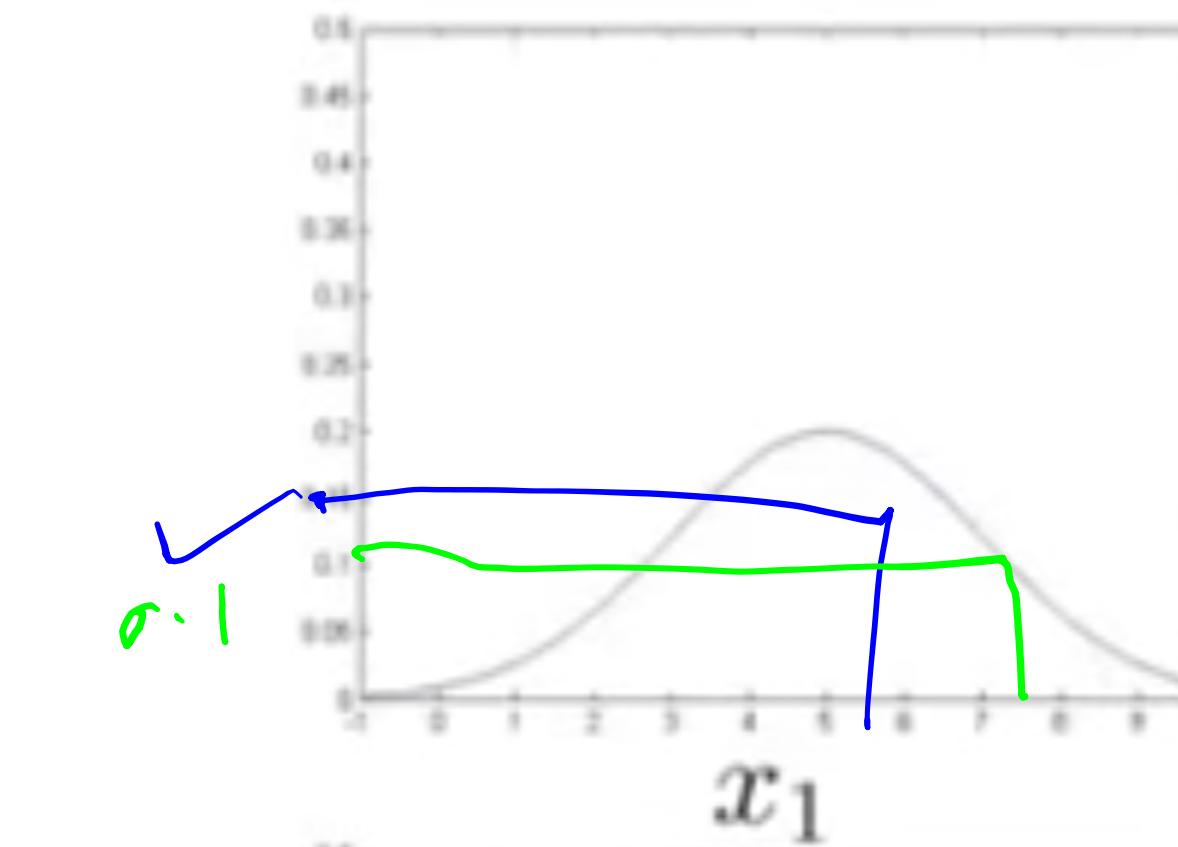


figure credit: Andrew Ng

Multivariate Gaussian Distribution

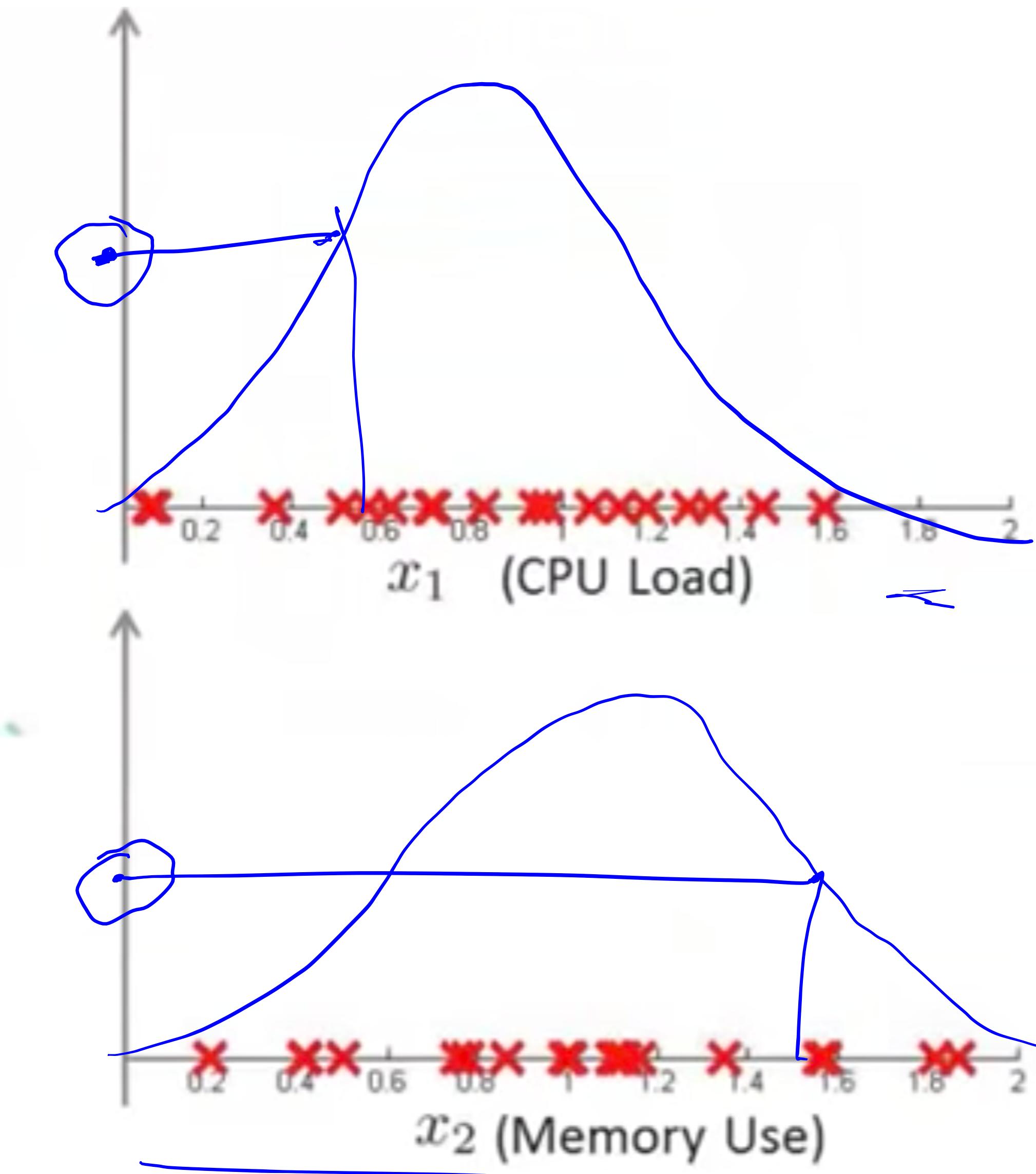
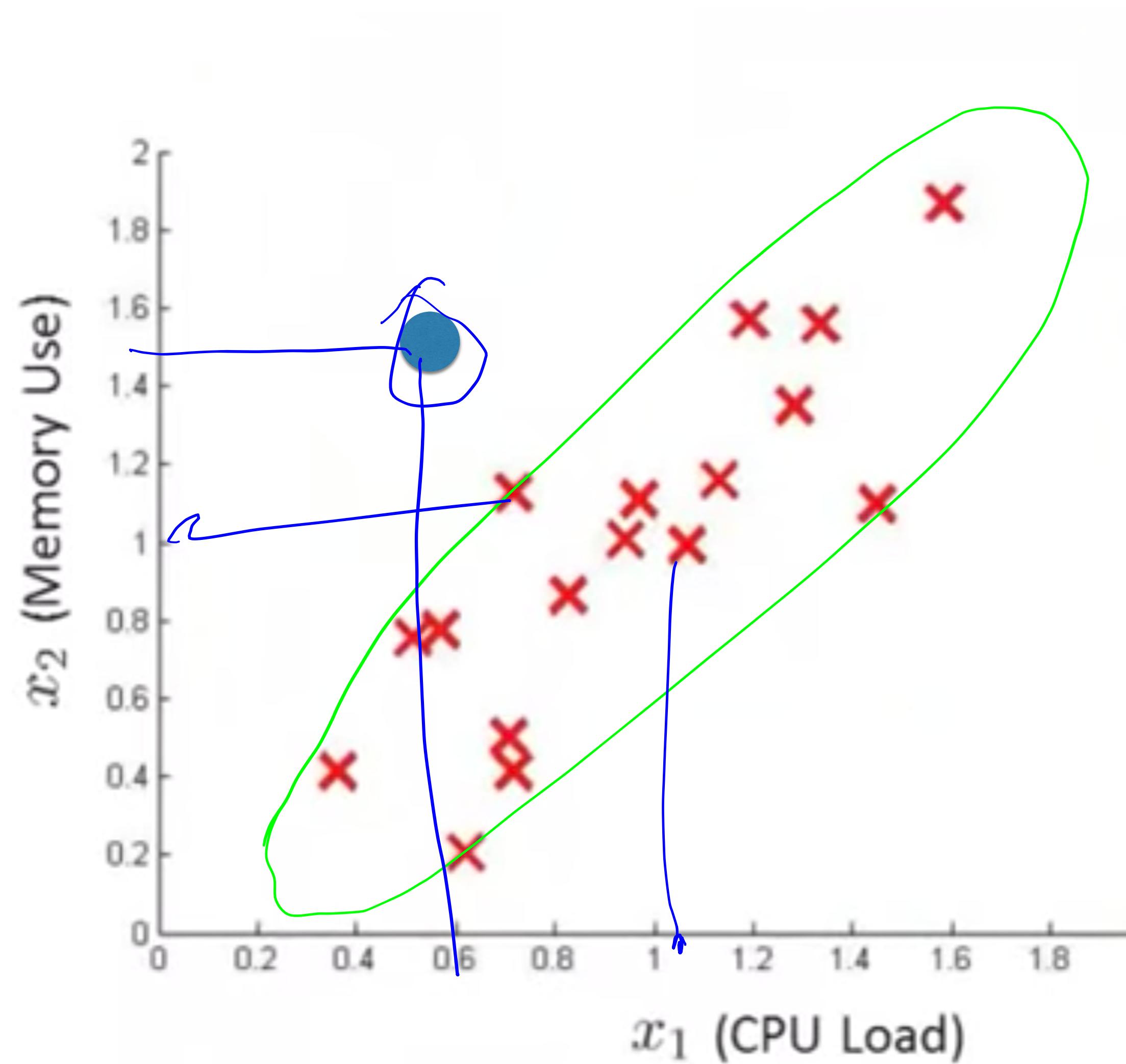


figure credit: Andrew Ng

Multivariate Gaussian Distribution

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^t \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$\boldsymbol{\mu} \equiv \mathcal{E}[\mathbf{x}] = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x}$$

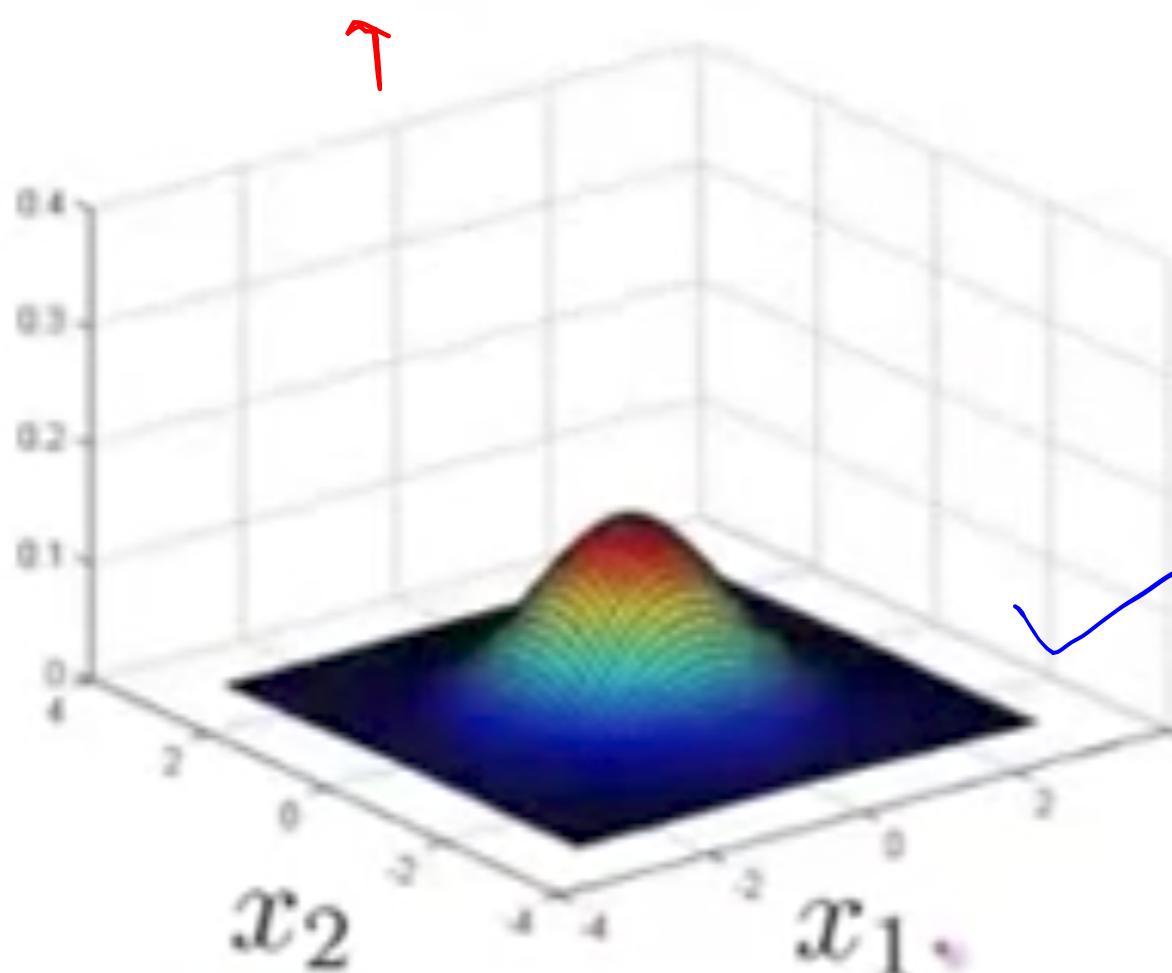
$$\boldsymbol{\Sigma} \equiv \mathcal{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t] = \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^t p(\mathbf{x}) d\mathbf{x},$$

$$\mu_i = \mathcal{E}[x_i]$$

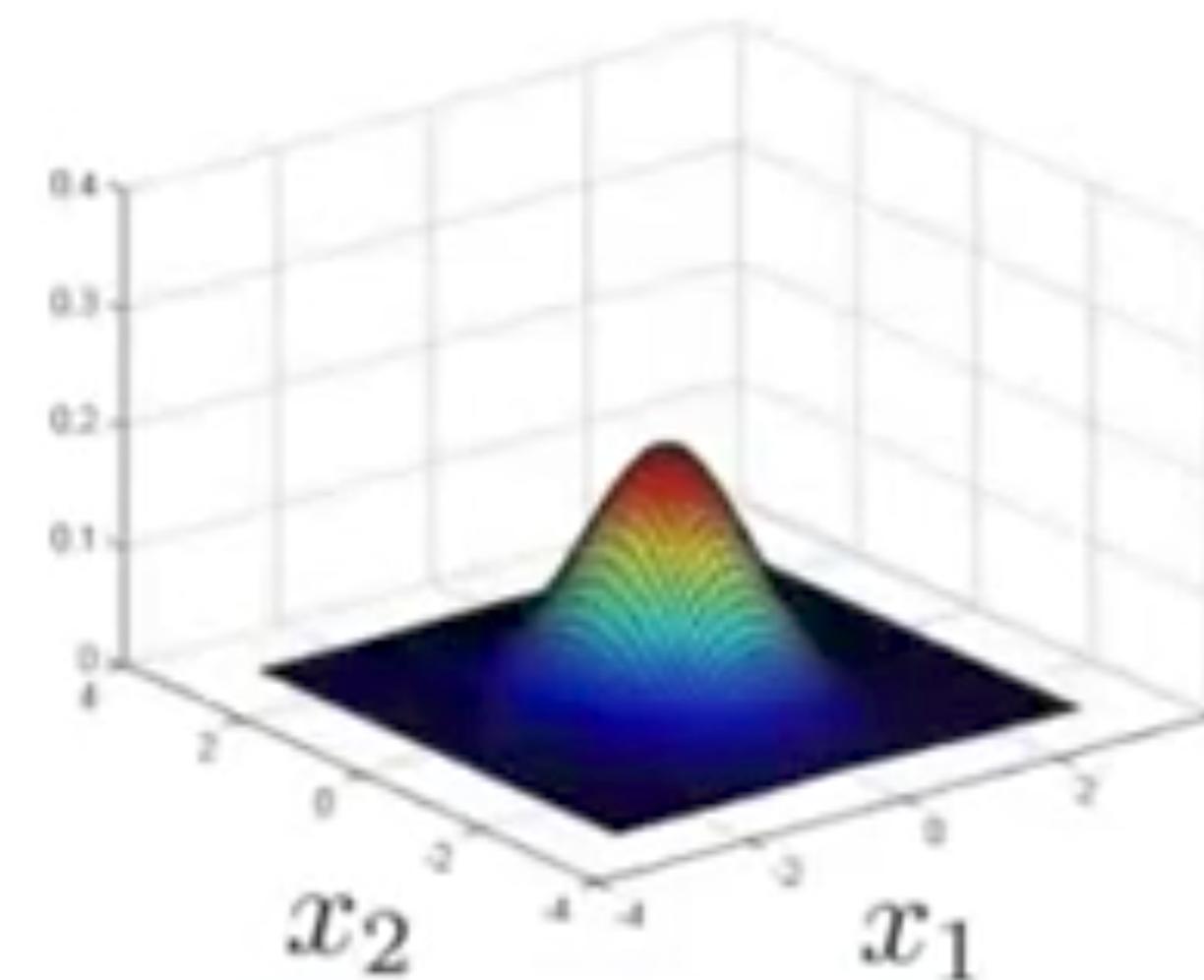
$$\sigma_{ij} = \mathcal{E}[(x_i - \mu_i)(x_j - \mu_j)].$$

