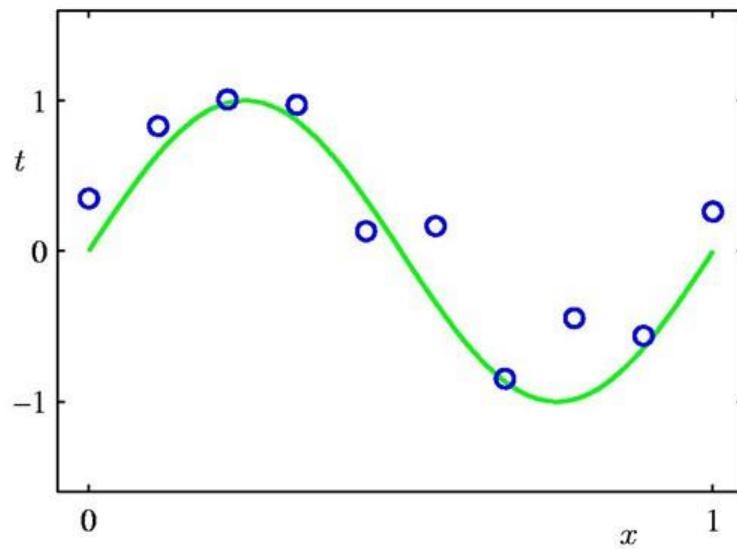


# Introduction to Probability distributions

---

# Polynomial Curve Fitting

---

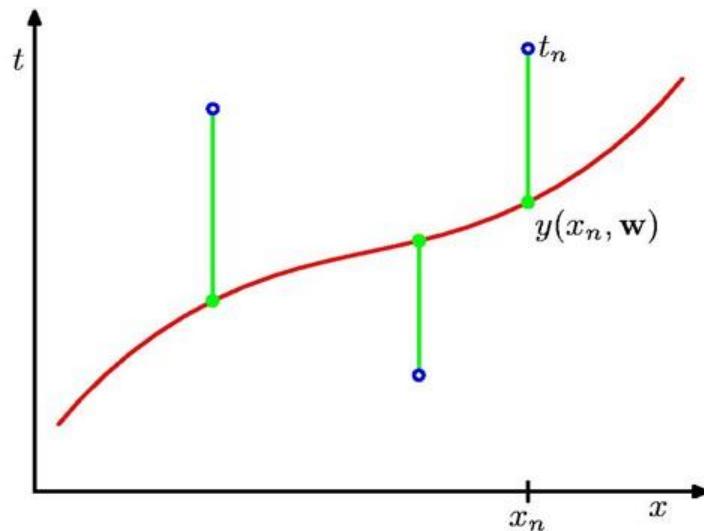


$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

---

# Sum-of-Squares Error Function

---

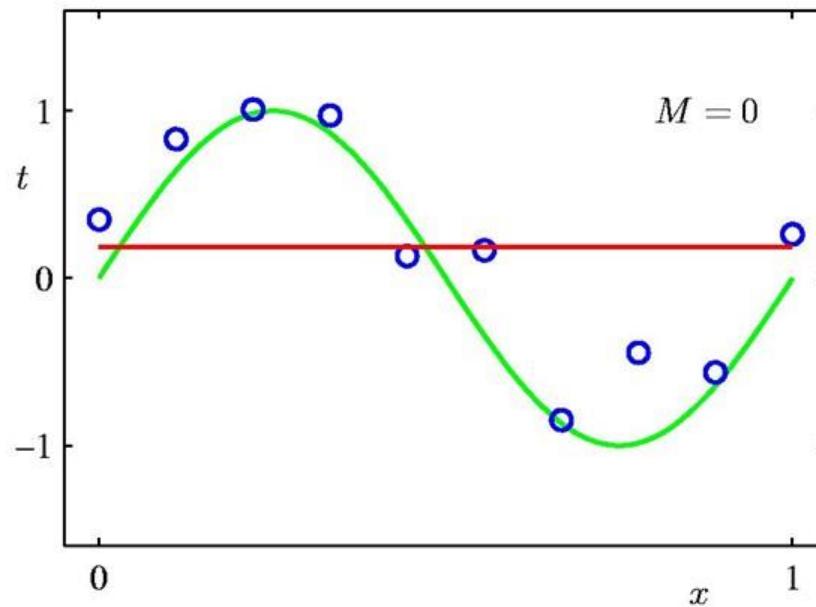


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

---

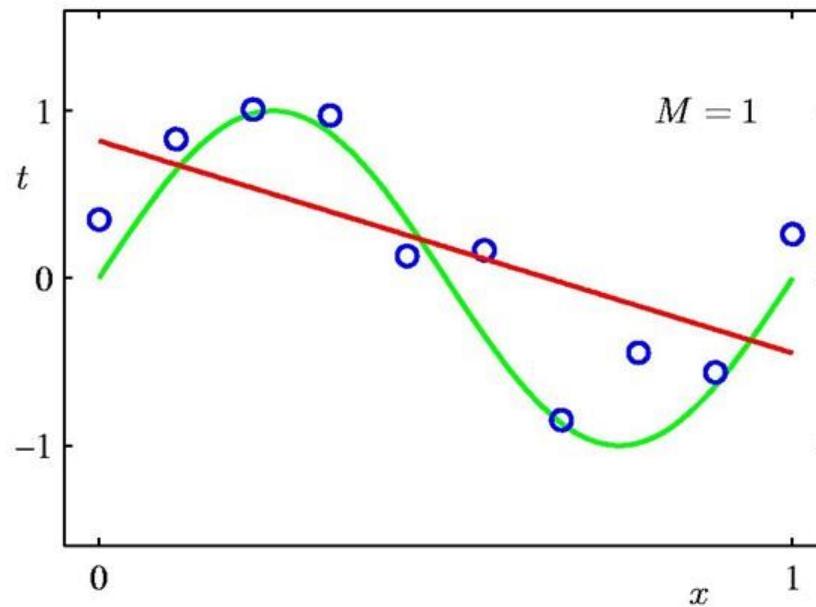
# $0^{\text{th}}$ Order Polynomial

---



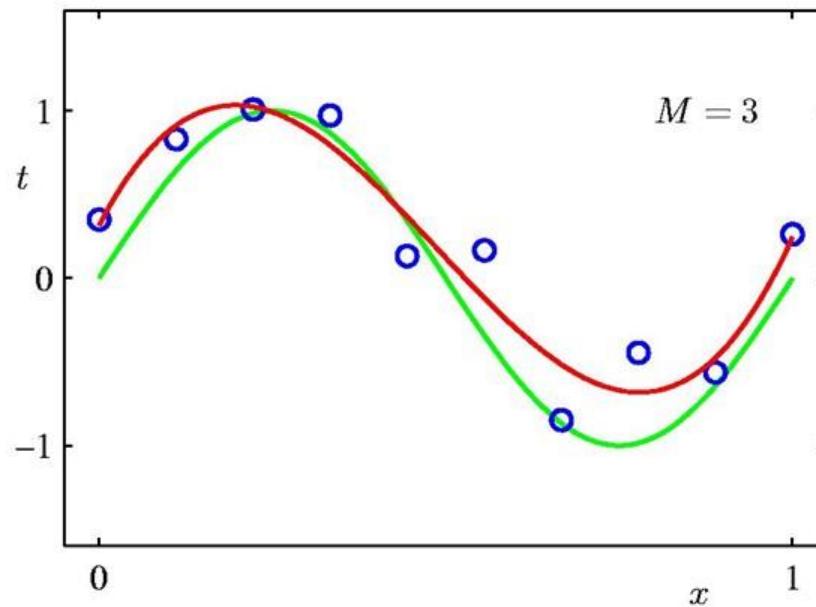
# 1<sup>st</sup> Order Polynomial

---



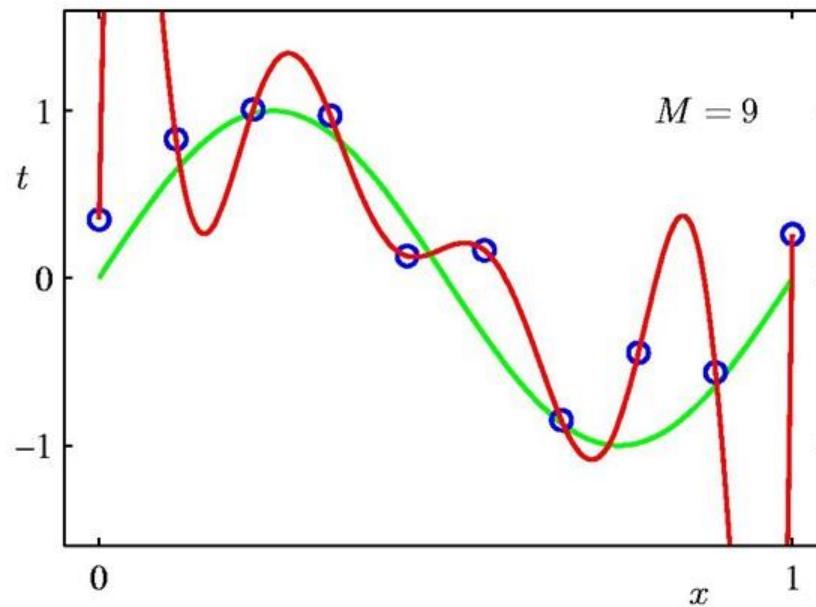
## 3<sup>rd</sup> Order Polynomial

---



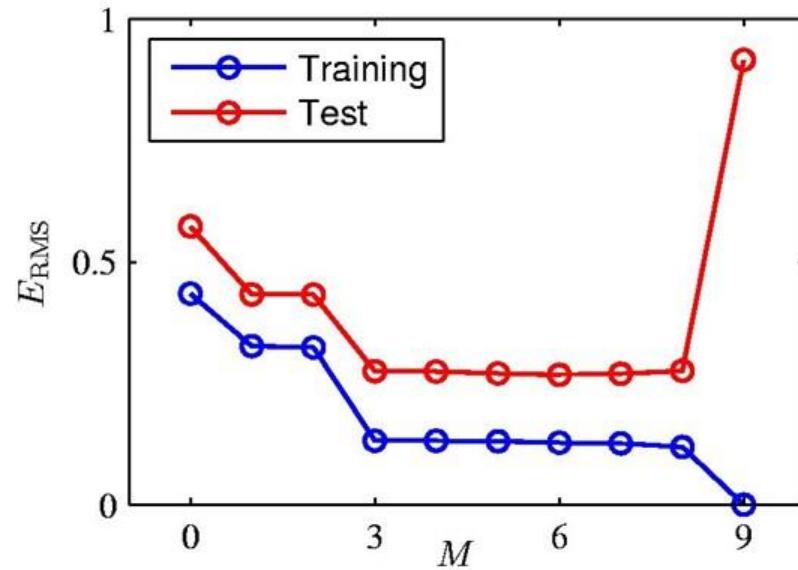
## 9<sup>th</sup> Order Polynomial

---



# Over-fitting

---



Root-Mean-Square (RMS) Error:  $E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$

---

# Polynomial Coefficients

---

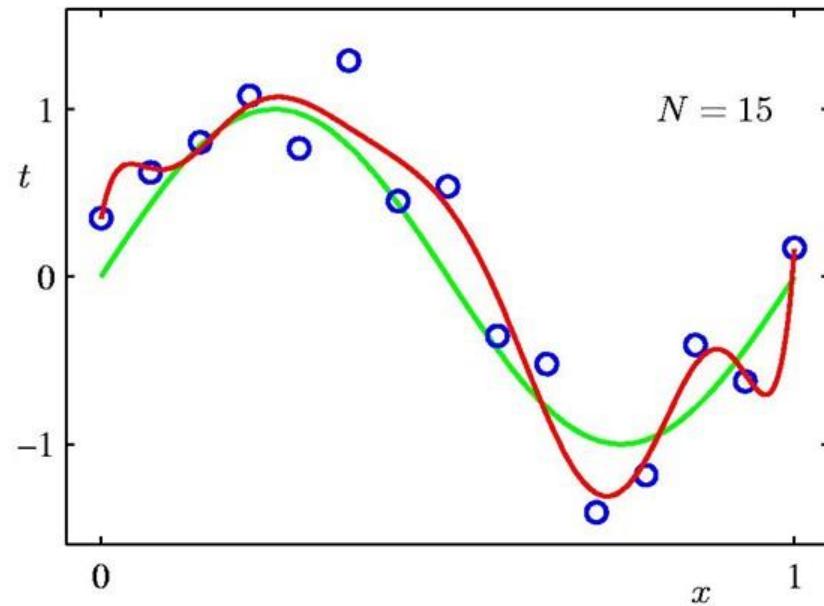
	$M = 0$	$M = 1$	$M = 3$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