Multivariate Gaussian Distribution

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



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



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

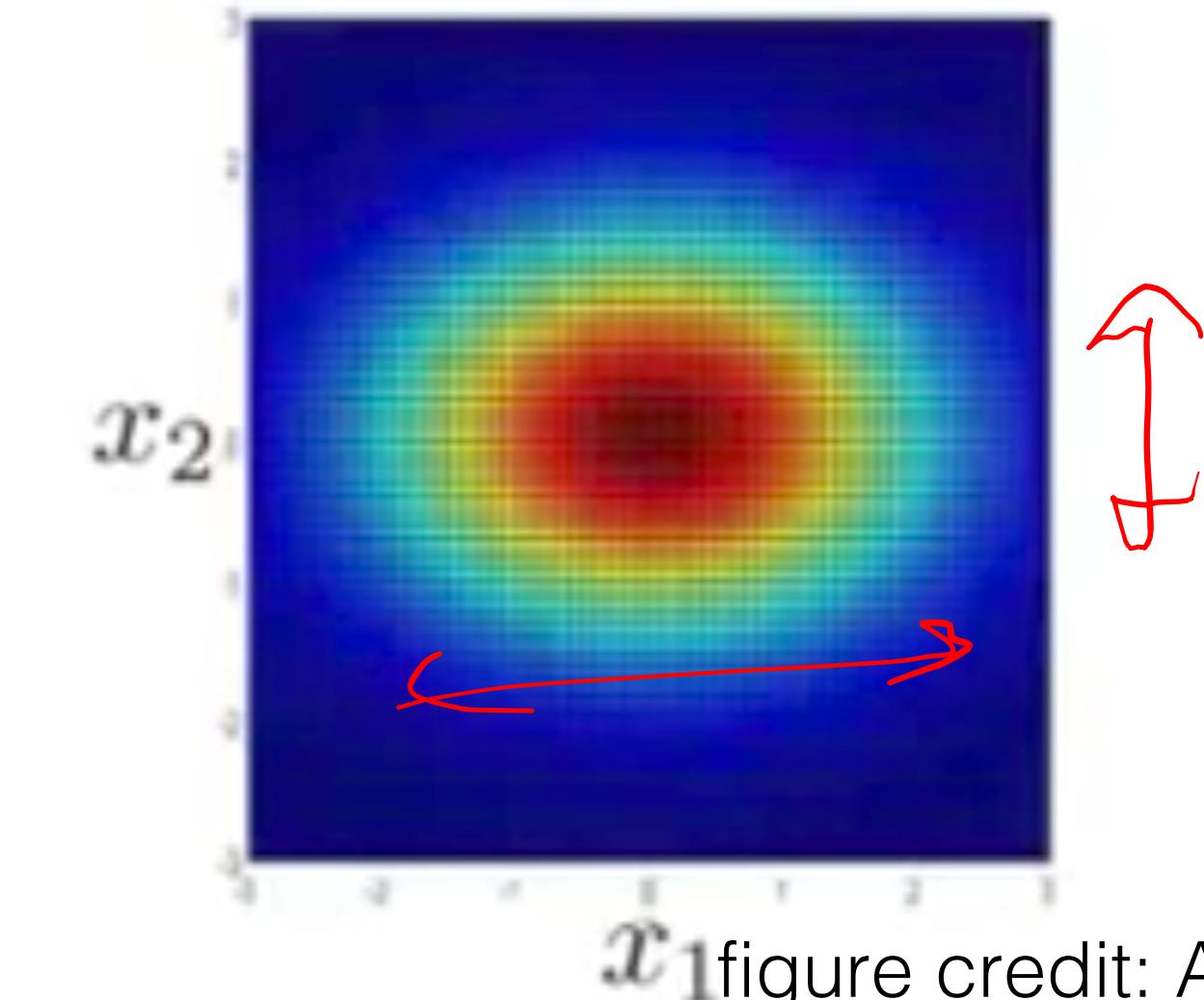
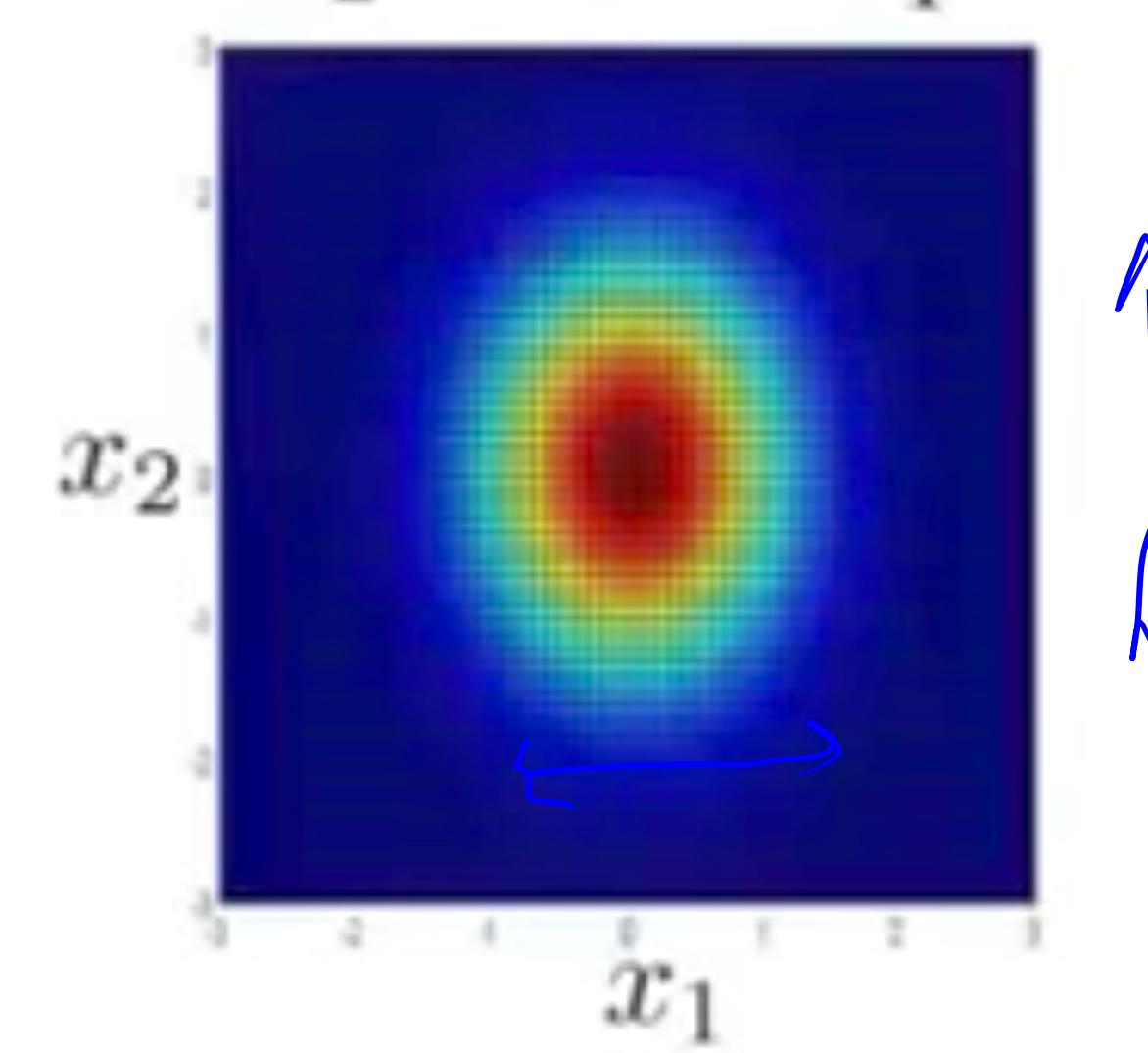
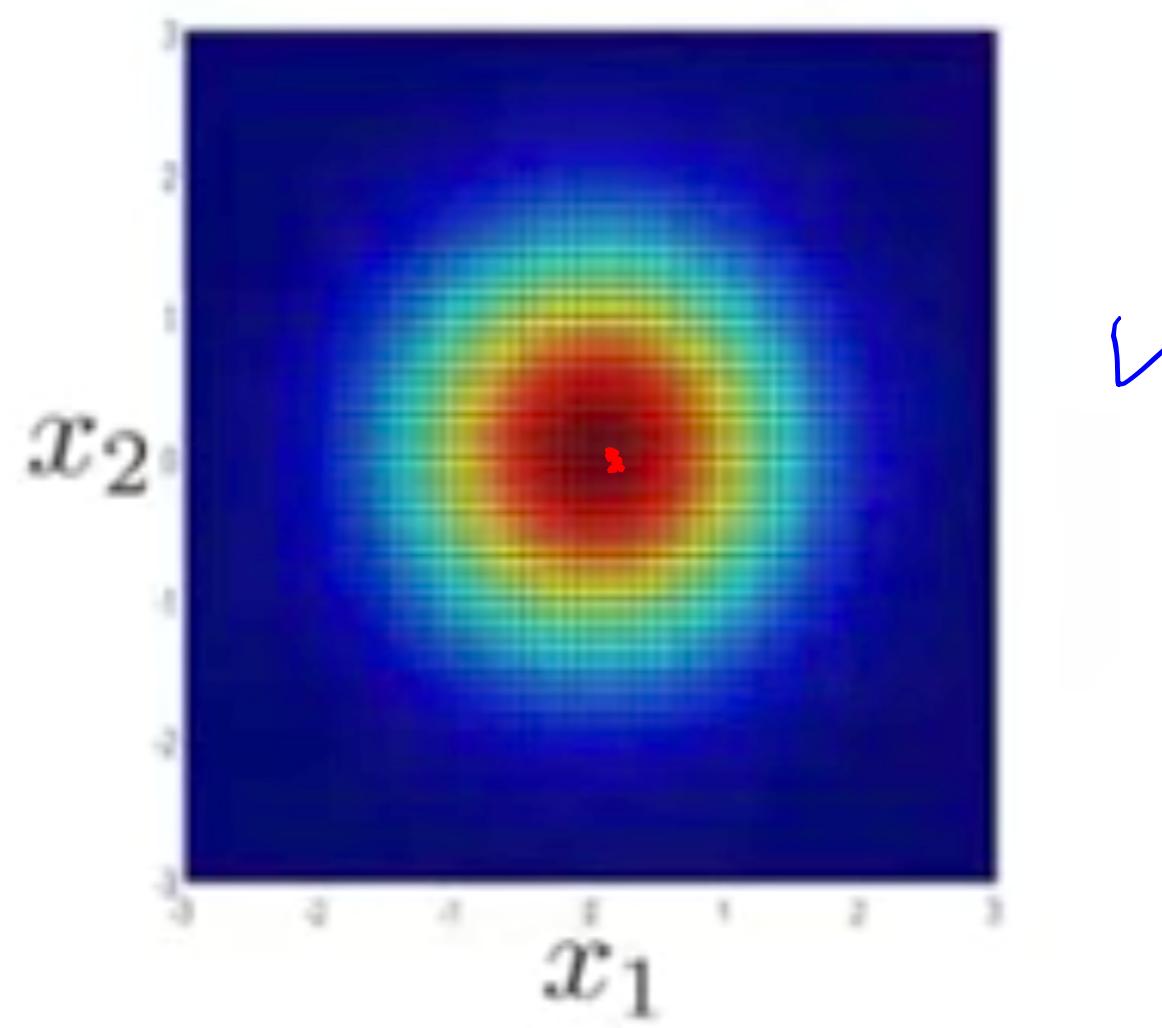
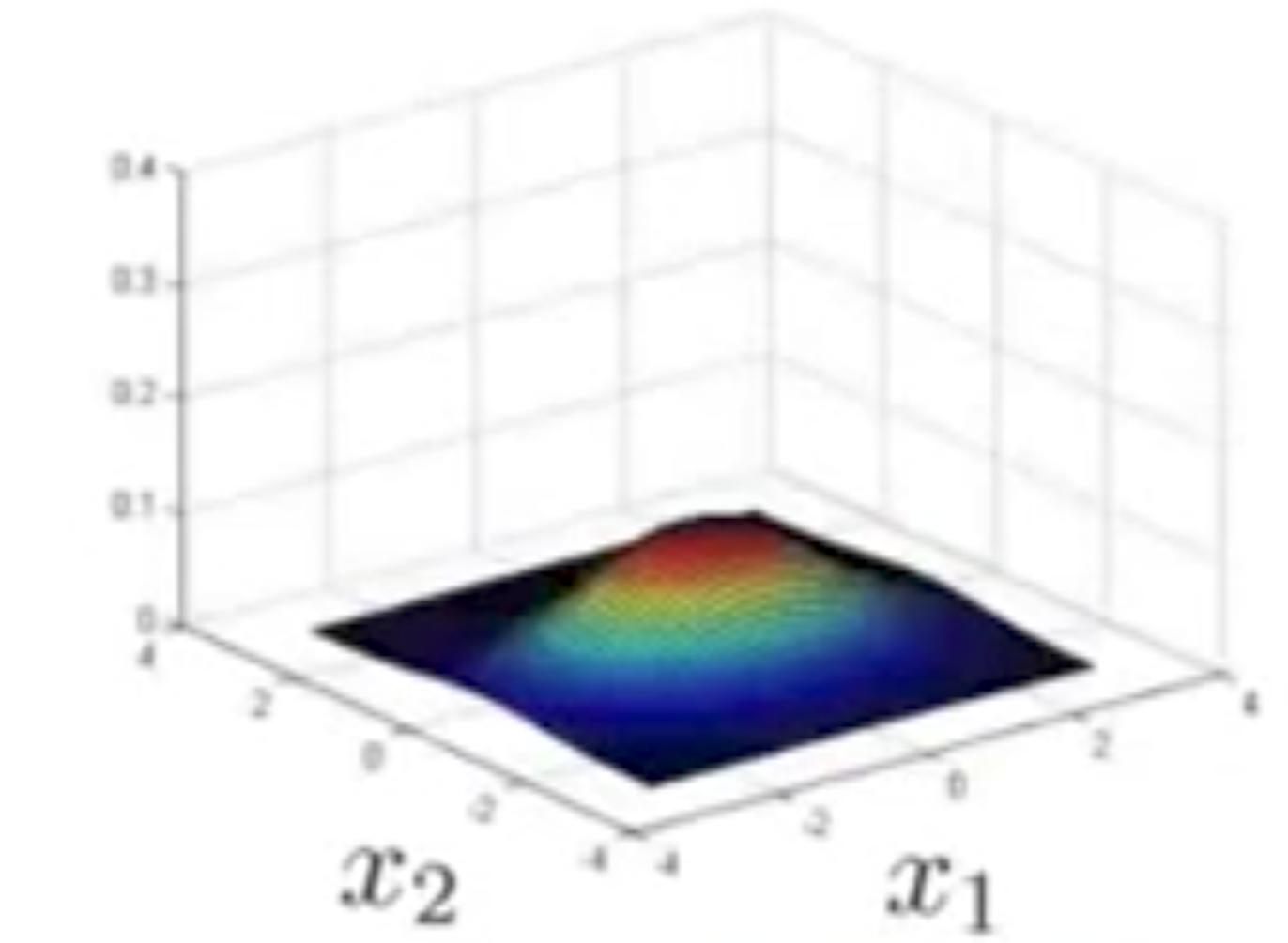
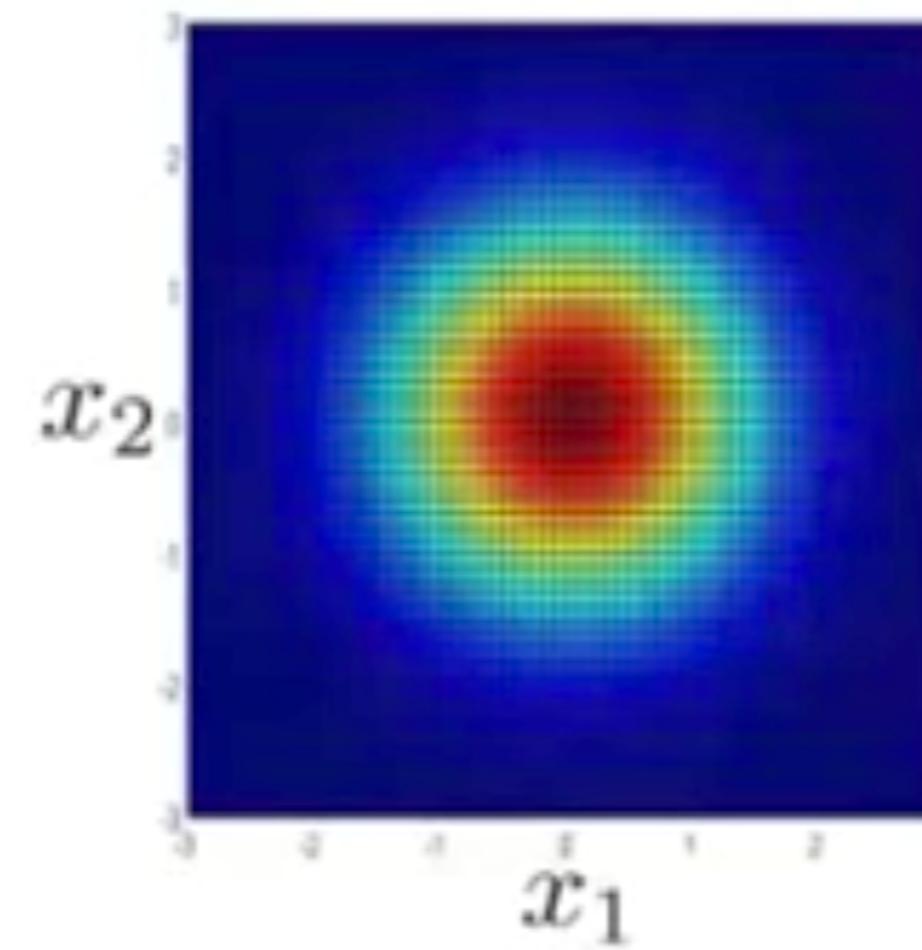
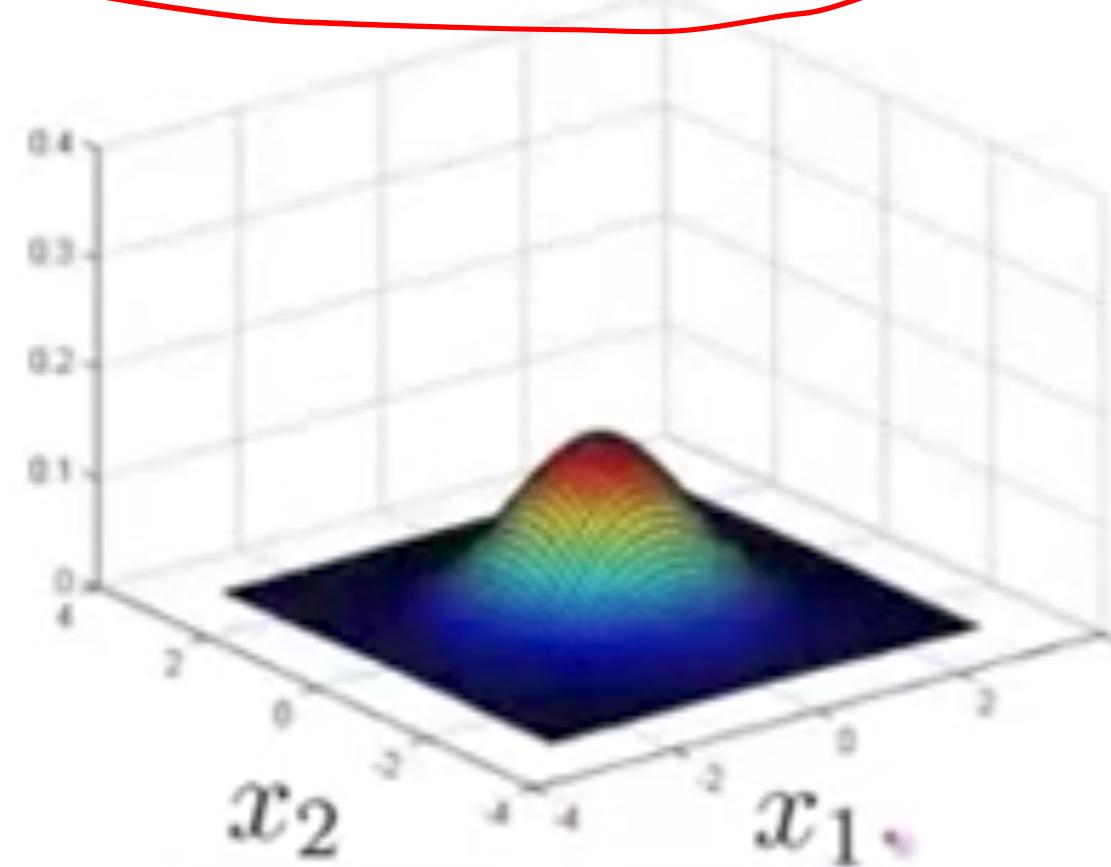