---

Data Set Size:  $N = 15$

---

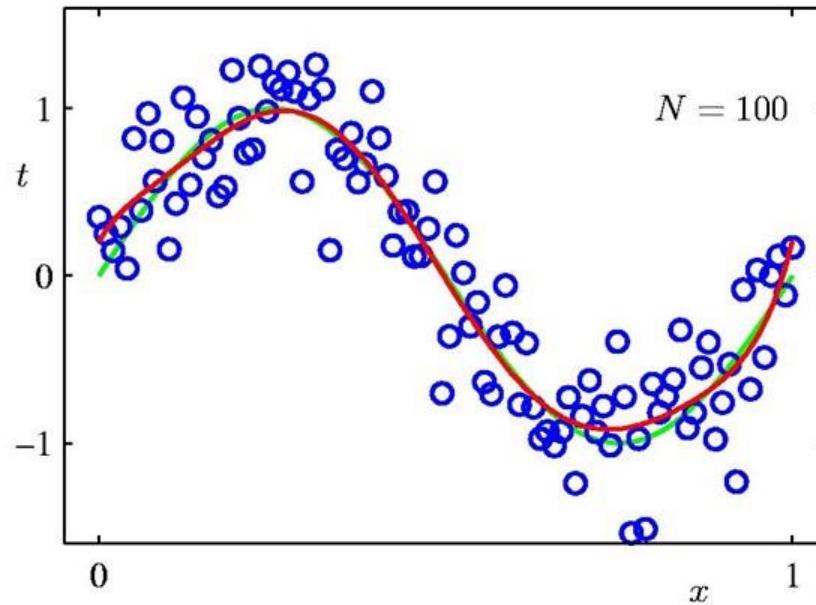
9<sup>th</sup> Order Polynomial



Data Set Size:  $N = 100$

---

9<sup>th</sup> Order Polynomial



# Regularization

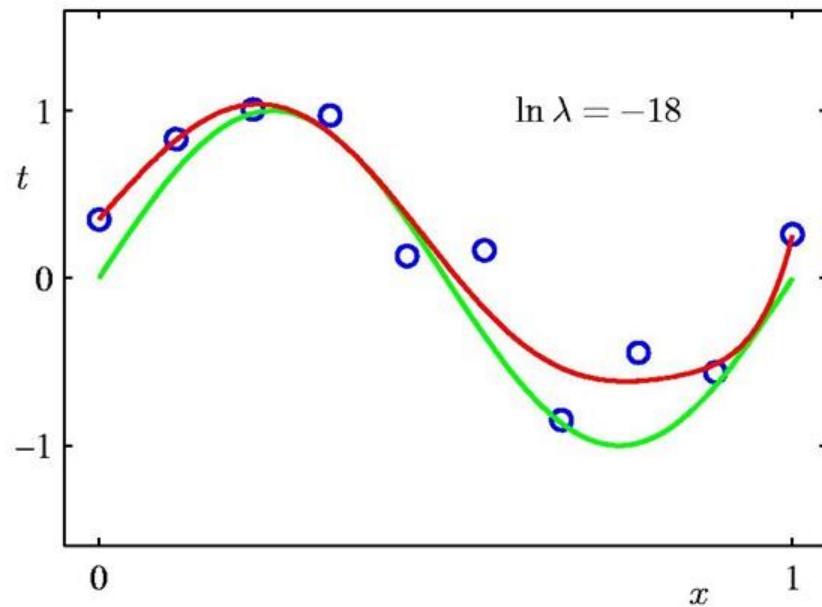
---

Penalize large coefficient values

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

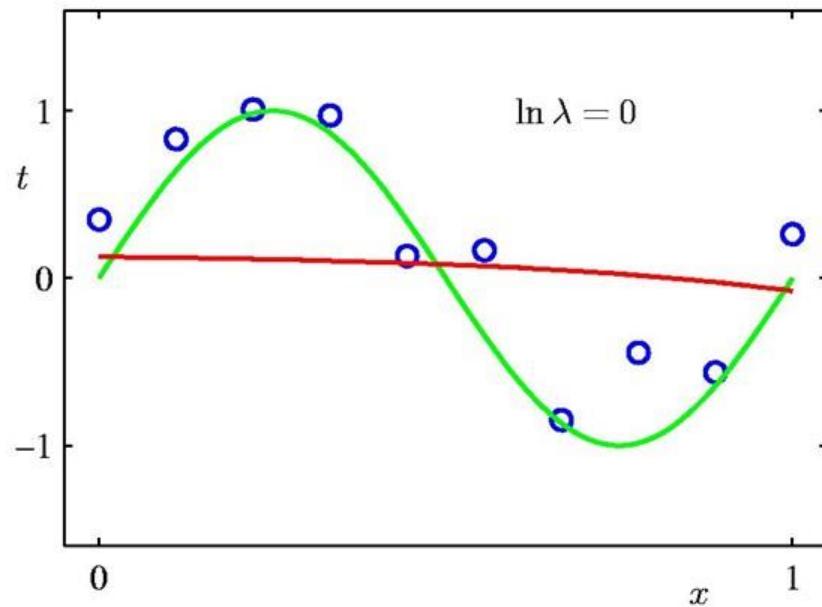
## Regularization: $\ln \lambda = -18$

---



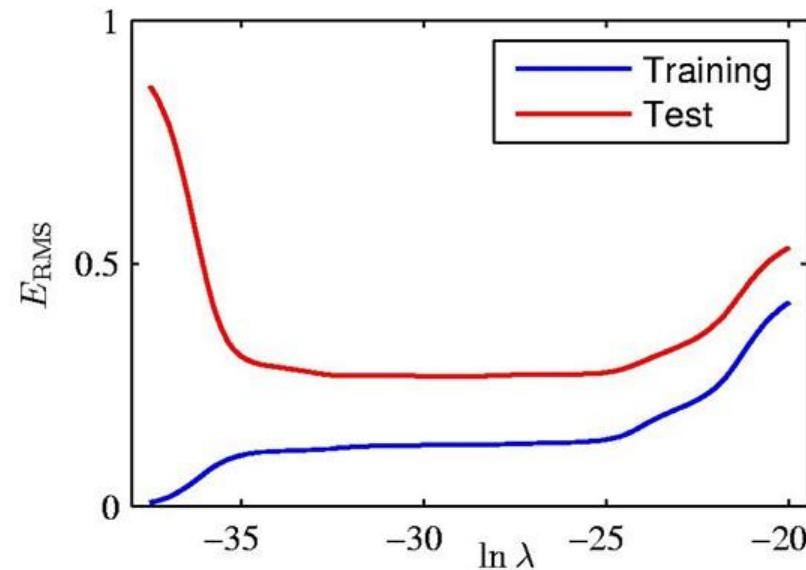
## Regularization: $\ln \lambda = 0$

---



## Regularization: $E_{\text{RMS}}$ vs. $\ln \lambda$

---



# Polynomial Coefficients

---

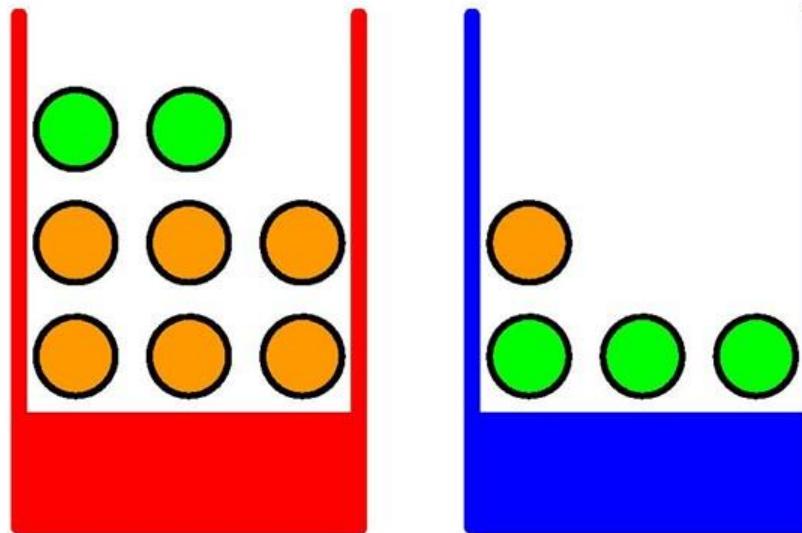
	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01

---

# Probability Theory

---

Apples and Oranges

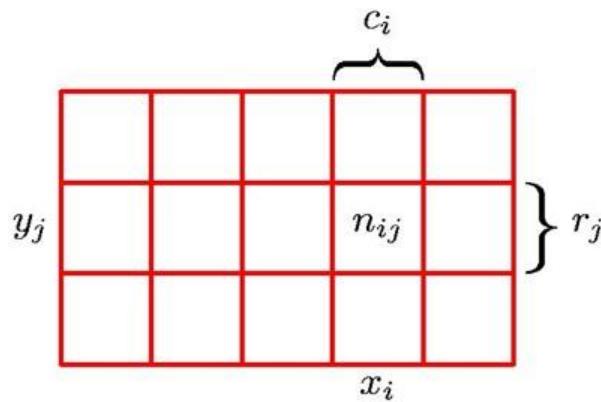


- 
- ❑ Consider two random variables  $X$  and  $Y$   
(which could for instance be the Box and Fruit variables considered above).
  - ❑ We shall suppose that
    - $X$  can take any of the values  $x_i$  where  $i = 1, \dots, M$ ,  
and
    - $Y$  can take the values  $y_j$  where  $j = 1, \dots, L$ .
  - ❑ Consider a total of  $N$  trials in which we sample both of the variables  $X$  and  $Y$ , and
-

- 
- Let the **number of such trials** in which  $X = x_i$  and  $Y = y_j$  be  $n_{ij}$  .
  - Also, let the **number of trials** in which  $X$  takes the value  $x_i$  (**irrespective of the value that  $Y$  takes**) be denoted by  $c_i$  ,
  - and similarly let the **number of trials** in which  $Y$  takes the value  $y_j$  be denoted by  $r_j$ .
-

# Probability Theory

---



Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}.$$

Joint Probability

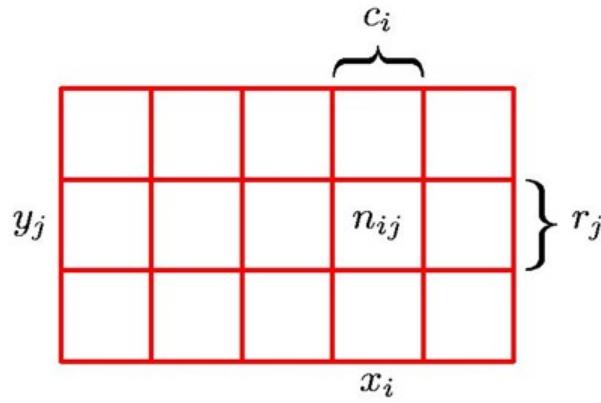
$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