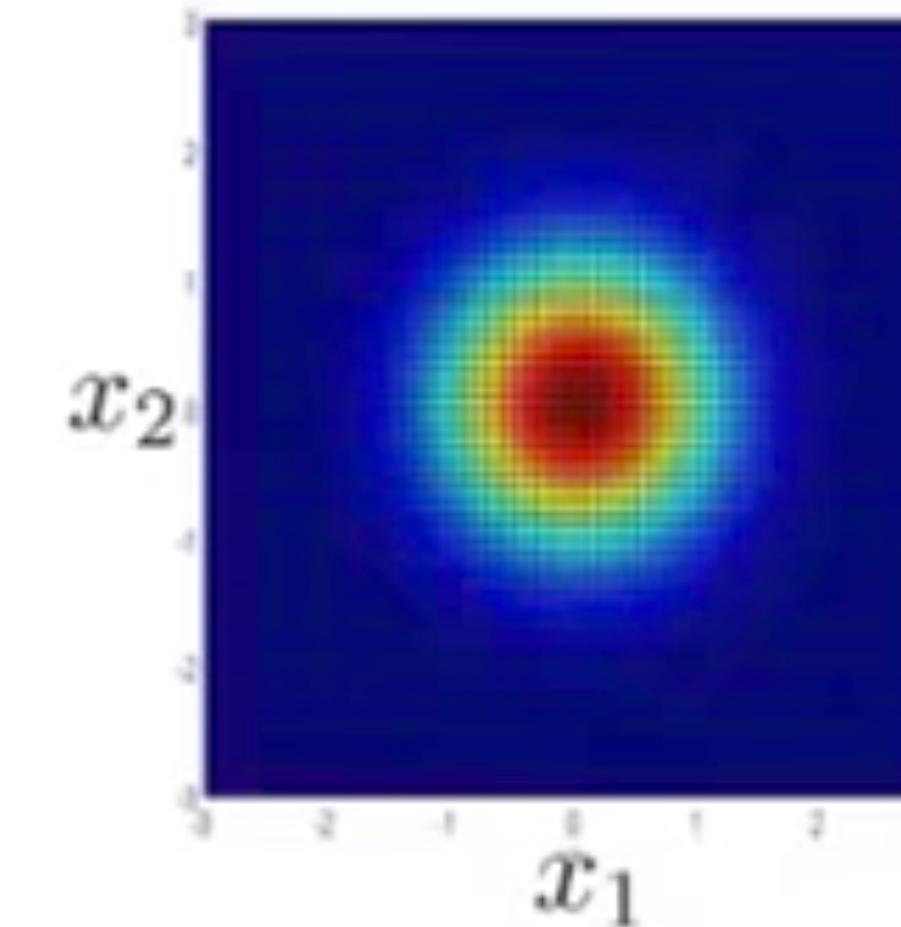
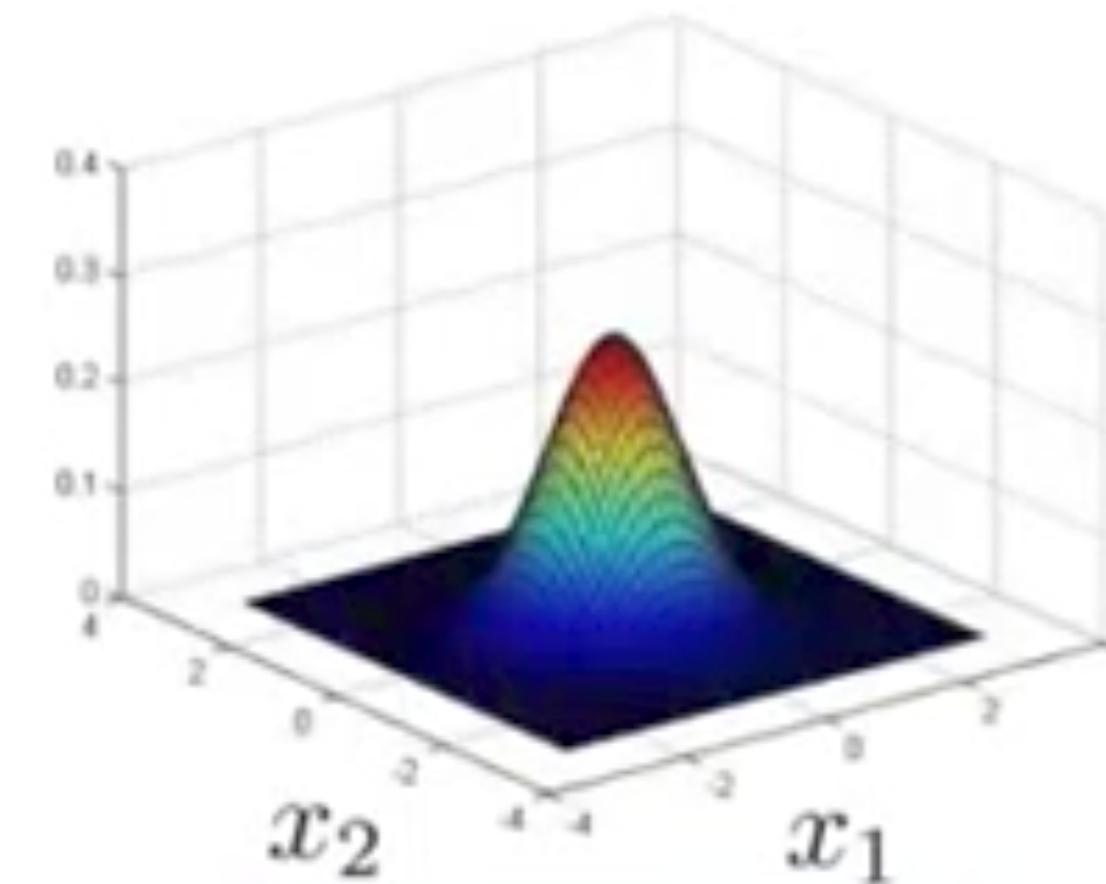
figure credit: Andrew Ng

Multivariate Gaussian Distribution

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



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



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

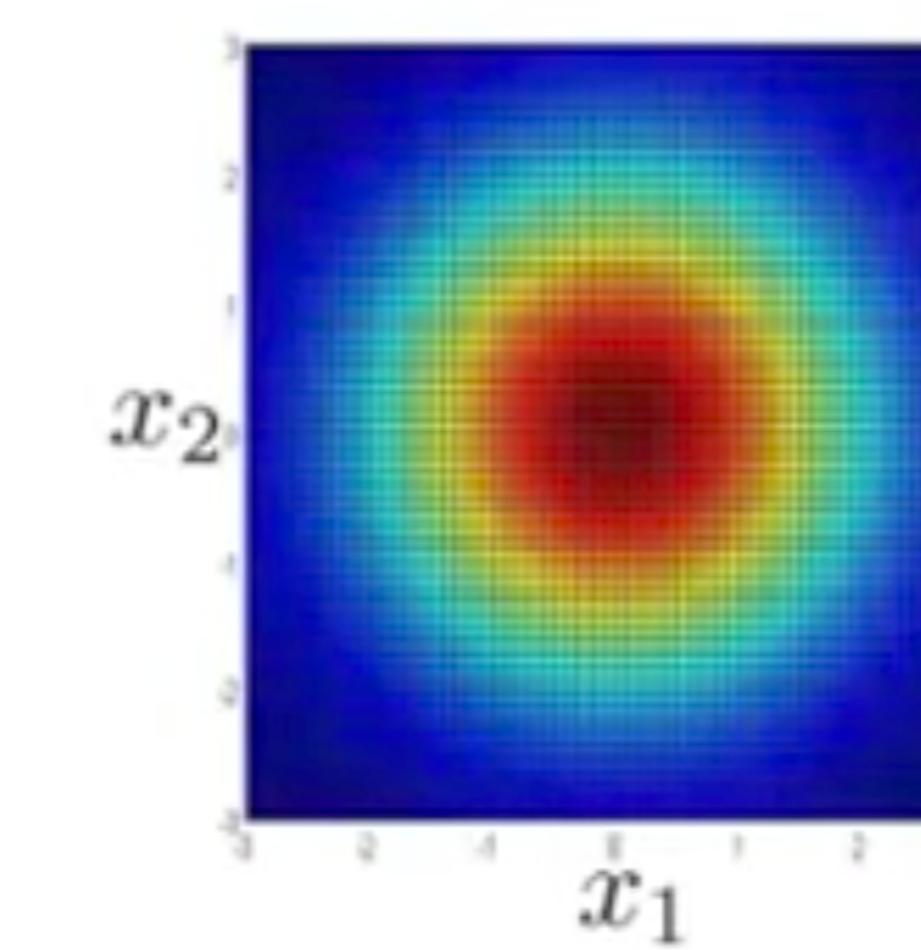
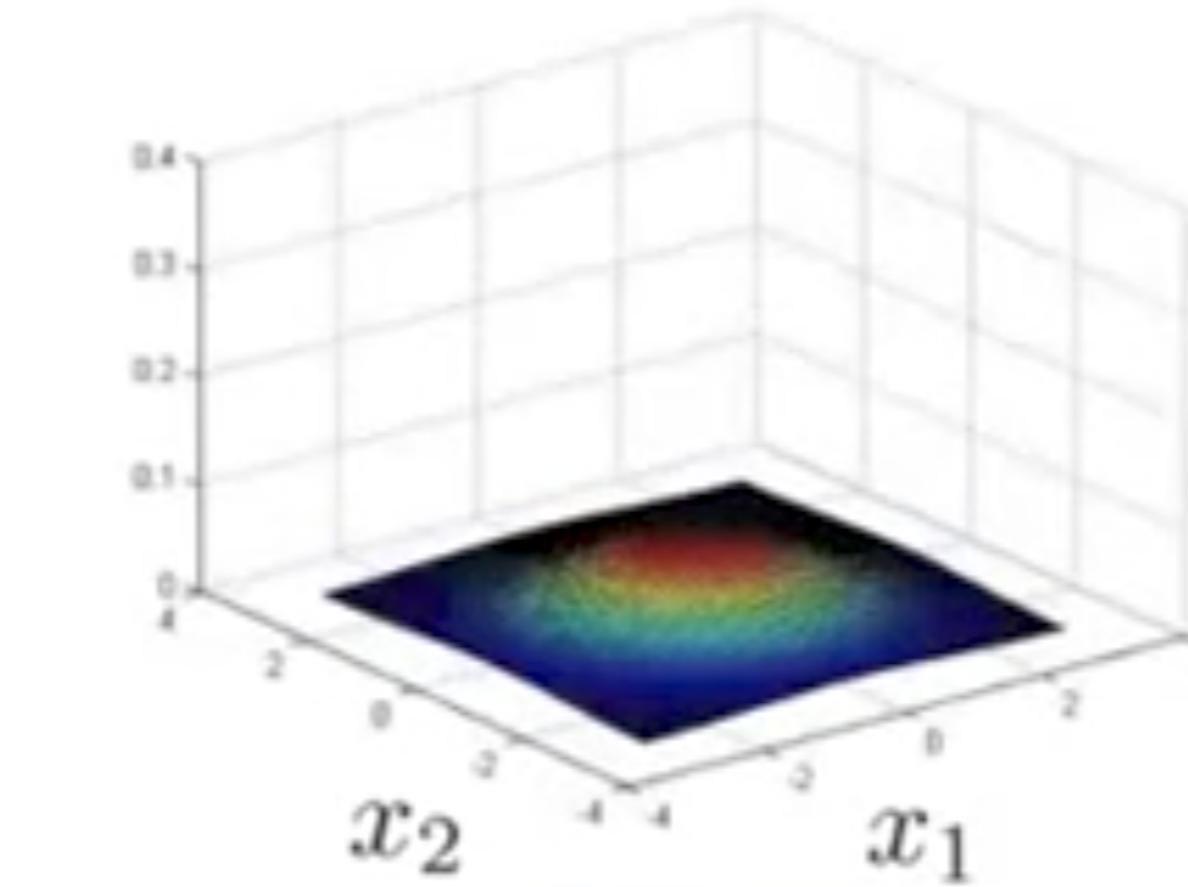
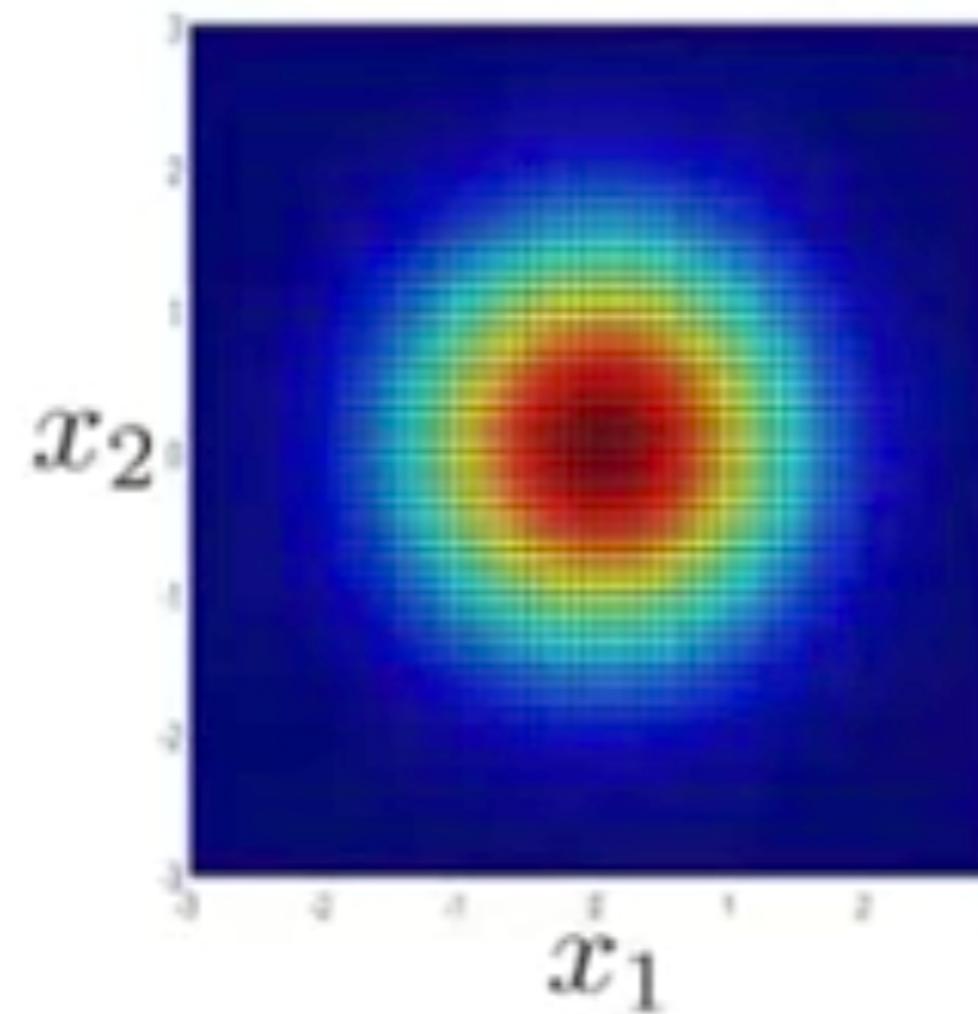
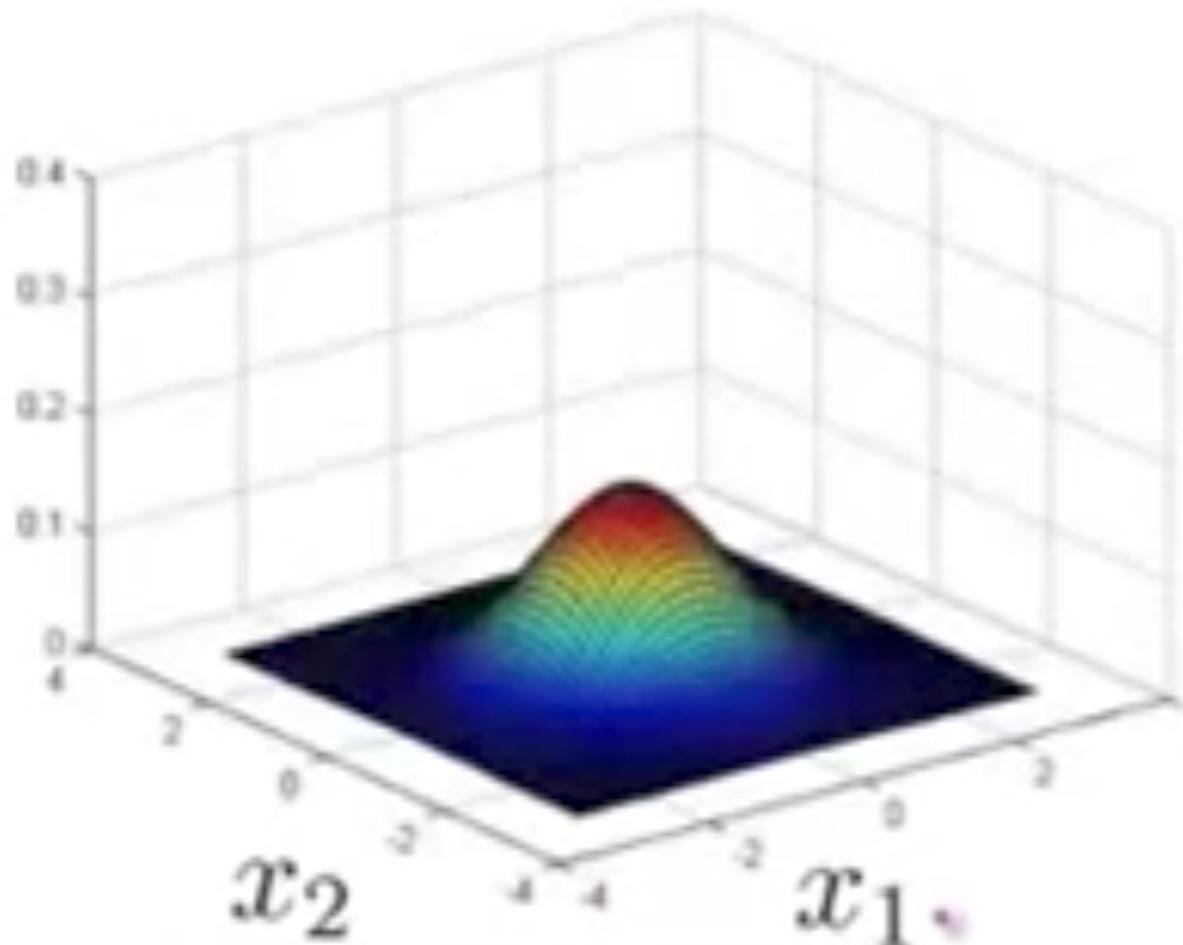


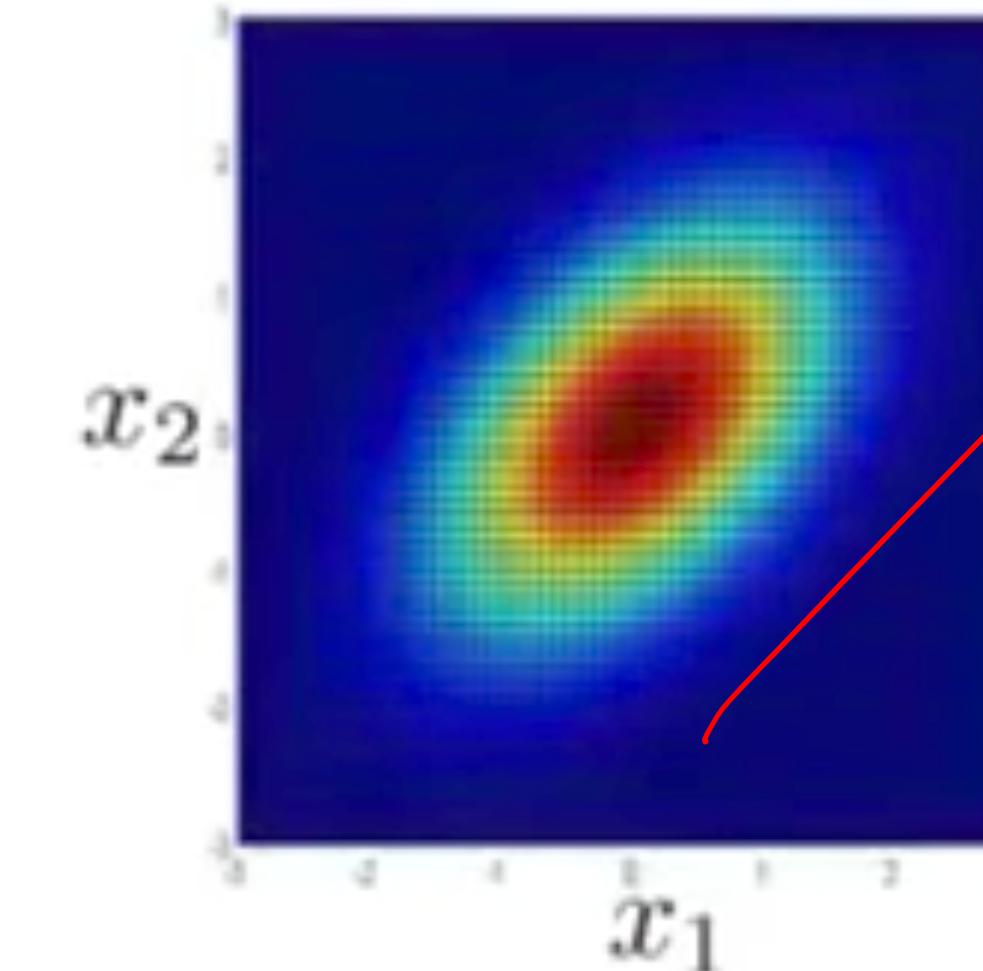
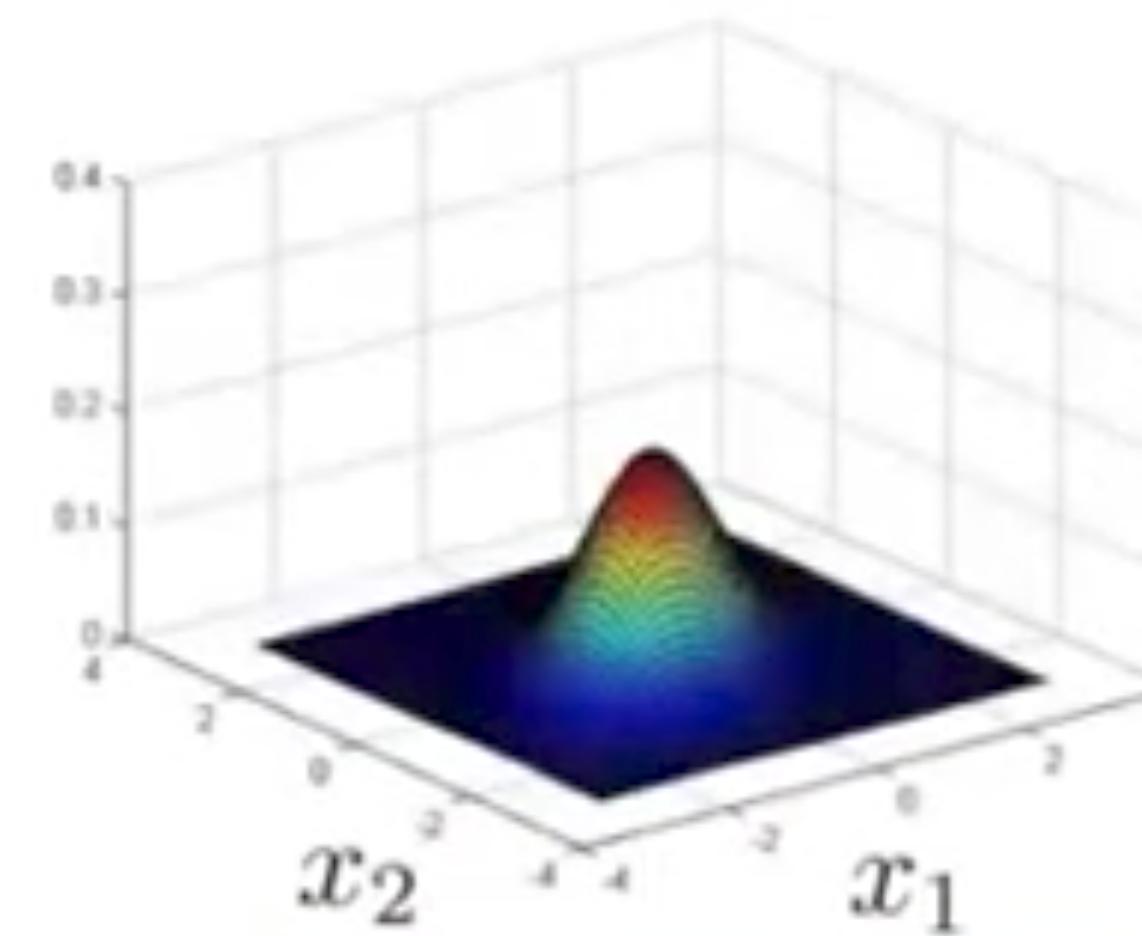
figure credit: Andrew Ng

Multivariate Gaussian Distribution

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



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

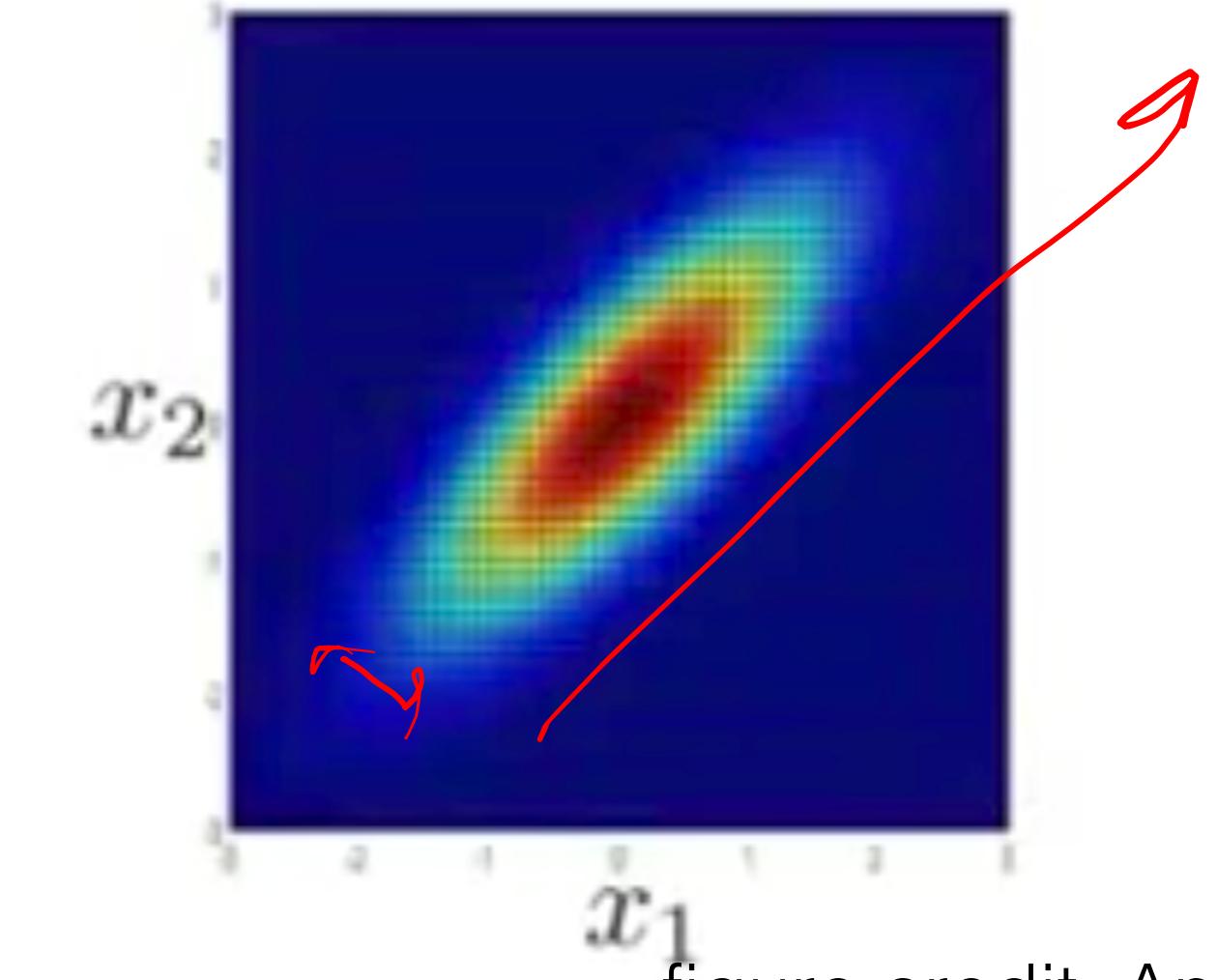
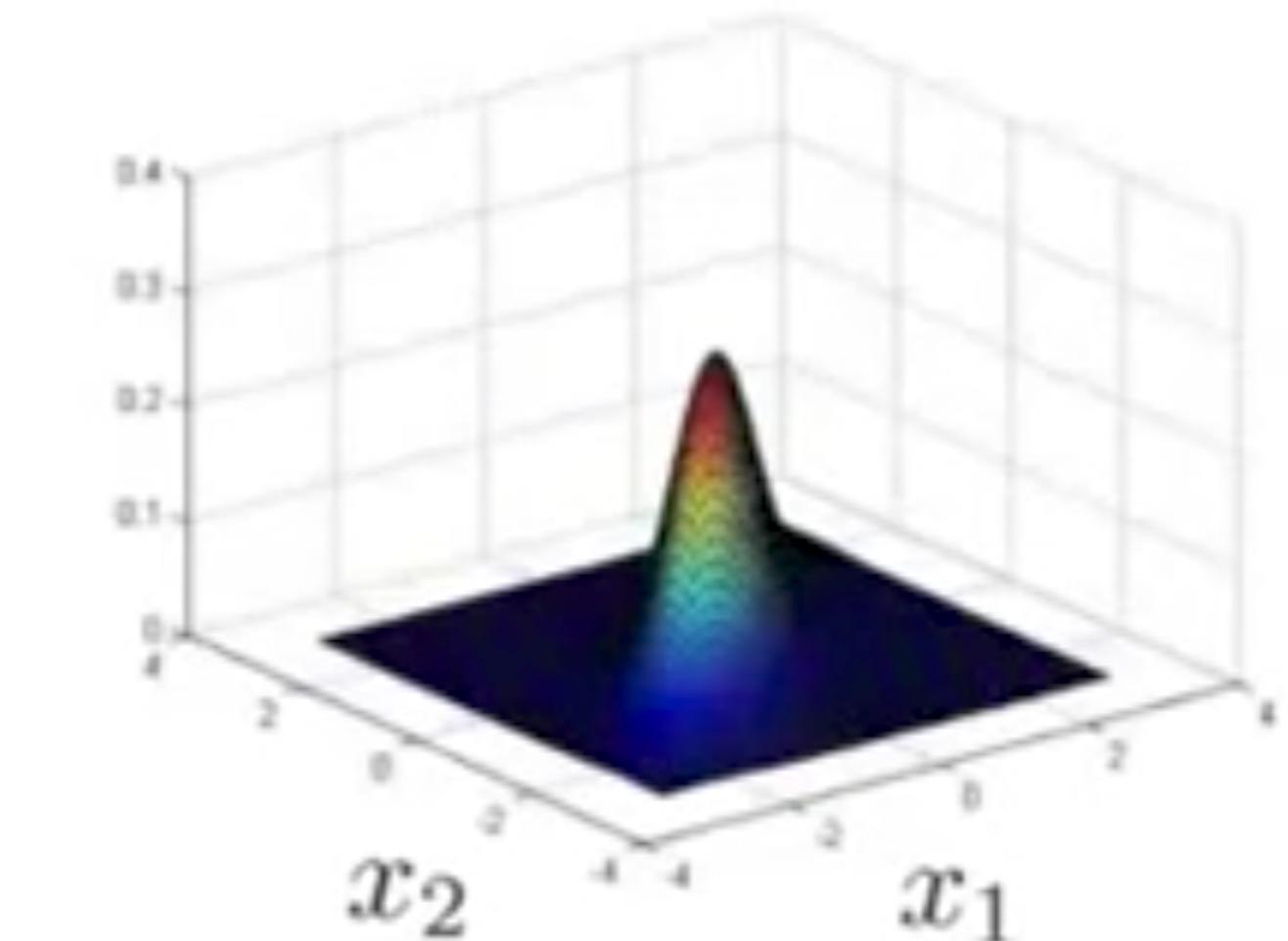
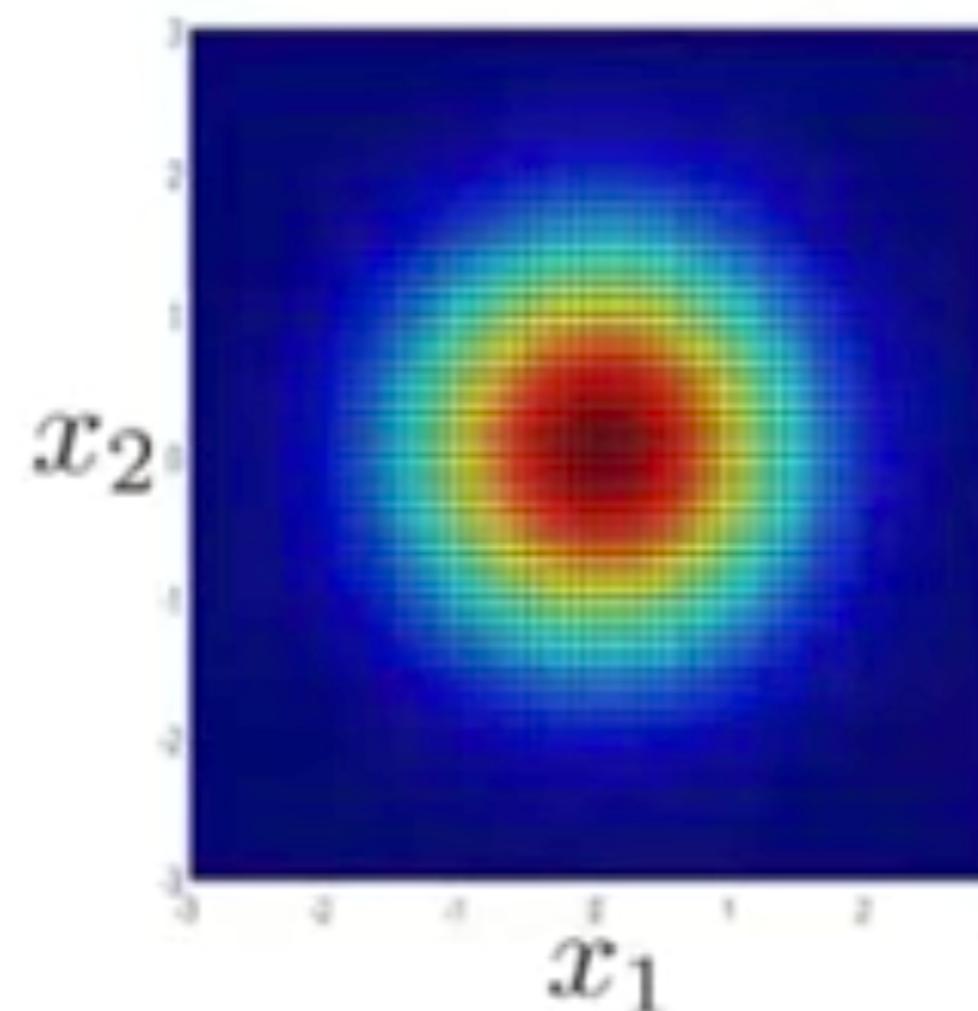
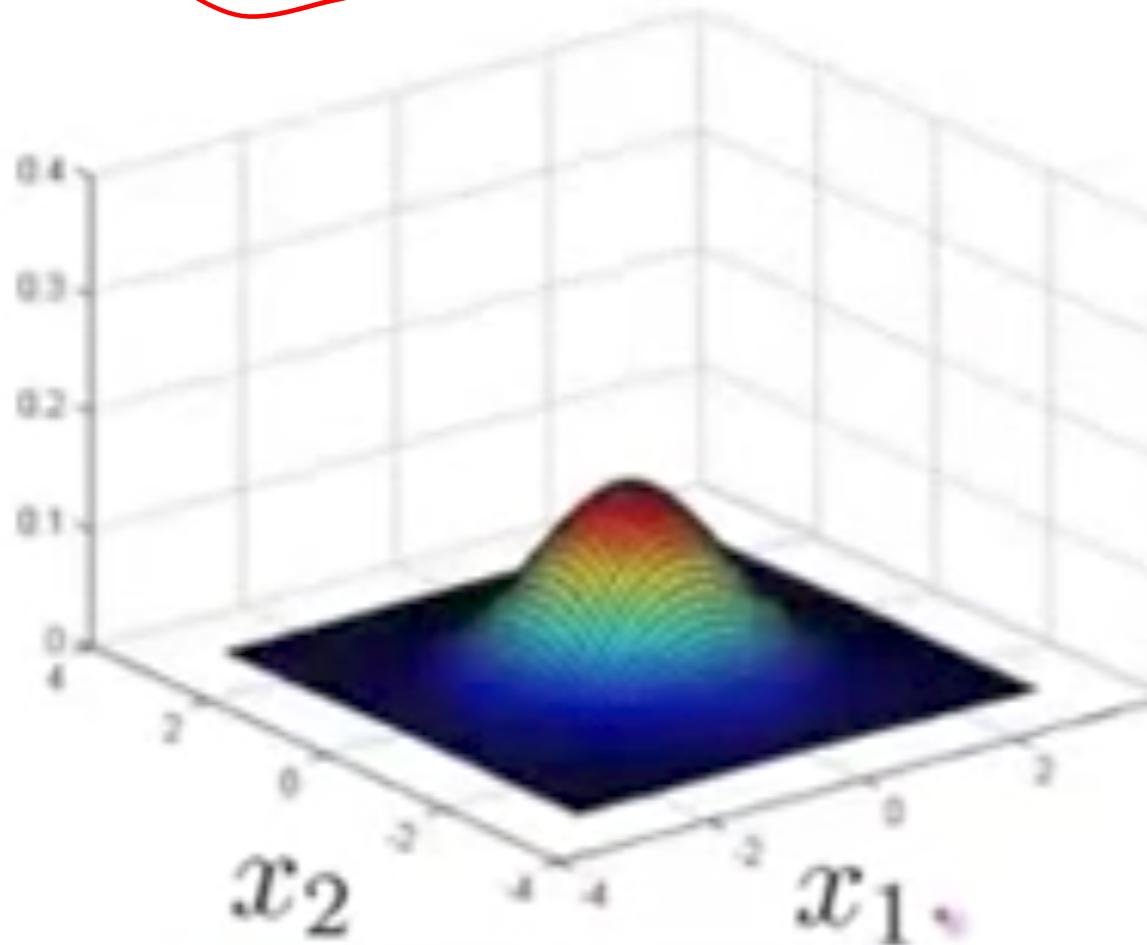