# Probability Theory

---



Sum Rule

$$p(X = x_i) = \frac{c_i}{N} = \frac{1}{N} \sum_{j=1}^L n_{ij}$$
$$= \sum_{j=1}^L p(X = x_i, Y = y_j)$$

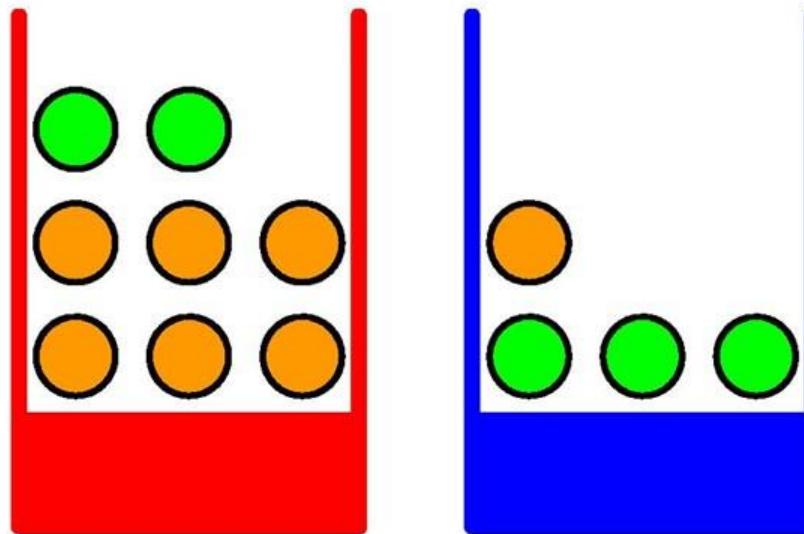
Product Rule

$$\begin{aligned} p(X = x_i, Y = y_j) &= \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} \\ &= p(Y = y_j | X = x_i) p(X = x_i) \end{aligned}$$

# Probability Theory

---

Apples and Oranges



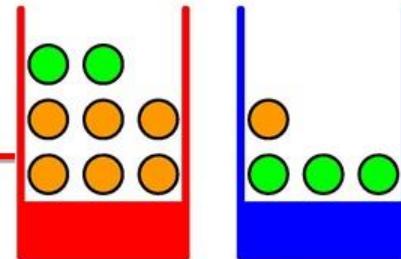
- 
- ❑ Let us suppose that we pick (with replacement)
  - ❑ the **red box 40% of the time** and we pick the **blue box 60% of the time**, and
  - ❑ that when we remove an item of fruit from a box we are **equally likely to select any of the pieces** of fruit in the box.
-

- 
- Let random variable **B** denotes the identity of Box chosen.
  - This random variable **B** can take two values:
    - r** (corresponding to the red box) or
    - b** (corresponding to the blue box).
  - Similarly, let random variable **F** denoted the identity of the fruits

random variable **F** can take values:

- a** (for apple) or
  - o** (for orange).
-

❑ Probability of choosing a box:

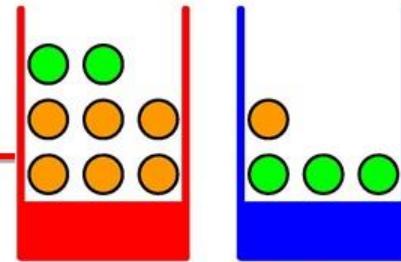


$$p(B = r) = 4/10 \text{ and } p(B = b) = 6/10$$

❑ Now Let us find answers of followings:

(i) “What is the overall probability of picking an apple?”,

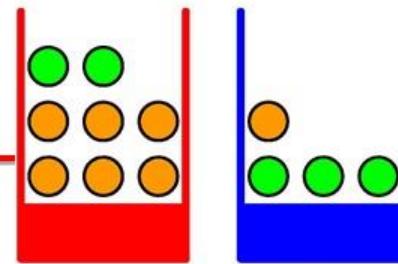
(ii) “Given that we have chosen an orange,  
what is the probability that it is from blue box?”



$$p(B = r) = 4/10$$

$$p(B = b) = 6/10$$

- Note that these satisfy  $p(B = r) + p(B = b) = 1$ .
  
  
  
  
  
  - Now suppose that we pick a box at random, and it turns out to be the blue box.
  
  
  
  
  
  - Then the probability of selecting an apple is just the fraction of apples in the blue box which is  $3/4$ , and so
- $p(F = a \mid B = b) = 3/4$ .
-



- ❑ In fact, we can write out all four **conditional probabilities** for each type of fruit, given the selected box:

$$p(F=a \mid B=r) = 1/4$$

$$p(F=o \mid B=r) = 3/4$$

$$p(F=a \mid B=b) = 3/4$$

$$p(F=o \mid B=b) = 1/4$$

- 
- ❑ Again, note that **these probabilities are normalized** so that

$$p(F = a | B = r) + p(F = o | B = r) = 1$$
$$p(F = a | B = b) + p(F = o | B = b) = 1.$$

and similarly

- ❑ We can now use **the sum and product rules** of probability to evaluate the **overall probability of choosing an apple**.
-

- 
- The overall probability of choosing an apple

$$p(F = a) = p(F = a | B = r)p(B = r) + p(F = a | B = b)p(B = b)$$

$$p(F=a)=11/20$$

- from which it follows, using the [sum rule](#), than

$$p(F = o) = 1 - 11/20 = 9/20.$$

- 
- Suppose instead we are told that:
    - a piece of fruit has been selected and it is an orange,
    - and we would like to know which box it came from.
  - This requires that we evaluate
    - the probability distribution of boxes given the identity of the fruit,
-

---

whereas the probabilities **calculated earlier** give

the **probability distribution of fruit given the identity of the box.**

We can solve the problem of **reversing the conditional probability** by **using Bayes' theorem** to give:

$$p(B = r | F = o) = p(F = o | B = r)p(B = r) / p(F = o)$$

$$p(B = r | F = o) = 2/3$$

From the sum rule, it then follows that  $p(B = b | F = o) = 1 - 2/3 = 1/3.$

- 
- ❑ So far we have been quite careful to make a distinction between:
    - a random variable (such as the box **B** in the fruit example)
    - and
    - the values that the random variable can take  
(for example **r** if the box were the red one).
  - ❑ Thus the probability that **B** takes the value **r** is denoted **p(B = r)**.
-

- 
- ❑ Although this helps to **avoid ambiguity**, it leads to a rather cumbersome notation, and
  - ❑ in many cases there will **be no need for such pedantry**.
- 
- ❑ Instead, we **may simply write  $p(B)$  to denote a distribution over the random variable  $B$ ,**  
or
  - ❑  **$p(r)$  to denote the distribution evaluated for the particular value  $r$ ,**  
provided that the interpretation is clear from the context.
-

### **Interpretation of Bayes' theorem:**

---

- If it is asked “which box had been chosen”  
before being told the identity of the selected item of fruit,  
  
then the most complete information we have is provided by  
  
the probability  $p(B)$ .
  
  - We call this the ***prior probability*** because  
  
it is the probability available before we observe the identity of  
the fruit.
-

### **Interpretation of Bayes' theorem:**

---

- Once we are told that the **fruit is an orange**,
  
  - we can then use Bayes' theorem  
to compute the probability  $p(B|F)$ ,
  
  - which we shall call the ***posterior probability*** because it is the probability obtained after we have observed F.
-

### **Interpretation of Bayes' theorem...**

---

- ❑ Note that in this example, the prior probability of selecting the red box was  $4/10$ ,

so that we were more likely to select the blue box than the red one.