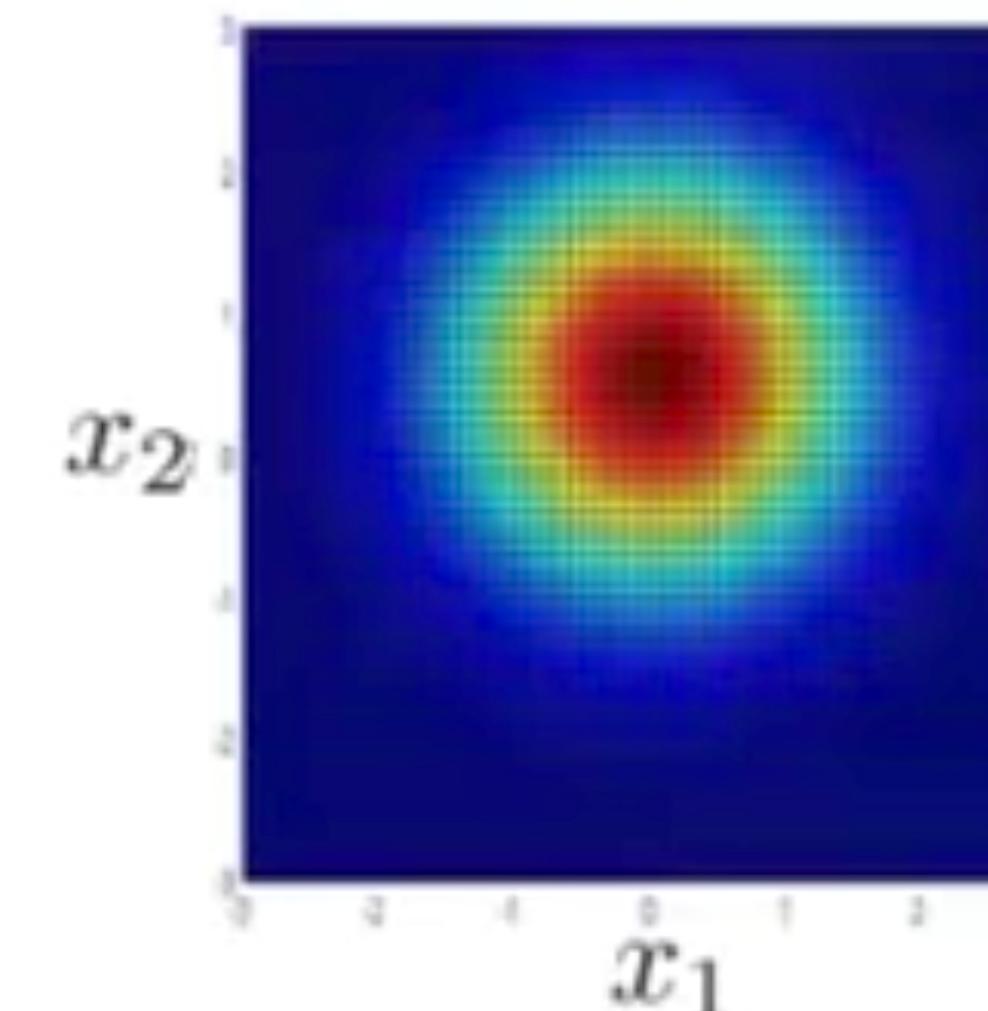
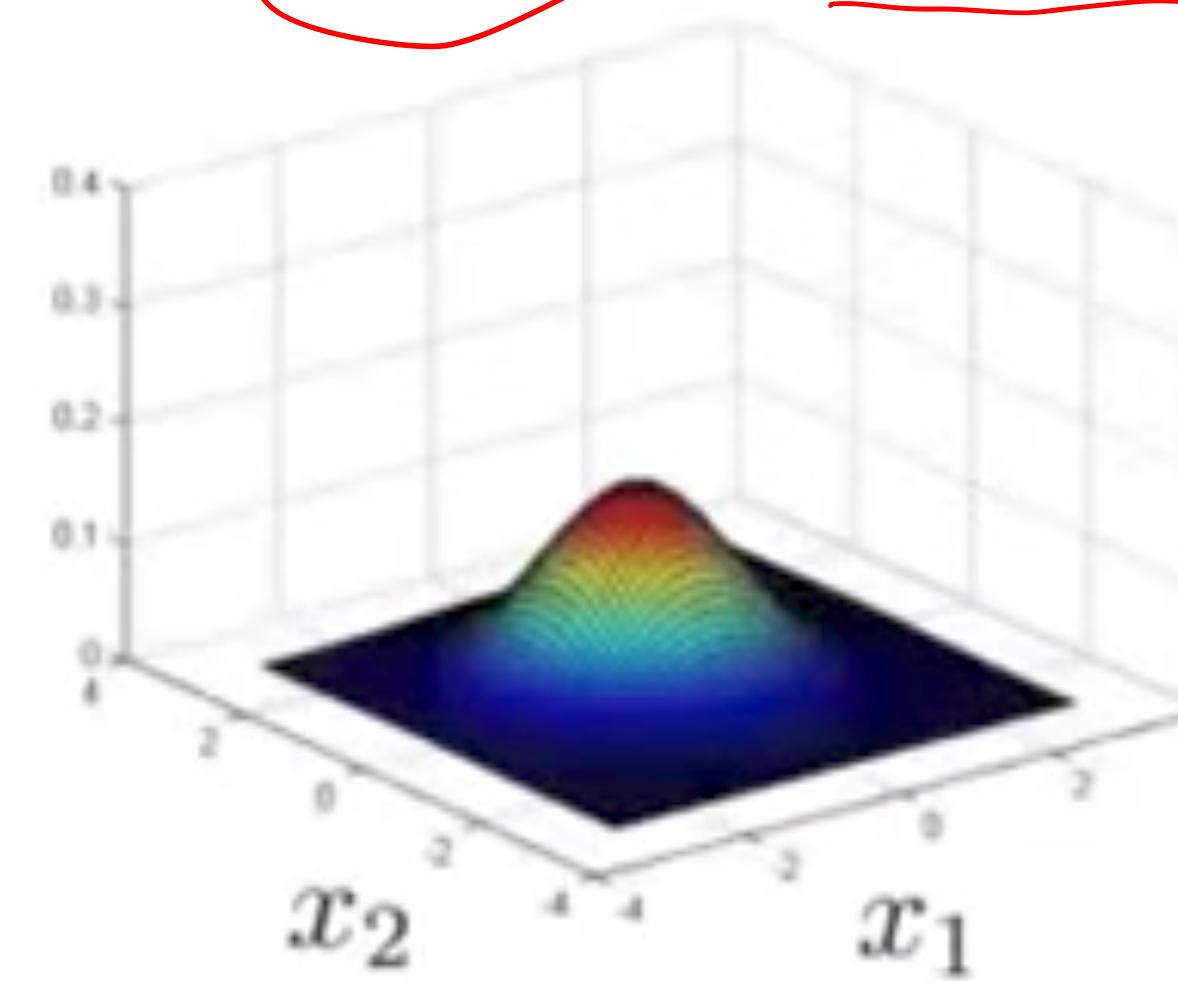
figure credit: Andrew Ng

Multivariate Gaussian Distribution

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



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



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

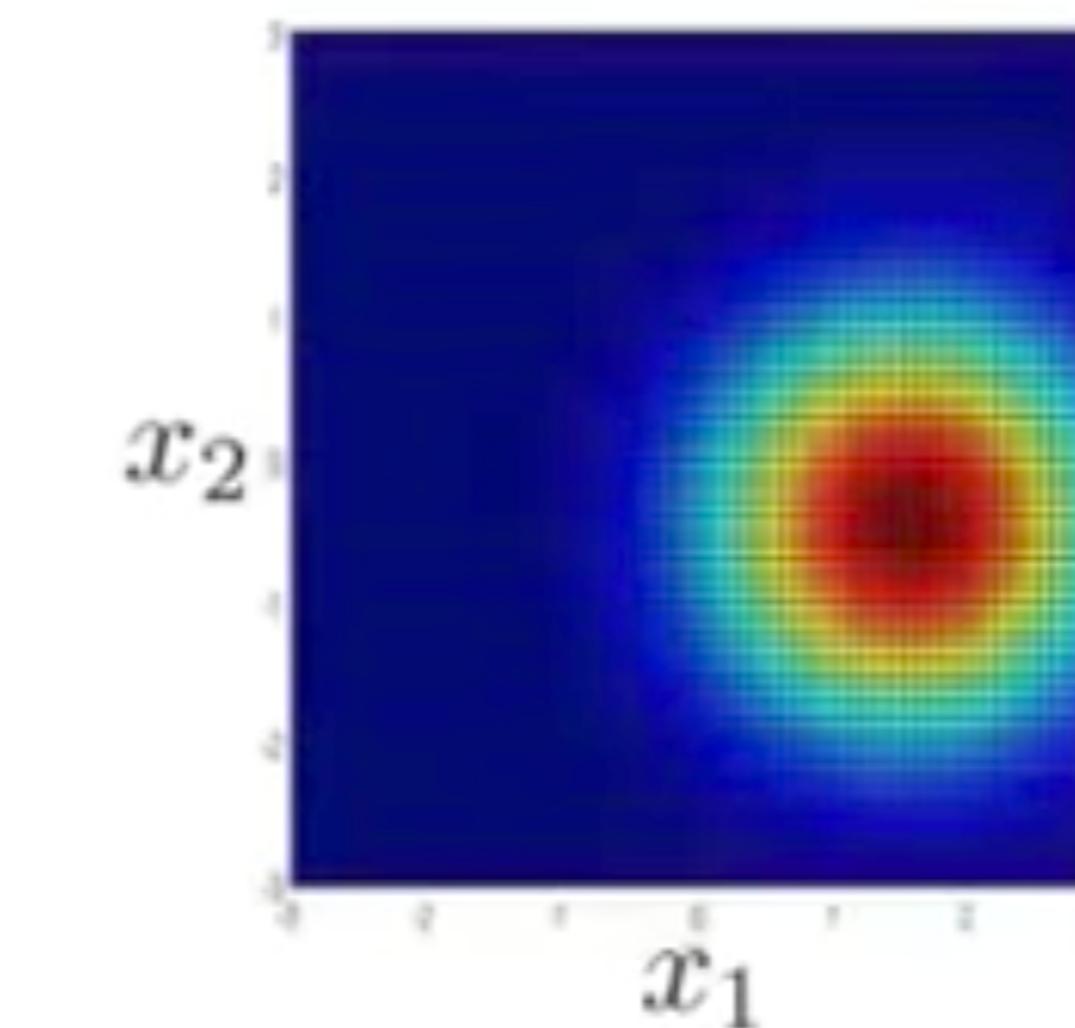
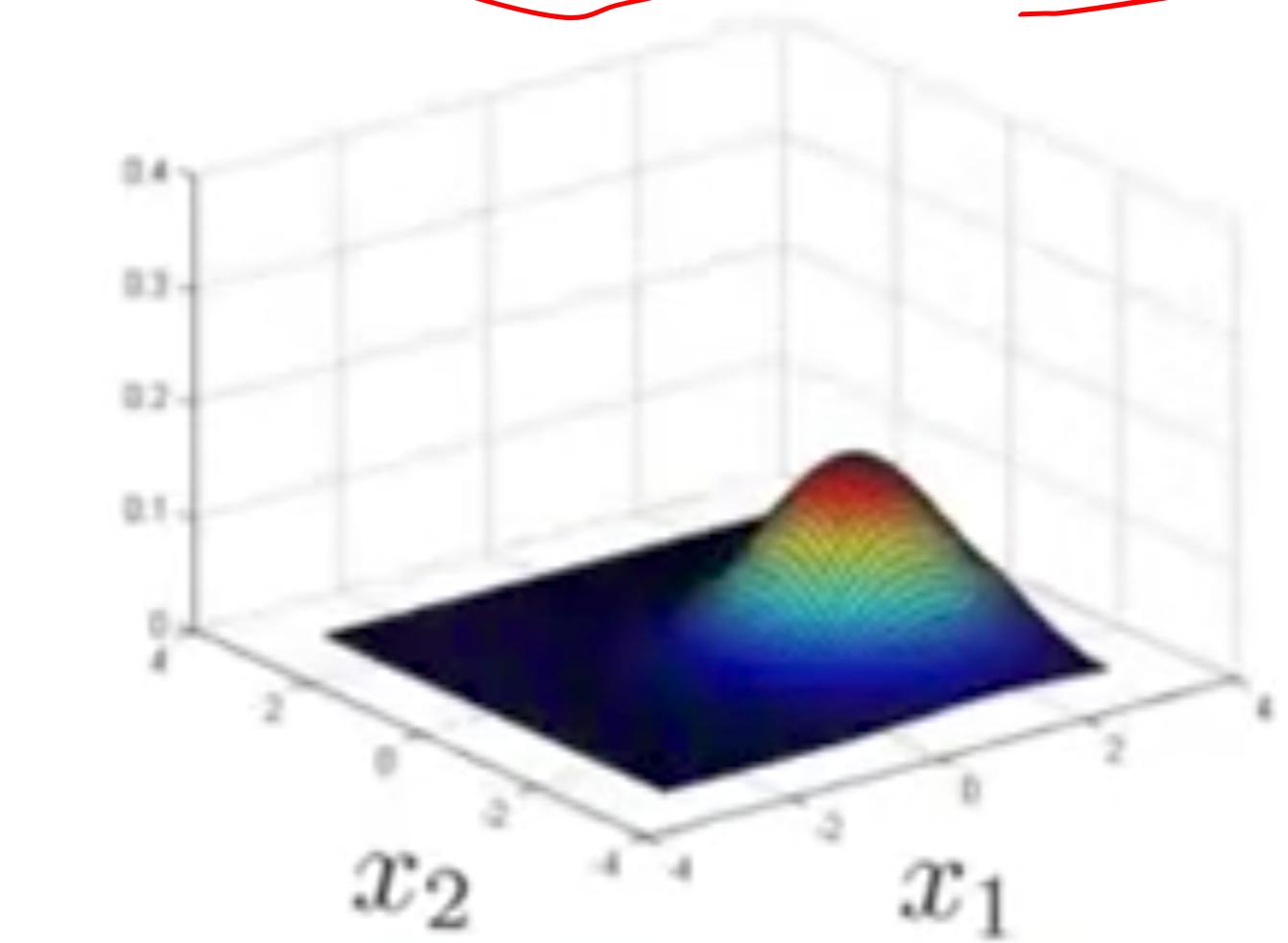