- ❑ However, once we have observed that the piece of selected fruit is an orange,

we find that the posterior probability of the red box is now  $2/3$ ,

- ❑ so that it is now more likely that the box we selected was in fact the red one.
-

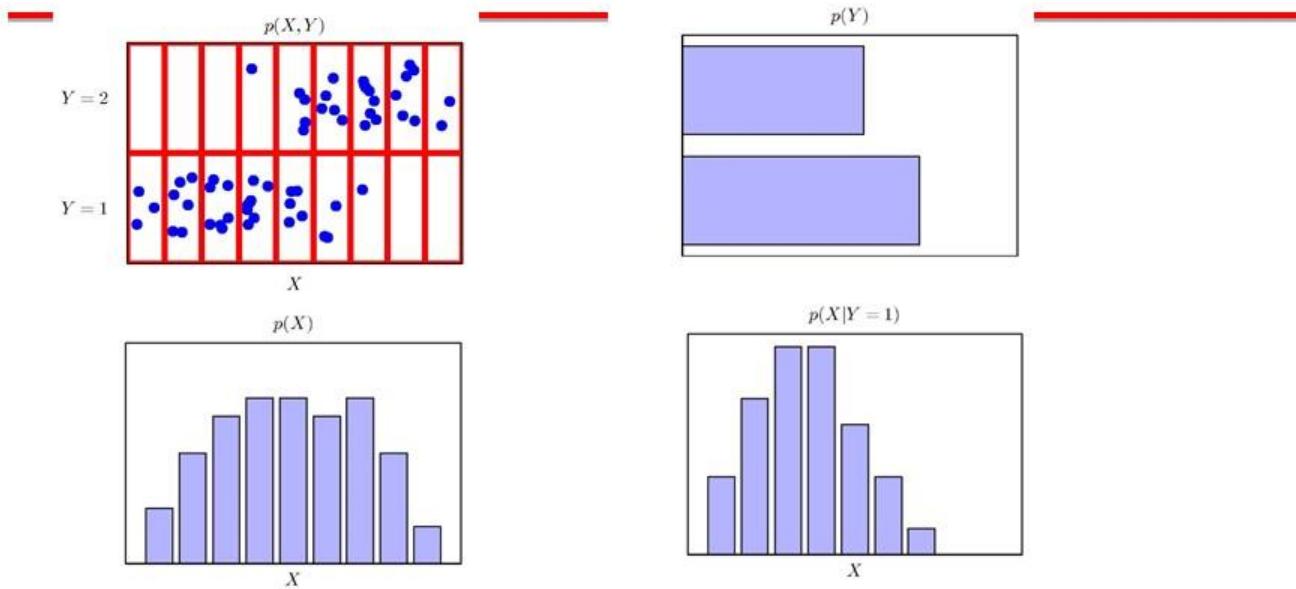
### **Interpretation of Bayes' theorem...**

---

- ❑ This result accords with our intuition,
  - ❑ as the **proportion of oranges is much higher in the red box than it is in the blue box**,
- and so the observation that the fruit was an orange provides significant evidence favouring the red box.
- ❑ In fact, the **evidence is sufficiently strong that it outweighs the prior** and

makes it more likely that the red box was chosen rather than the blue one.

---



An illustration of a distribution over two variables,  $X$ , which takes 9 possible values, and  $Y$ , which takes two possible values. The top left figure shows a sample of 60 points drawn from a joint probability distribution over these variables. The remaining figures show histogram estimates of the marginal distributions  $p(X)$  and  $p(Y)$ , as well as the conditional distribution  $p(X|Y=1)$  corresponding to the bottom row in the top left figure.

---

# The Rules of Probability

---

**Sum Rule**

$$p(X) = \sum_Y p(X, Y)$$

**Product Rule**

$$p(X, Y) = p(Y|X)p(X)$$

# Bayes' Theorem

---

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

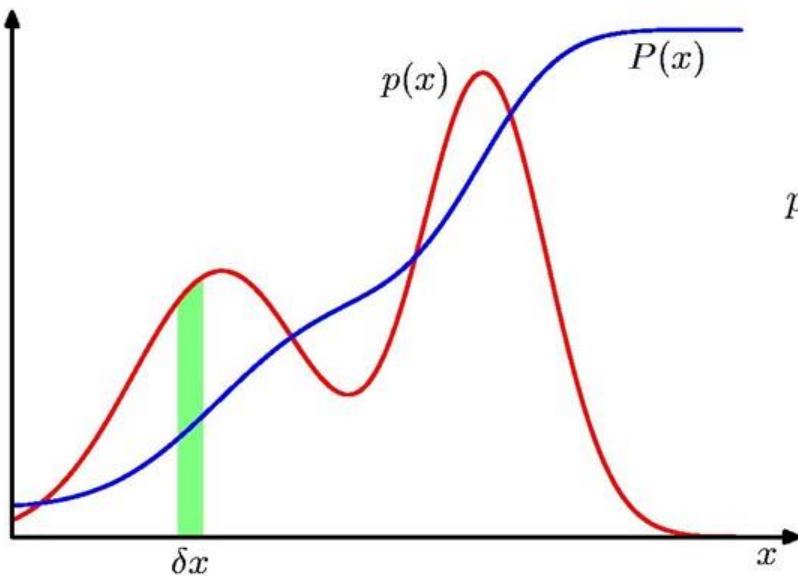
$$p(X) = \sum_Y p(X|Y)p(Y)$$

posterior  $\propto$  likelihood  $\times$  prior

---

# Probability Densities

---



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

$$P(z) = \int_{-\infty}^z p(x) dx$$

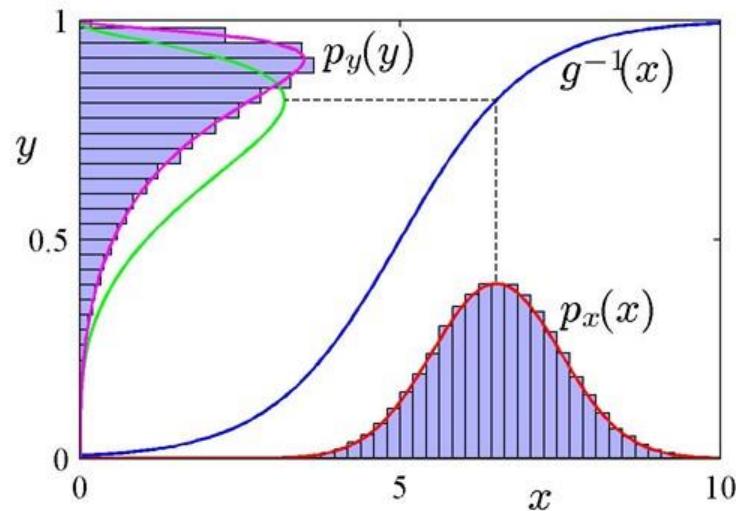
$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

---

# Transformed Densities

---



$$\begin{aligned} p_y(y) &= p_x(x) \left| \frac{dx}{dy} \right| \\ &= p_x(g(y)) |g'(y)| \end{aligned}$$

# Expectations

---

$$\mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x) dx$$

$$\mathbb{E}_x[f|y] = \sum_x p(x|y)f(x)$$

↑  
-----> x

Conditional Expectation  
(discrete)

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^N f(x_n)$$

Approximate Expectation  
(discrete and continuous)

# Variances and Covariances

---

$$\text{var}[f] = \mathbb{E} \left[ (f(x) - \mathbb{E}[f(x)])^2 \right] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

$$\begin{aligned}\text{cov}[x, y] &= \mathbb{E}_{x,y} [\{x - \mathbb{E}[x]\} \{y - \mathbb{E}[y]\}] \\ &= \mathbb{E}_{x,y}[xy] - \mathbb{E}[x]\mathbb{E}[y]\end{aligned}$$

$$\begin{aligned}\text{cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}_{\mathbf{x},\mathbf{y}} [\{\mathbf{x} - \mathbb{E}[\mathbf{x}]\}\{\mathbf{y}^T - \mathbb{E}[\mathbf{y}^T]\}] \\ &= \mathbb{E}_{\mathbf{x},\mathbf{y}}[\mathbf{x}\mathbf{y}^T] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^T]\end{aligned}$$

---

## Curve Fitting: Bayesian Approach

---

Consider the example of polynomial curve fitting discussed earlier.

Lets model uncertainty in selecting appropriate parameters  $w$  (Bayesian perspective)

Recall that in the boxes of fruit example, the **observation of the identity of the fruit** provided relevant information that **altered** the probability that the **chosen box was the red one**.

In that example, Bayes' theorem was used to **convert a prior** probability into a **posterior** probability **by incorporating the evidence** provided by the observed data.

As we shall see in detail later, we can adopt a similar approach when making inferences about quantities such as the parameters  $w$  in the polynomial curve fitting example.

---

---

$$P(c | x) = \frac{P(x | c)P(c)}{P(x)}$$

Likelihood                      Class Prior Probability

↓                                  ↓

Posterior Probability      Predictor Prior Probability

---

$$P(c | X) = P(x_1 | c) \times P(x_2 | c) \times \cdots \times P(x_n | c) \times P(c)$$

---

## 2.1 Decision Using Posterior Probability

### ■ Posterior Probabilities

Define  $P(\omega_j|x) = \left( \text{the probability of } \omega \text{ being } \omega_j \text{ given that } x \text{ has been measured} \right)$

Bayes rule derives

$$P(\omega_j|x) = \frac{p(x|\omega_j)P(\omega_j)}{p(x)} \quad (1)$$

### ■ Decision Rule (1) Minimizing error probability

$$\begin{array}{lll} \textbf{Decide} & \omega_1 & \text{if } P(\omega_1|x) > P(\omega_2|x) \\ & \omega_2 & \text{if } P(\omega_1|x) < P(\omega_2|x) \end{array} \quad (2)$$

### ■ Decision Rule (2) *Likelihood ratio*

$$\begin{array}{lll} \textbf{Decide} & \omega_1 & \text{if } \frac{p(x|\omega_1)}{p(x|\omega_2)} > \frac{P(\omega_1)}{P(\omega_2)} \\ & \omega_2 & \text{if } \frac{p(x|\omega_1)}{p(x|\omega_2)} < \frac{P(\omega_1)}{P(\omega_2)} \end{array} \quad (3)$$

independent of  
observation  $x$       10

## Continuous Random Variable

---

- Continuous random variables can have an **infinite number (uncountable)** of possible values within their defined interval or range.
  - Unlike discrete random variables, which can only assume distinct, countable values.
  - Continuous random variables are characterized by probability density functions (PDFs).
-

## Continuous Random Variable

---

- PDFs **describe the likelihood or density** of the random variable assuming different values within the range.
  - The PDF **specifies the relative likelihood of observing the variable at any specific value** or within any specific interval.
-

## Probability of a Specific Value vs Range of Values (interval)

---

- In continuous probability distributions, the probability of a specific, exact value is technically zero.
  - This is because continuous random variables can take on an infinite number of values within a given range, and
  - Thus, the probability associated with any single, isolated point is infinitesimal.
-

## Probability of a Specific Value vs Range of Values (interval)

---

- For example: You measured the height of students ( $H$ ) in a University campus with very high precision of measurement.
  - The height can be 170 cm, 170.01 cm, and 170.00006 cm etc.
  - Indeed between any two value say 170 cm and 170.01, there are infinite possible values of height.
-

## Probability of a Specific Value vs Range of Values (interval)

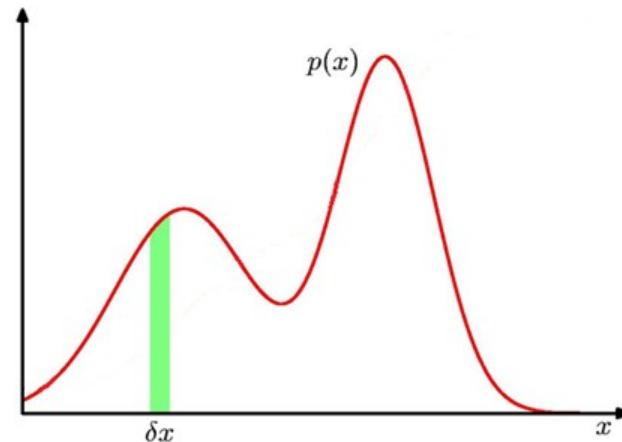
---

- In case of continuous random variable, it is meaningless to talk about probability of a specific value say  $p(H = 170 \text{ cm})$ .
  - As  $p(H = 170 \text{ cm}) = 0$ ,  $H = 170 \text{ cm}$  is one value against infinite possible values.
  - Thus, we calculate the **probability of the variable falling within a specified interval** or range of values in case of continuous random variables.
-

## Probability Density Function (PDF)

---

- It describes the probability distribution of a continuous random variable and
- It provides information about how the **values of that variable are likely to be distributed over a specified range.**
- PDF is always non negative.
- The total area under the PDF curve over the entire range is equal to 1.



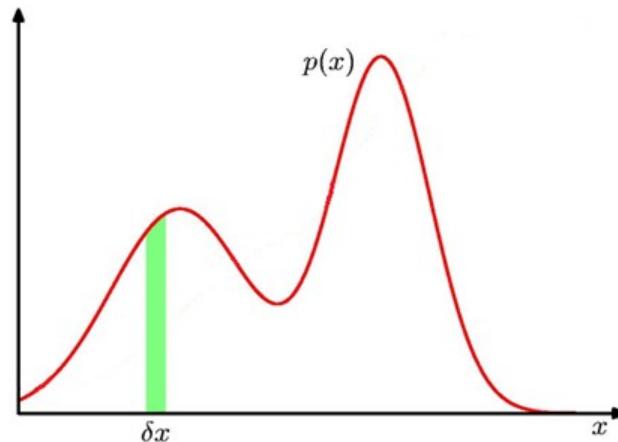
$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

## Probability in Intervals

---

- To find the probability that the continuous random variable falls within a specific interval, **integrate the PDF over that interval.**
- Probability of a continuous random variable falling within a specified interval or range of values, say  $(a, b)$  is