figure credit: Andrew Ng

Anomaly detection: using multivariate Gaussian distribution

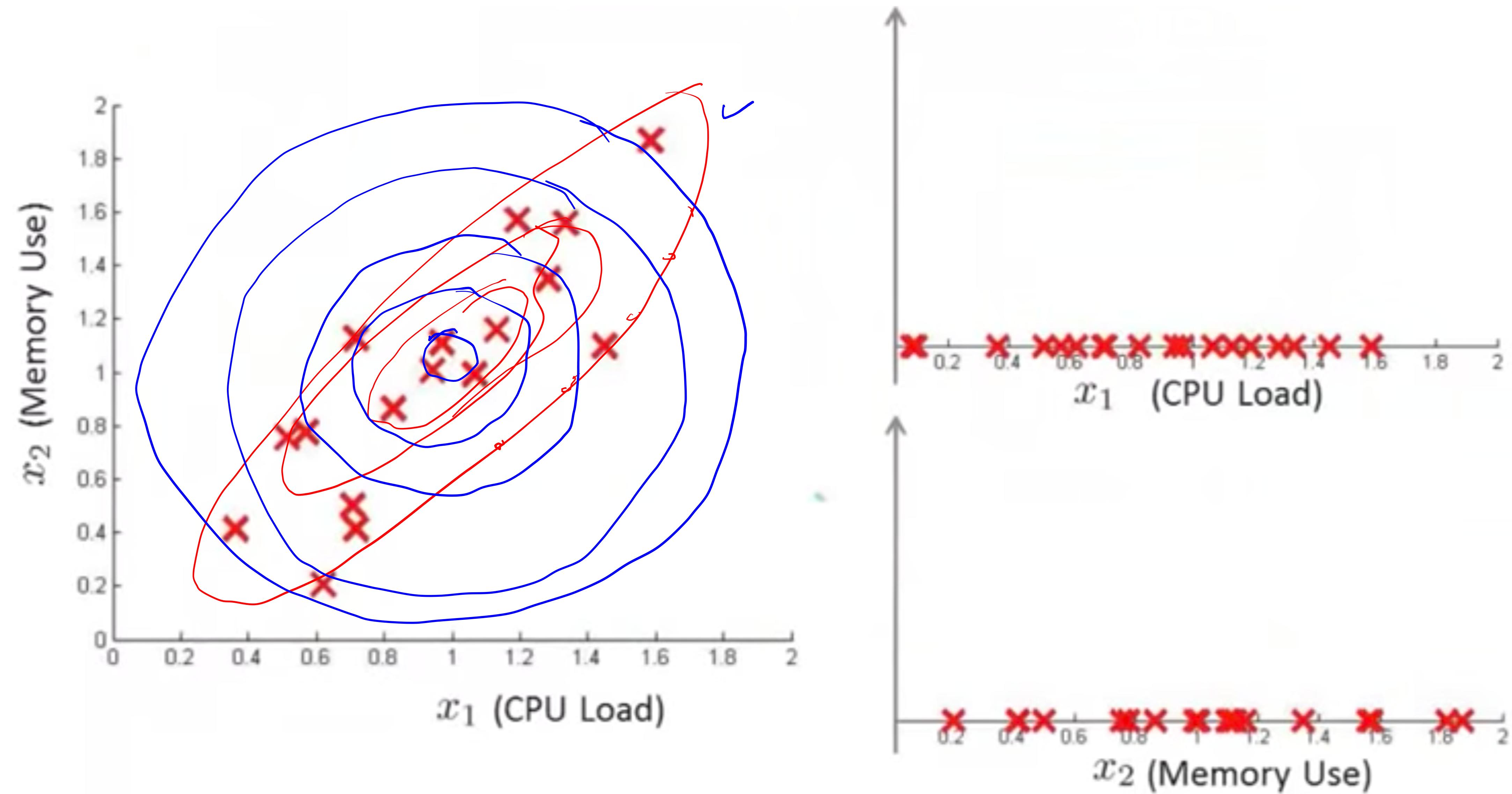


figure credit: Andrew Ng

Anomaly detection

vs.

Supervised learning

- Very small number of positive examples ($y = 1$). (0-20 is common).
- Large number of negative ($y = 0$) examples. $p(x) \leftarrow$
- Many different “types” of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like; future anomalies may look nothing like any of the anomalous examples we’ve seen so far.

Large number of positive and negative examples.

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set.

Bayesian Classification

- Goal: learning functions $f(x) \rightarrow y$
 - $y \rightarrow$ one of the k classes ✓
 - $x = x_1, x_2, \dots, x_d$ - values of attributes (numerical or categorical)
- Probabilistic classification:
 - most probable class given observations: $\hat{y} = \underset{y}{\operatorname{argmax}} P(y | x)$ ✓
- Bayesian probability of a class:

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

Likelihood

Prior

Normalization → ranking

$$P(x) = \sum_{y'} P(x|y')P(y')$$

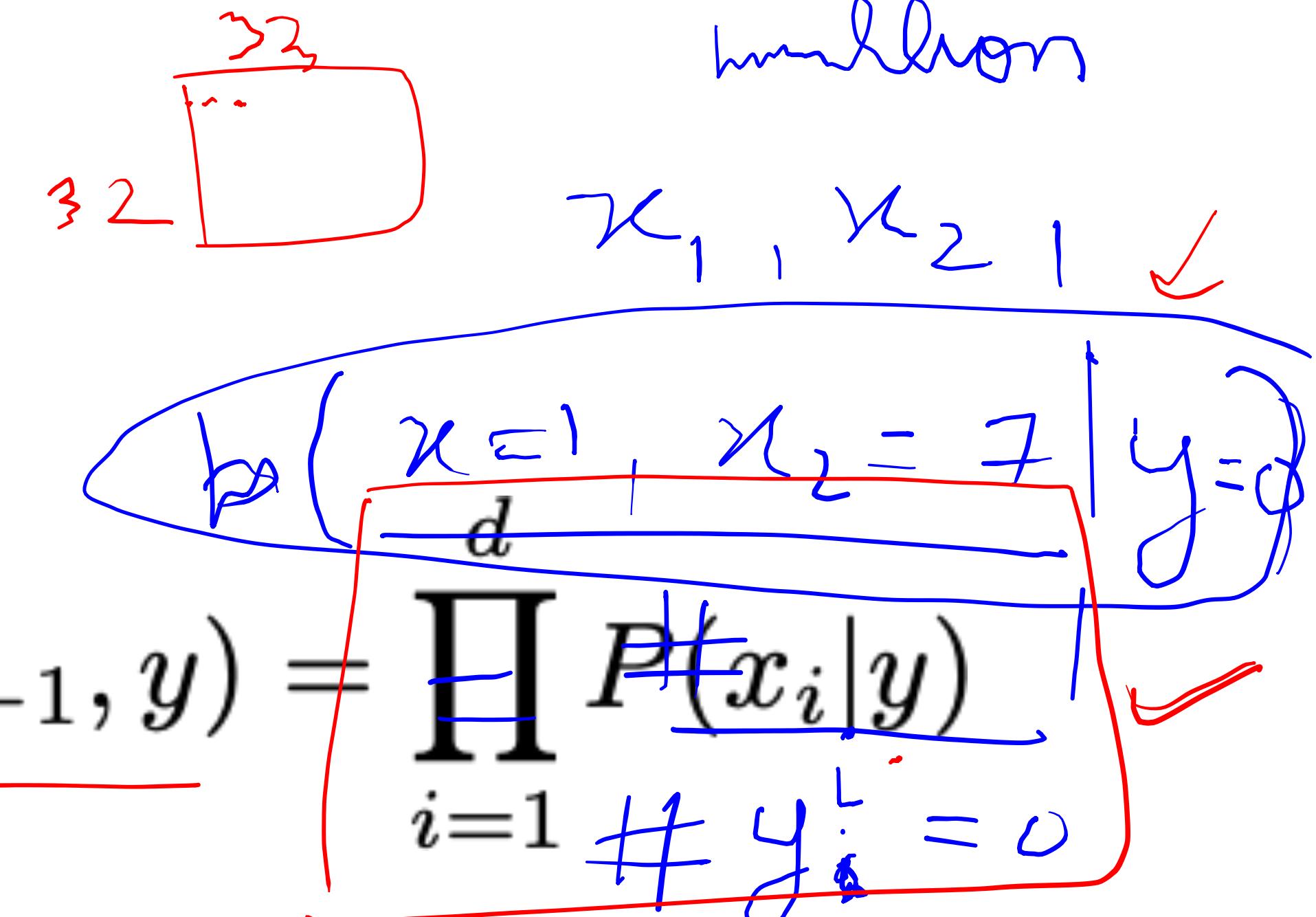
Normalization → ranking

Independence assumption: Naive Bayes

- We need to compute $P(x|y)$ for Bayesian classification
- x is multidimensional, which makes counting like easy probability estimation difficult

Idea: Make independence assumption

$$P(x_1, \dots, x_d | y) = \prod_{i=1}^d P(x_i | x_1, \dots, x_{i-1}, y)$$



assuming that the class value explains all the dependence between attributes

Continuous Example

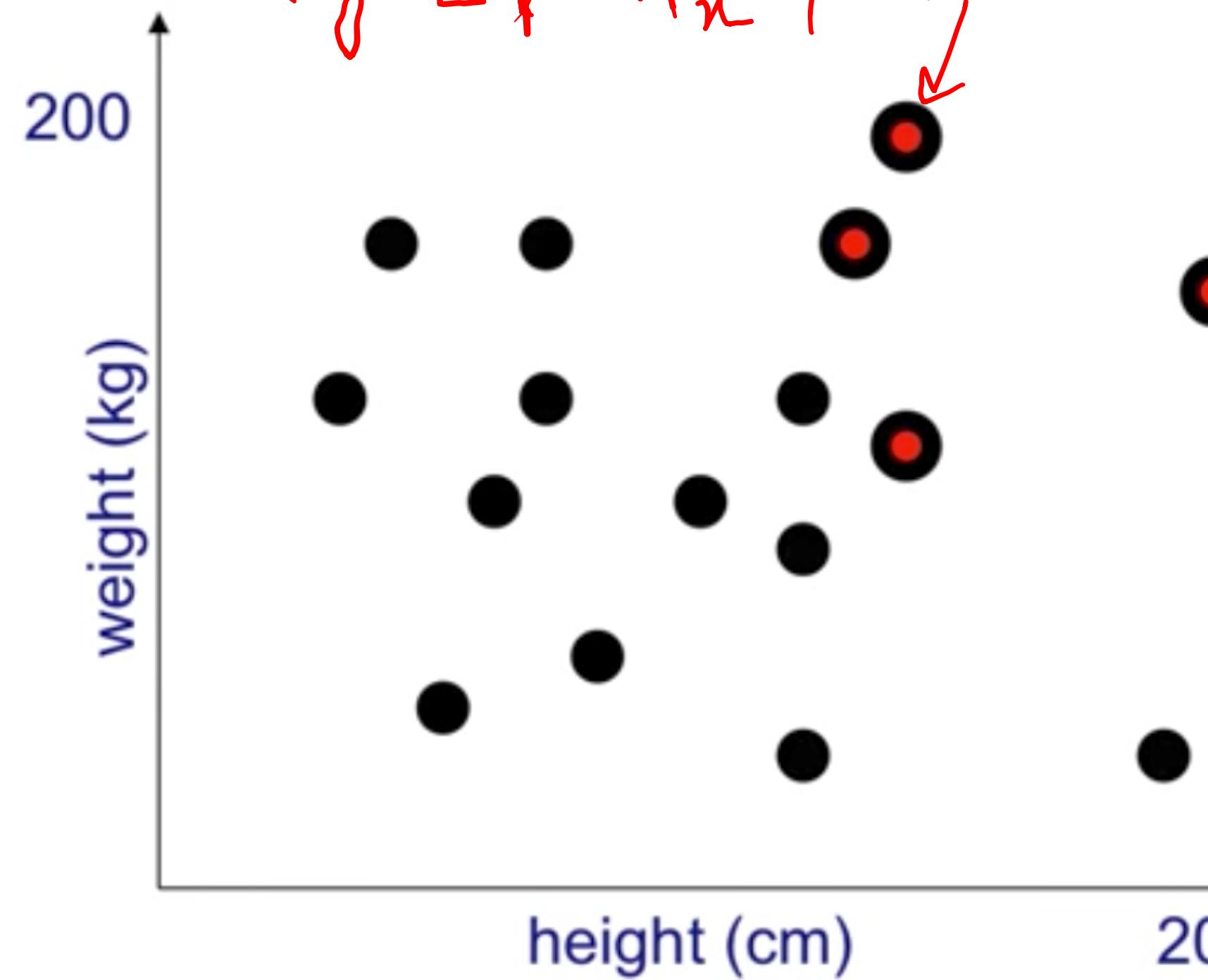
$$\begin{pmatrix} \nu_1 \\ \nu_2 \end{pmatrix}$$

- Distinguish children from adults based on height and weight
 - classes: {a,c} , attributes: height (cm), weight (kg)
 - training examples: $\{h_i, w_i, y_i\}$, 4 adults and 12 children
- Prior probabilities: $P(a) = \underline{4/16}$ and $P(c) = \underline{12/16}$
- Model: gaussians using mean and variances

$$p(y=a | h_x, w_x)$$

$$(h_x, w_x)$$

$$p(y=c | h_x, w_x)$$



Continuous Example

$$p(a) = \frac{4}{7+1} = \frac{4}{11}$$

$$p(b) = \frac{12}{15}$$

$$p(h_x, w_x | y=a) = p(h_x | \mu_{h,a}, \sigma_{h,a}) \cdot p(w_x | \mu_{w,a}, \sigma_{w,a})$$

$$p(h_x | a) = \frac{1}{\sqrt{2\pi} \sigma_{h,a}} \cdot e^{-\frac{1}{2} \frac{(h_x - \mu_{h,a})^2}{\sigma_{h,a}^2}}$$

$$p(w_x | a) = \frac{1}{\sqrt{2\pi} \sigma_{w,a}} \cdot e^{-\frac{1}{2} \frac{(w_x - \mu_{w,a})^2}{\sigma_{w,a}^2}}$$

$$p(y=a | h_x, w_x) = p(h_x, w_x | y=a) p(y=a)$$

Decision boundary

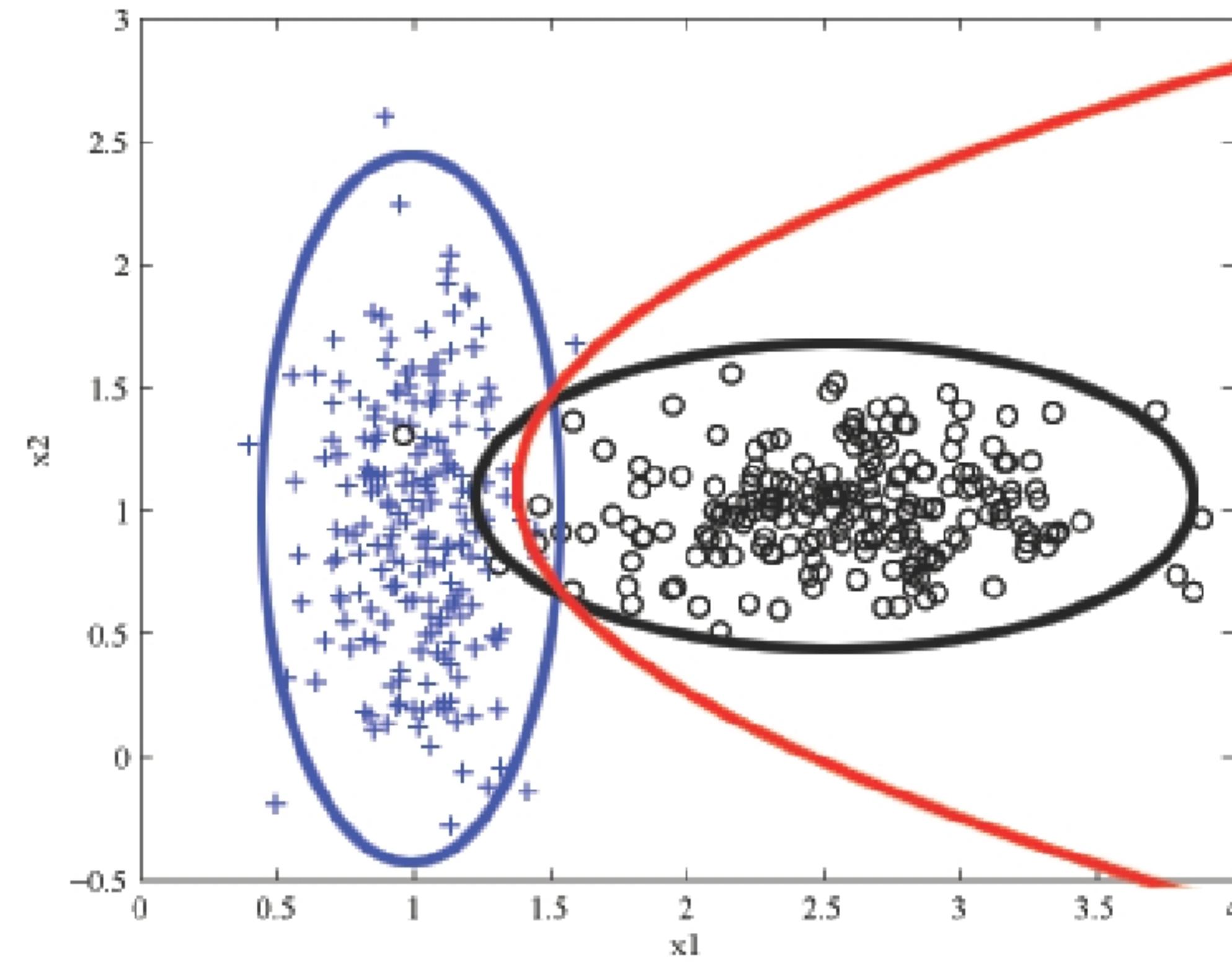
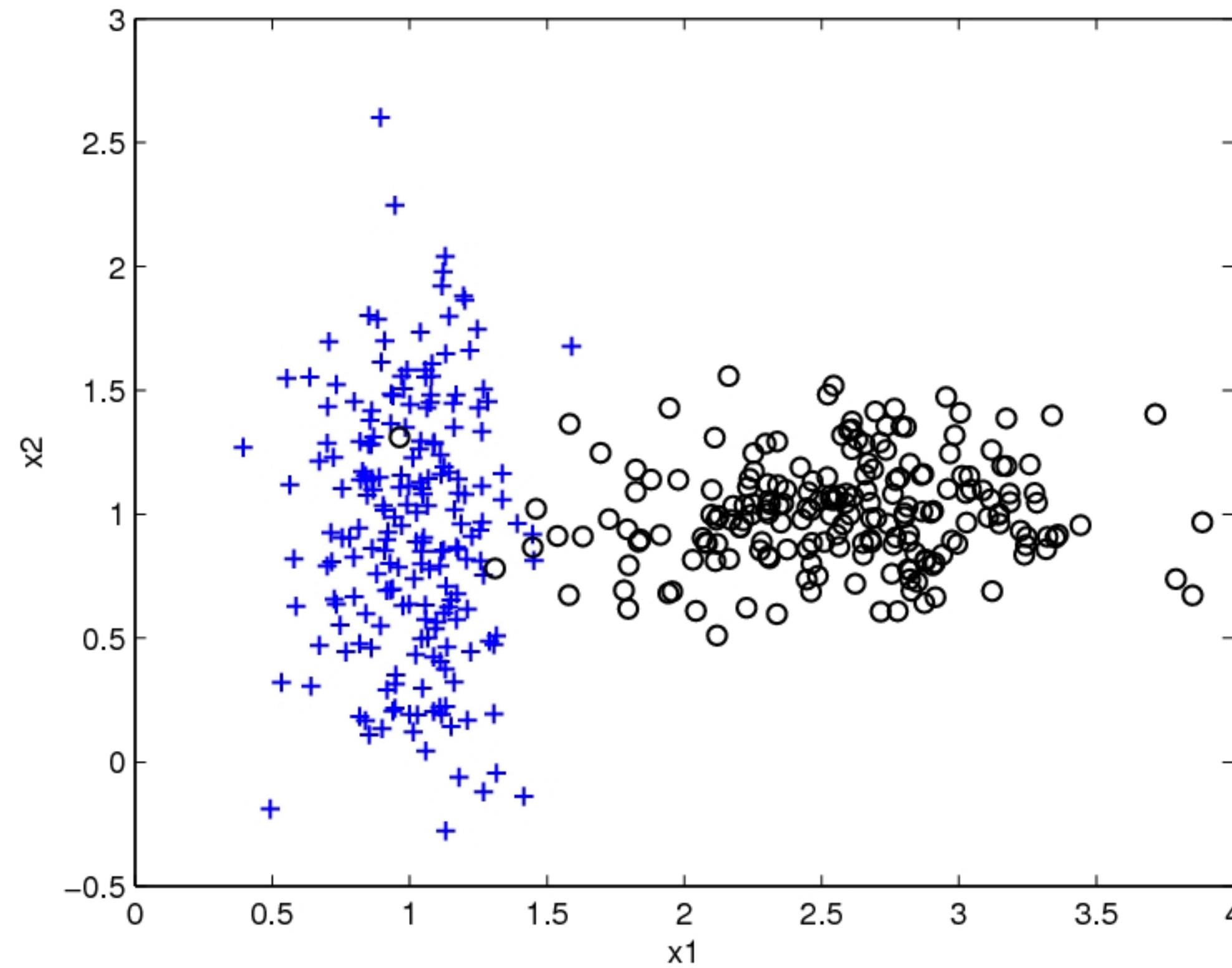