$$p(x) \geq 0$$

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

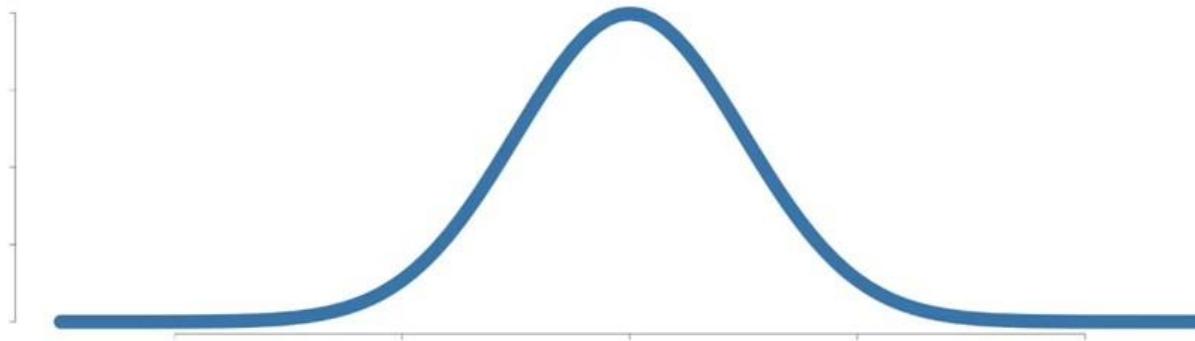
$$\int_{-\infty}^{\infty} p(x) dx = 1$$

---

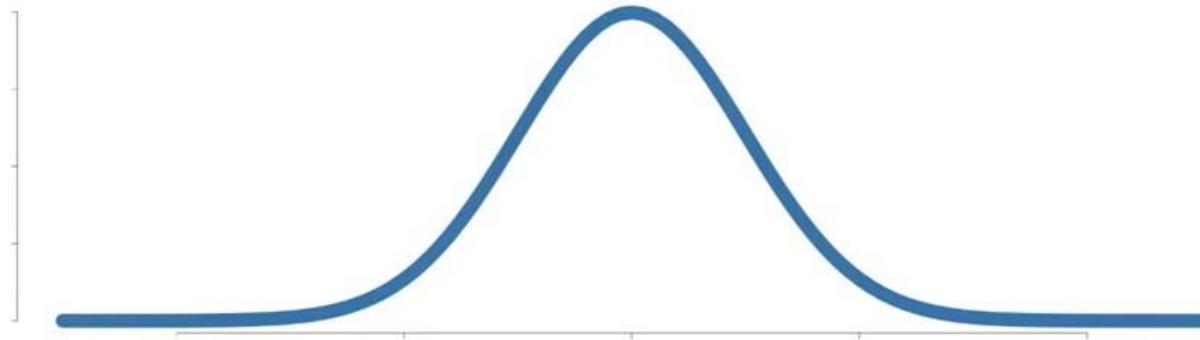
# **Likelihood**

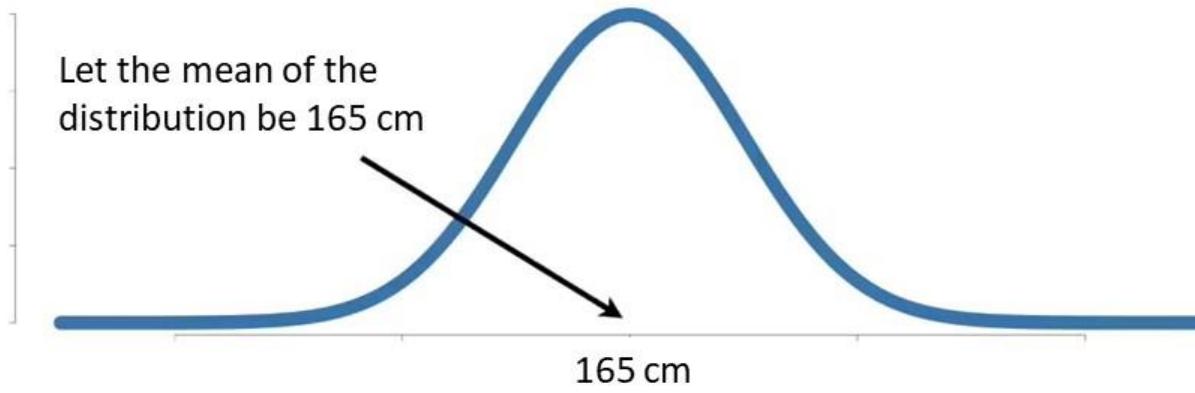
---

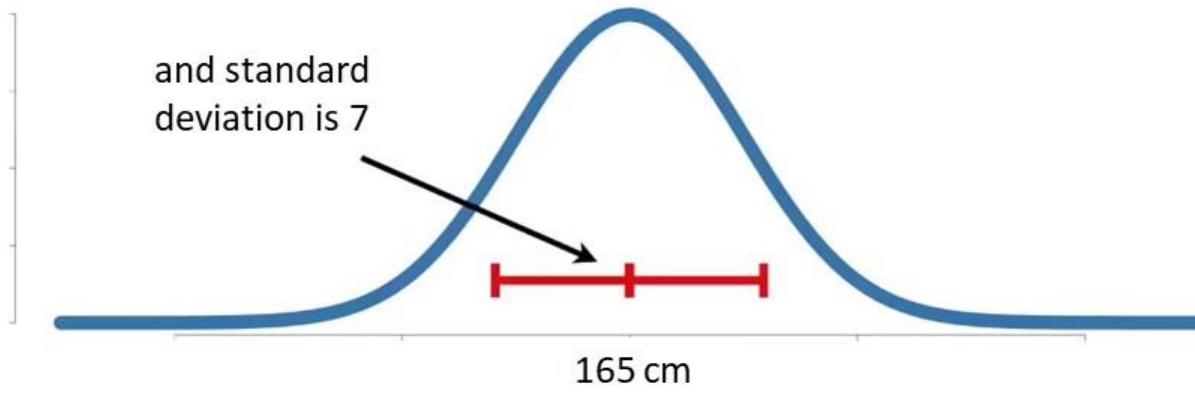
So let's start by looking at probability with respect to a normal distribution (keeping in mind that this concept applies to all continuous distributions).

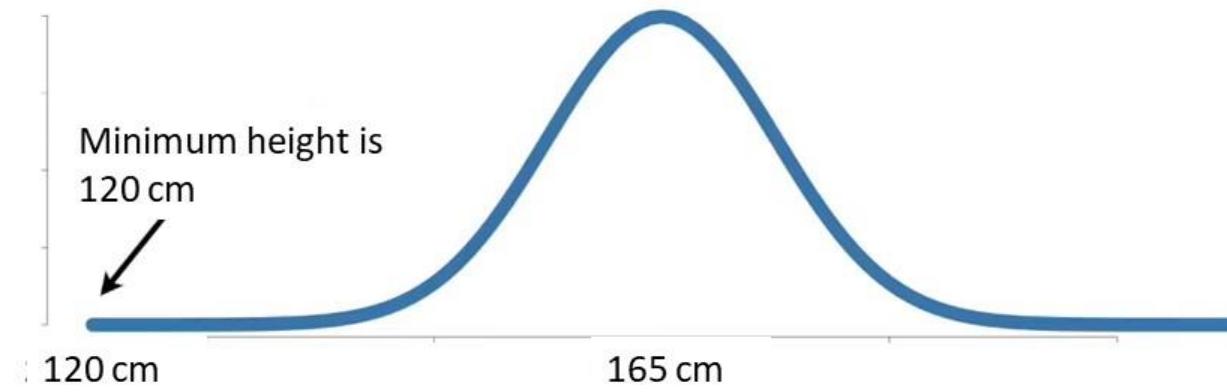


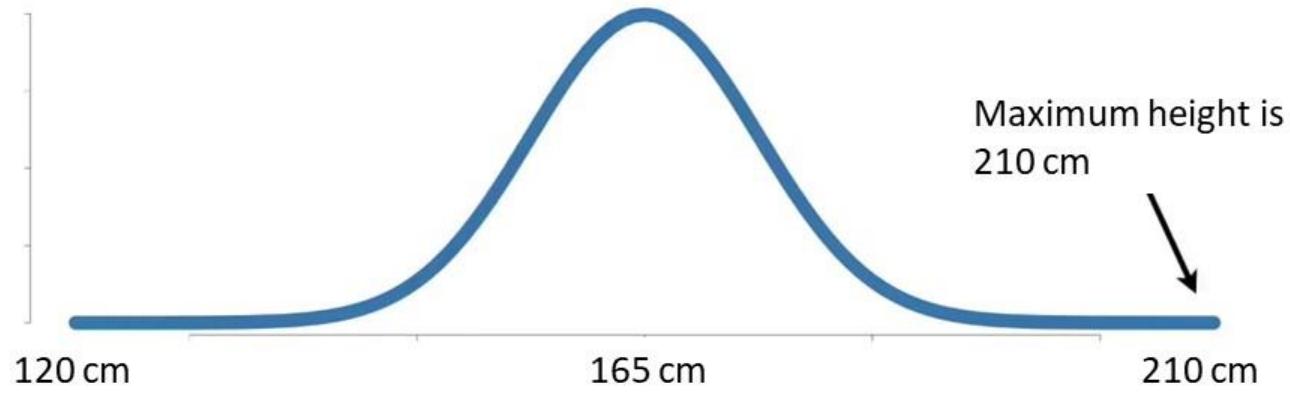
Let it be a distribution of students' height (measured with very high precision) of a University Campus