figure credit: Ben Tasker

Problem with Naive Bayes

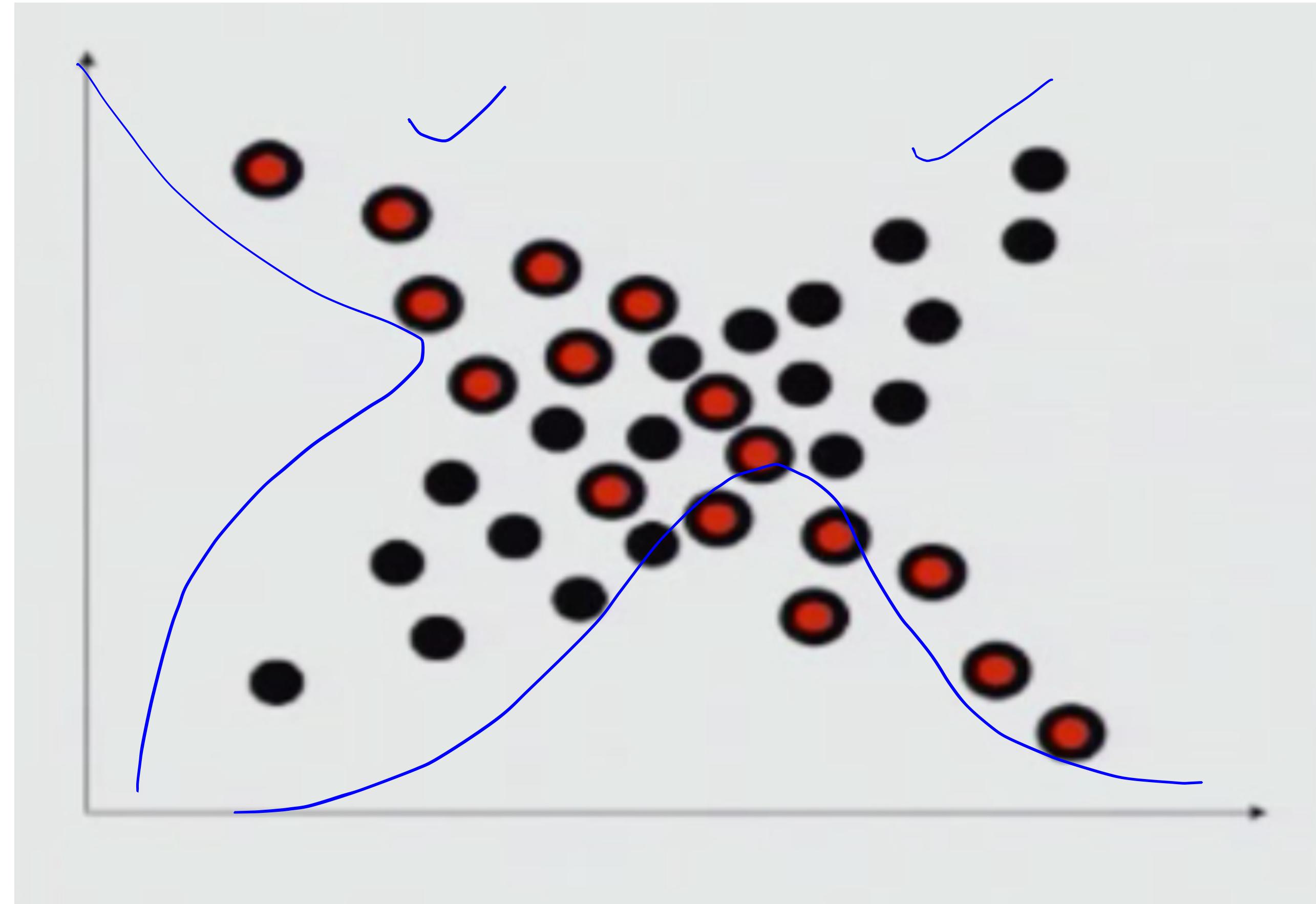


figure credit: Victor Lavrenko

$p(\text{spam} \mid \text{review us})$

Discrete example

$p(\sigma \mid \text{ham}) = 2/2$

$= p(\text{review us} \mid \text{spam}) p(\text{spam})$

$P(\text{spam}) = 4/6 \quad P(\text{ham}) = 2/6$

- D1: "send us your password" - spam
- D2: "send us your review" - ham
- D3: "review password" - ham
- D4: "review us" - spam
- D5: "send your password" - spam
- D6: "send us your account" - spam

	spam	ham	
spam	2/4	1/2	password $\sim \kappa_1$
ham	1/4	2/2	review $\sim \kappa_2$
send	3/4	1/2	send
us	3/4	1/2	us
your	3/4	1/2	your
account	1/4	0/2	account $\sim \kappa_3$

\times new email: "review us now"

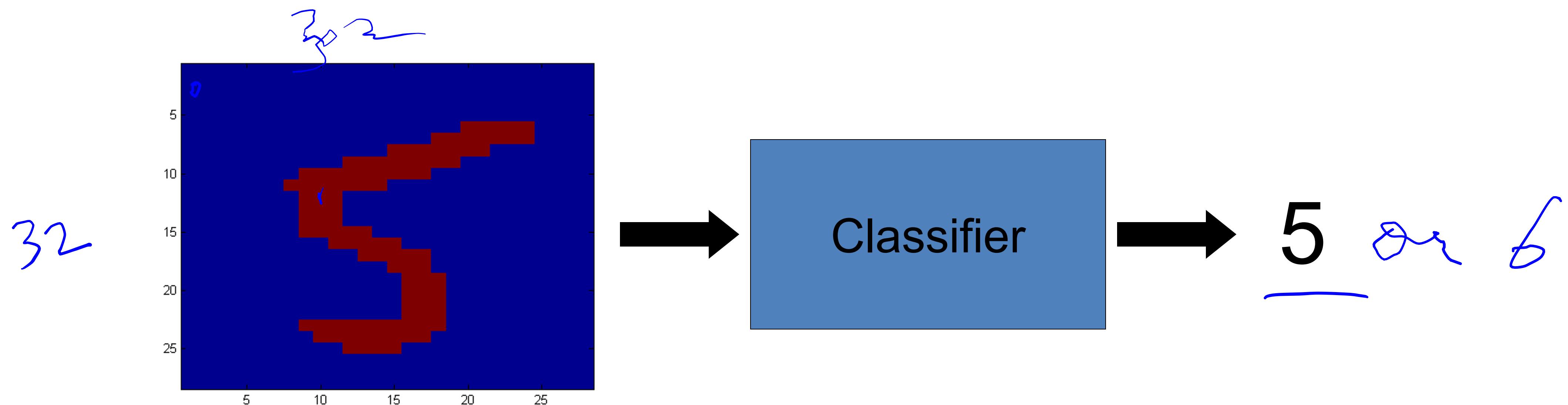
$2^{50,000}$

$p(\text{review us} \mid \text{spam}) = p(0, 1, 0, 1, 0, 0 \mid \text{spam})$

$= \left(1 - \frac{2}{4}\right) \left(\frac{1}{4}\right) \left(1 - \frac{3}{4}\right) \left(\frac{3}{4}\right) \left(1 - \frac{3}{4}\right) \left(1 - \frac{1}{4}\right)$

Another Application

- **Digit Recognition**



- $X_1, \dots, X_n \in \underline{\{0,1\}}$ (Black vs. White pixels)
- $Y \in \{5,6\}$ (predict whether a digit is a 5 or a 6)

The Bayes Classifier

- A good strategy is to predict:

$$\arg \max_Y P(Y=5 | X_1, \dots, X_n)$$

- (for example: what is the probability that the 32×32 image represents a 5 given its pixels?)
- How do we compute that?

The Bayes Classifier

- Use Bayes Rule!

$$P(Y|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y)P(Y)}{P(X_1, \dots, X_n)}$$

Likelihood

Prior

Normalization Constant

- Why did this help? Well, we think that we might be able to specify how features are “generated” by the class label

The Bayes Classifier

- Let's expand this for our digit recognition task:

$$P(Y = 5|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 5)P(Y = 5)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$
$$P(Y = 6|X_1, \dots, X_n) = \frac{P(X_1, \dots, X_n|Y = 6)P(Y = 6)}{P(X_1, \dots, X_n|Y = 5)P(Y = 5) + P(X_1, \dots, X_n|Y = 6)P(Y = 6)}$$

- To classify, we'll simply compute these two probabilities and predict based on which one is greater

2 ^{32 × 32}

Model Parameters

- For the Bayes classifier, we need to “learn” two functions, the likelihood and the prior
- How many parameters are required to specify the prior for our digit recognition example?

Model Parameters

- How many parameters are required to specify the likelihood?
 - (Supposing that each image is 30x30 pixels)

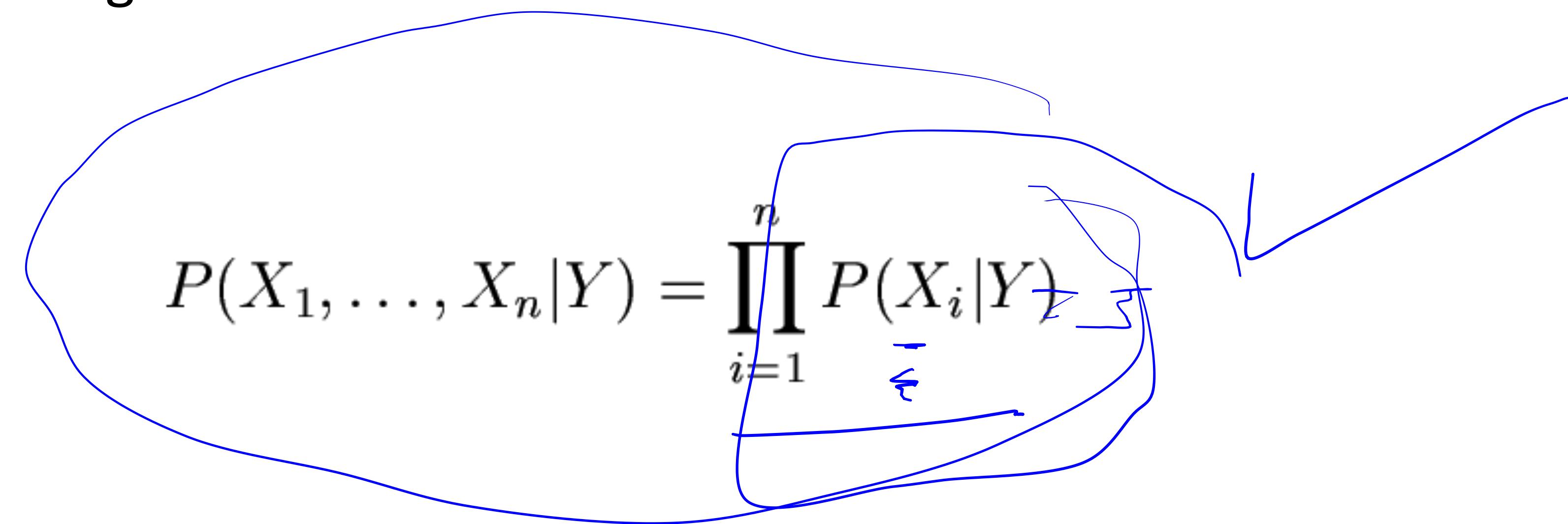
?

Model Parameters

- The problem with explicitly modeling $P(X_1, \dots, X_n | Y)$ is that there are usually way too many parameters:
 - We'll run out of space
 - We'll run out of time
 - And we'll need tons of training data (which is usually not available)

The Naïve Bayes Model

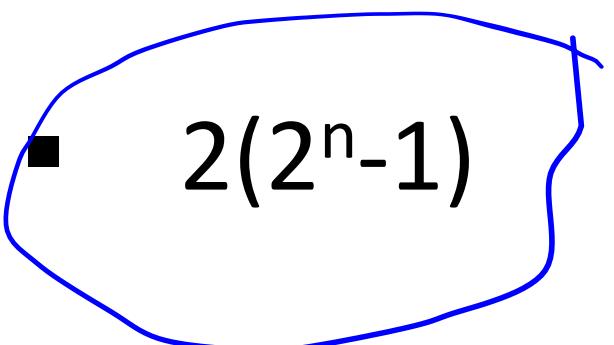
- The *Naïve Bayes Assumption*: Assume that all features are independent **given the class label Y**
- Equationally speaking:

$$P(X_1, \dots, X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$


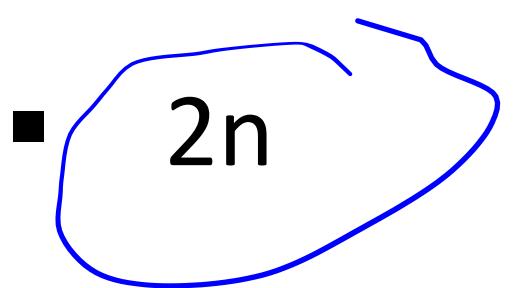
- (We will discuss the validity of this assumption later)

Why is this useful?

- # of parameters for modeling $P(X_1, \dots, X_n | Y)$:

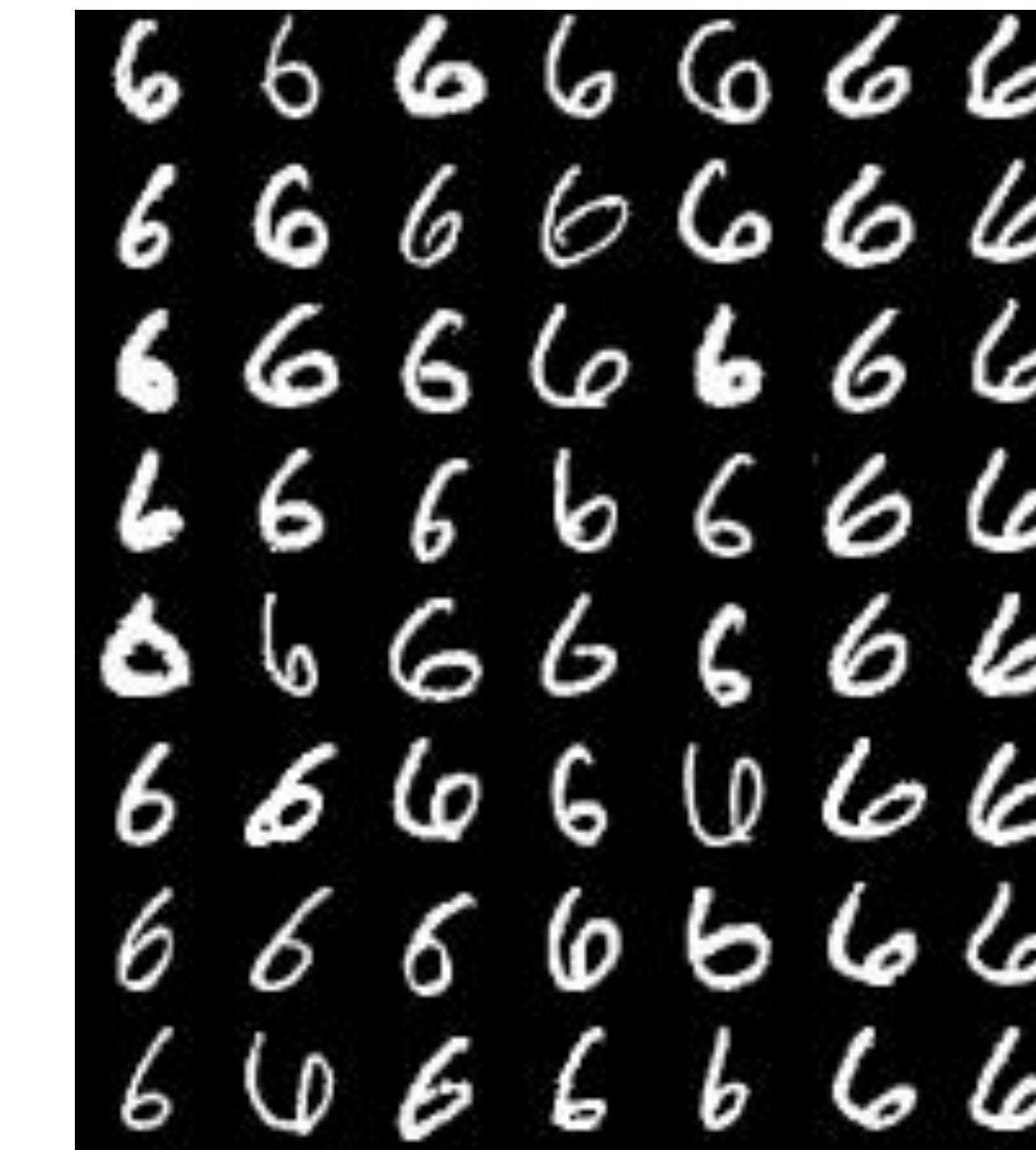
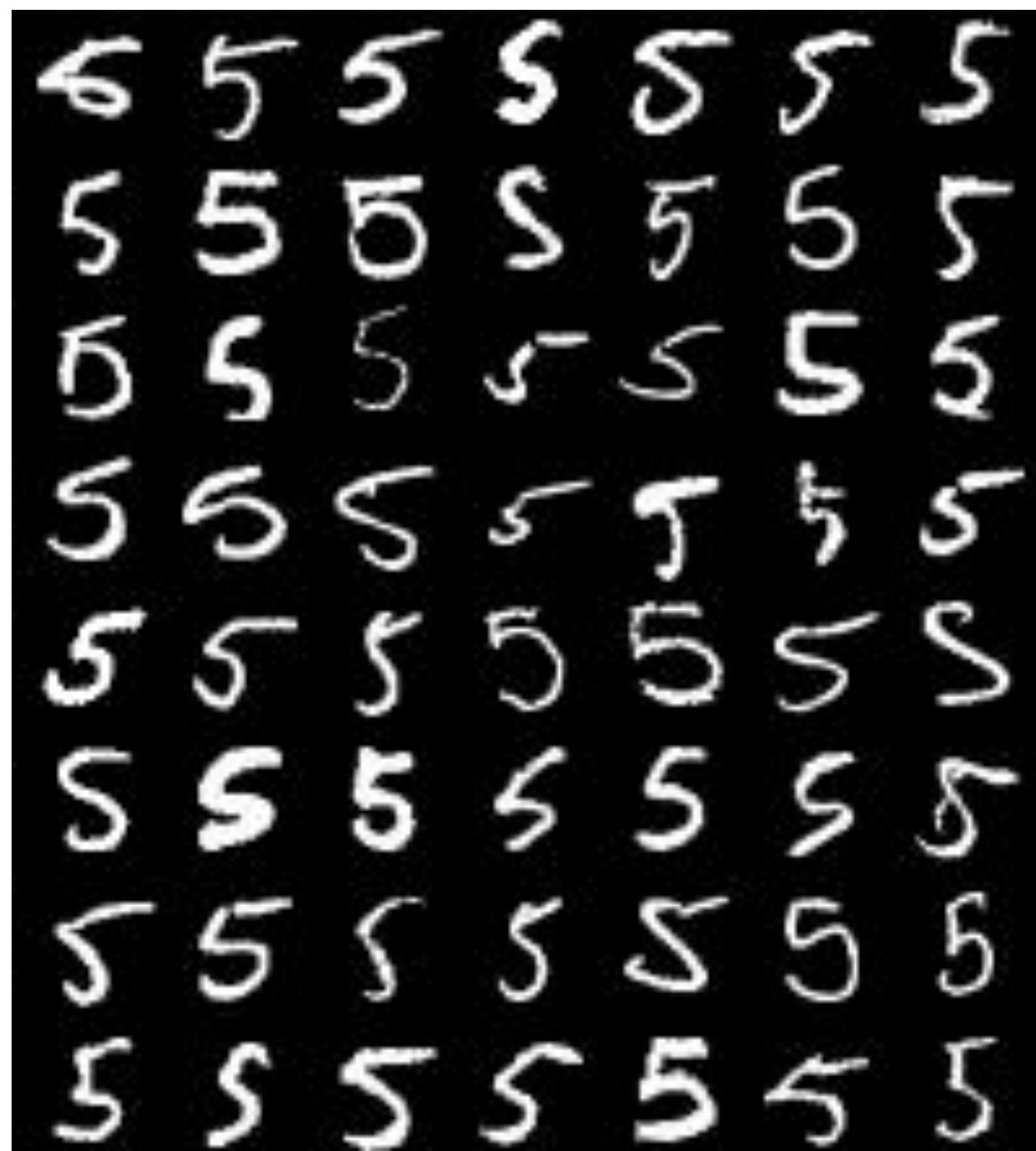

$$2(2^n - 1)$$

- # of parameters for modeling $P(X_1 | Y), \dots, P(X_n | Y)$


$$2n$$

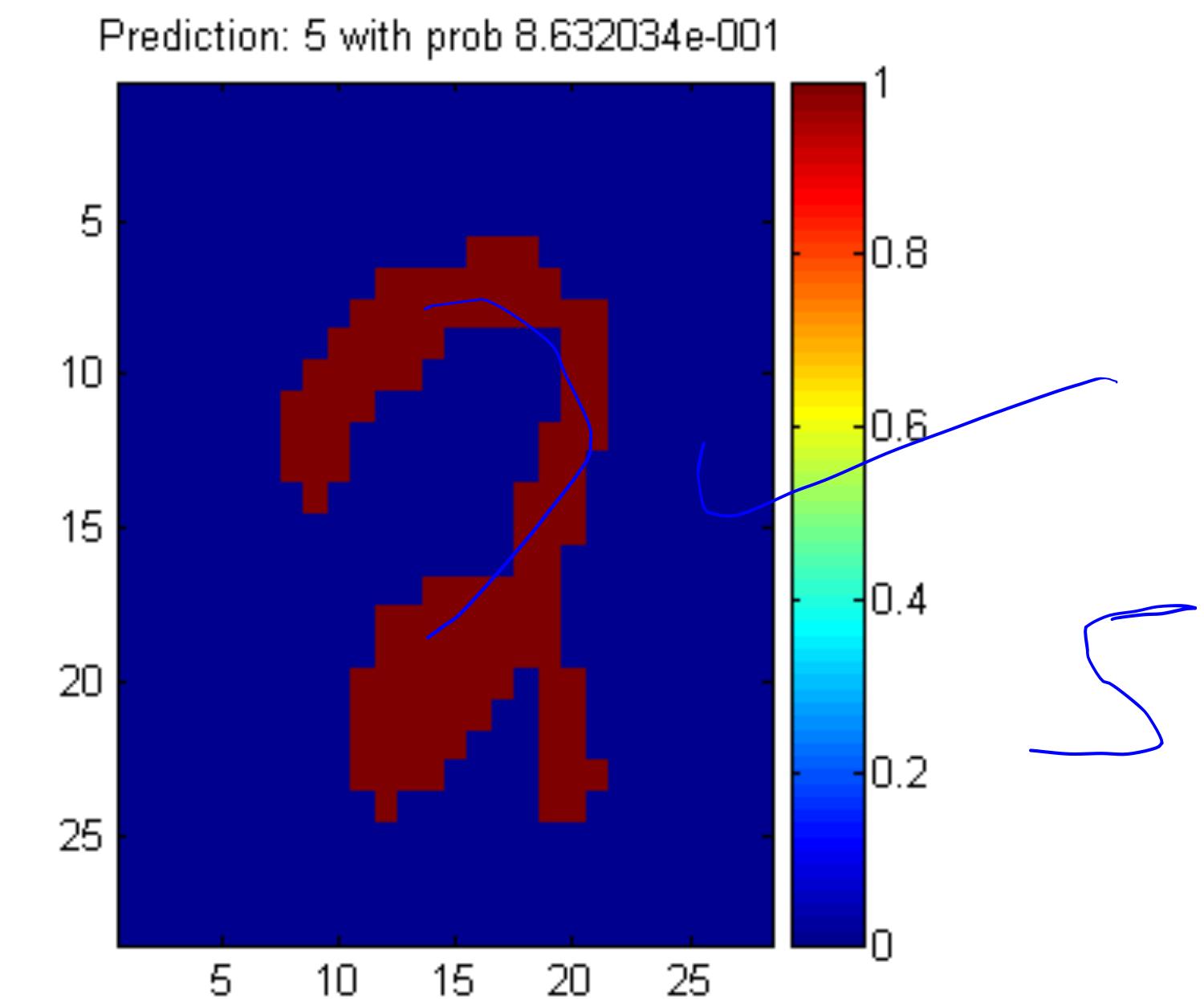
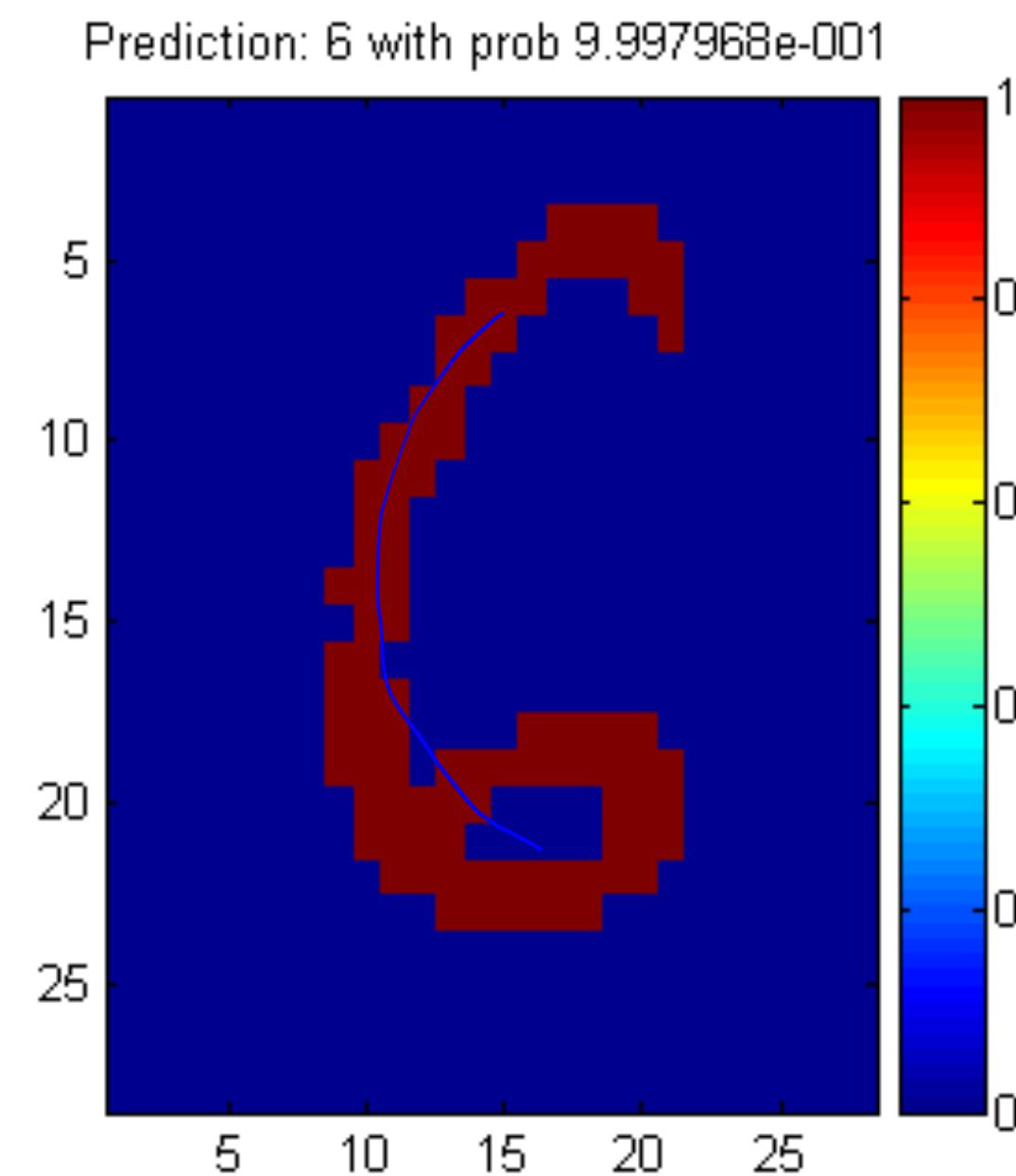
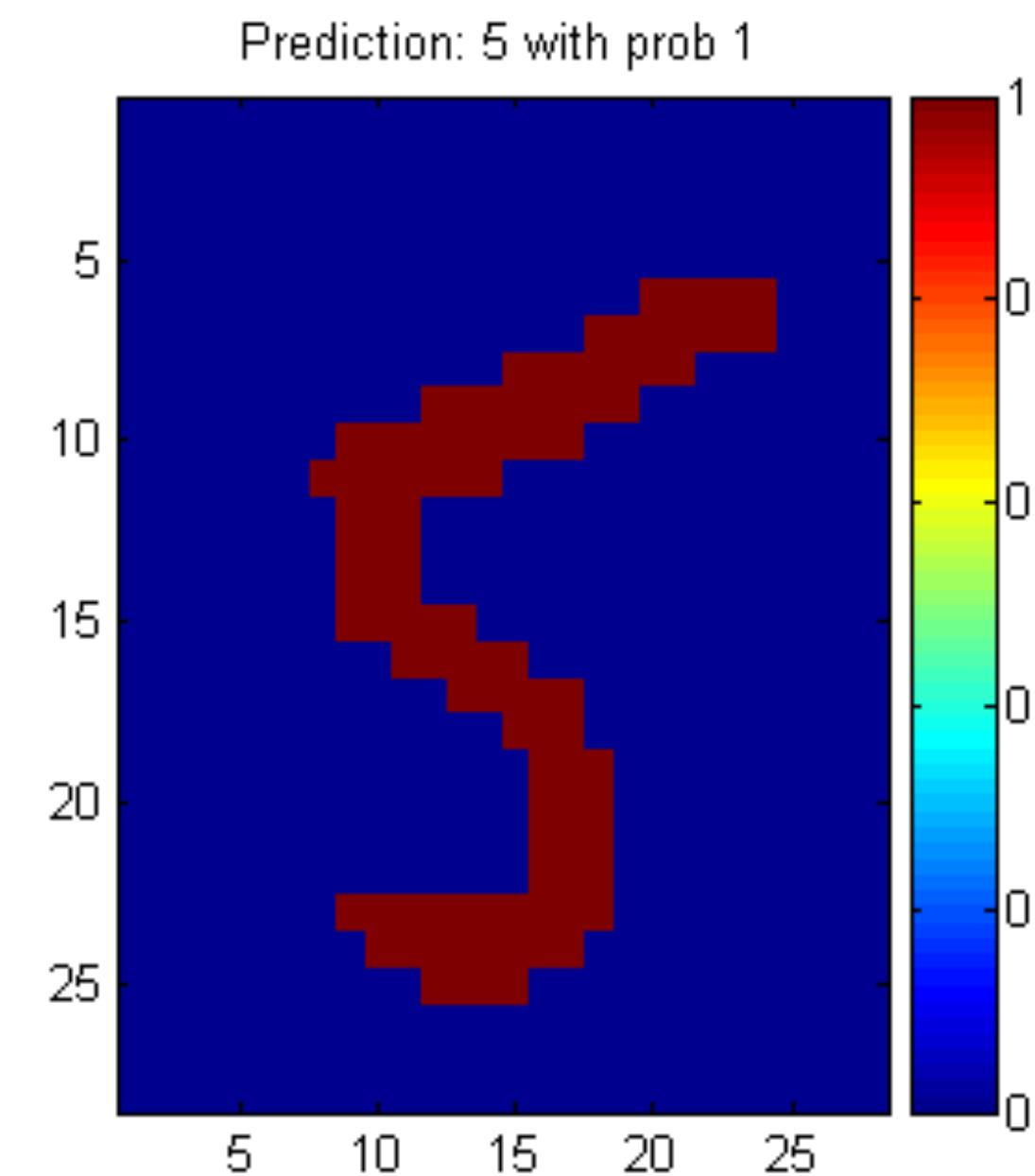
Naïve Bayes Training

- Now that we've decided to use a Naïve Bayes classifier, we need to train it with some data:



MNIST Training Data

Naïve Bayes Classification



Naïve Bayes Training

- Training in Naïve Bayes is **easy**:

- Estimate $P(Y=v)$ as the fraction of records with $Y=v$

$$P(Y = v) = \frac{\text{Count}(Y = v)}{\# \text{ records}}$$

- Estimate $P(X_i=u | Y=v)$ as the fraction of records with $Y=v$ for which $X_i=u$

$$P(X_i = u | Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v)}{\text{Count}(Y = v)}$$

Naïve Bayes (zero issue)

- In practice, some of these counts can be zero
- Fix this by adding “virtual” counts:

$$P(X_i = u|Y = v) = \frac{\text{Count}(X_i = u \wedge Y = v) + 1}{\text{Count}(Y = v) + 2}$$

- (This is like putting a prior on parameters and doing MAP estimation instead of MLE)
- This is called *Smoothing*

Underflow Prevention

- Multiplying lots of probabilities
→ floating-point underflow.

- Recall: $\log(xy) = \log(x) + \log(y),$

→ better to sum logs of probabilities rather than multiplying probabilities.

Conditional Independence

- Probabilities of going to the beach and getting a heat stroke are not independent $P(B, S) > P(B) P(S)$
- May be independent, if we know the weather is hot $P(B, S|H) = P(B|H) P(S|H)$
- Hot weather explains the dependence between beach and heatstroke