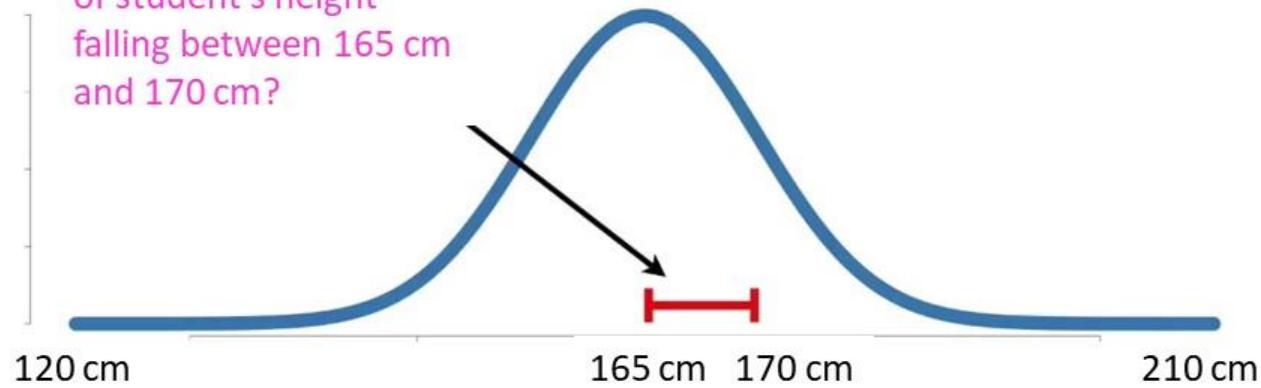




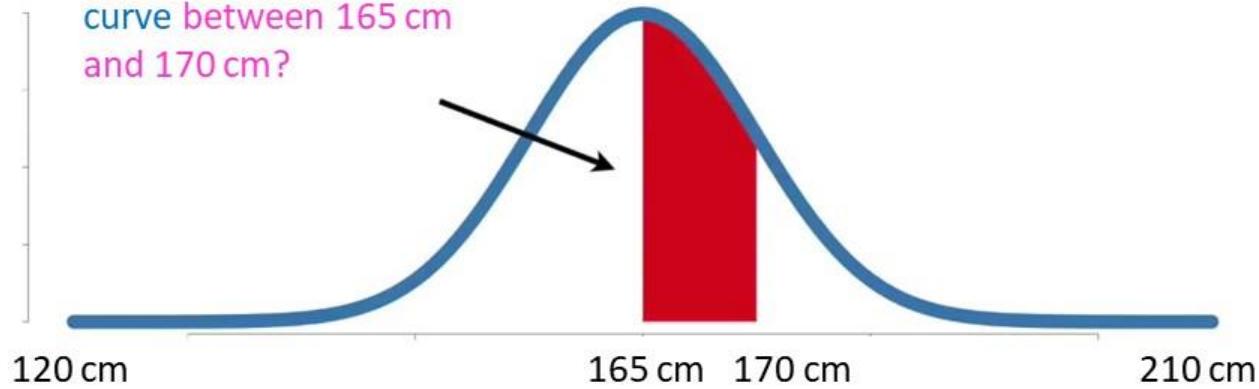


If we select a random student then

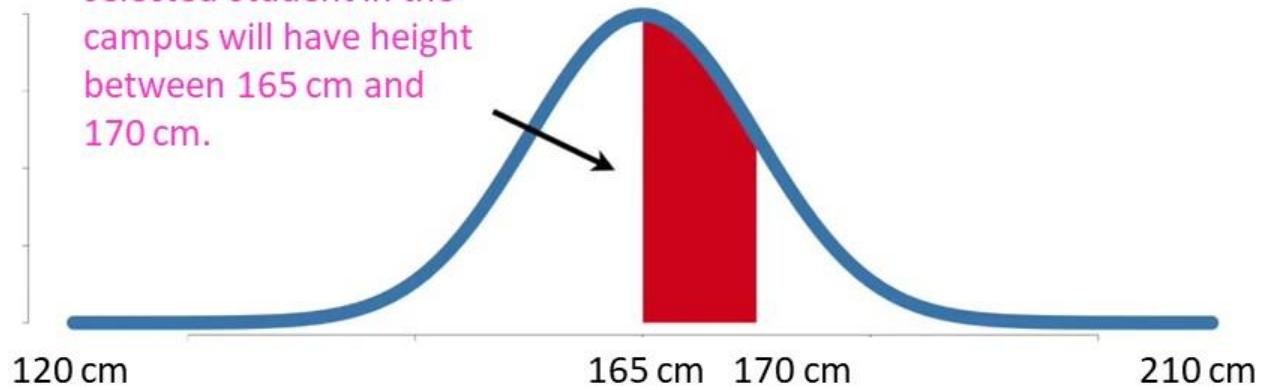
what is the probability of student's height falling between 165 cm and 170 cm?



It is the area under  
curve between 165 cm  
and 170 cm?

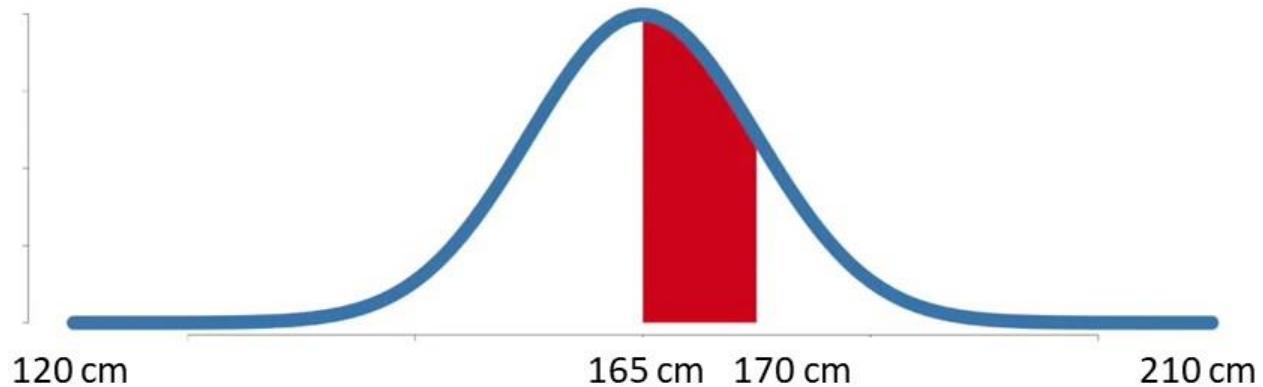


In this case, the area under curve = 0.26 means there is a 26% chance a randomly selected student in the campus will have height between 165 cm and 170 cm.

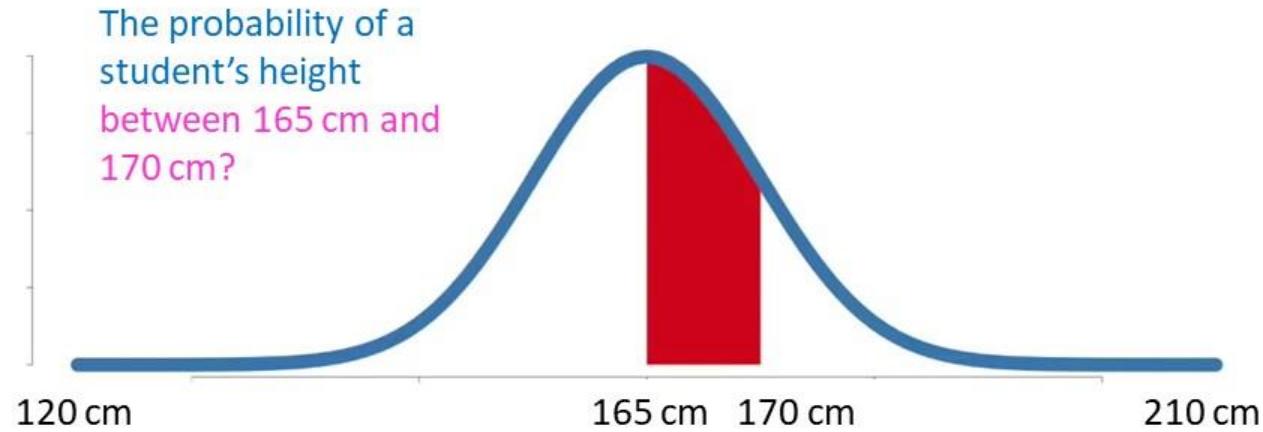


It is denoted as:

$$p(\text{height between } 165 \text{ and } 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$$

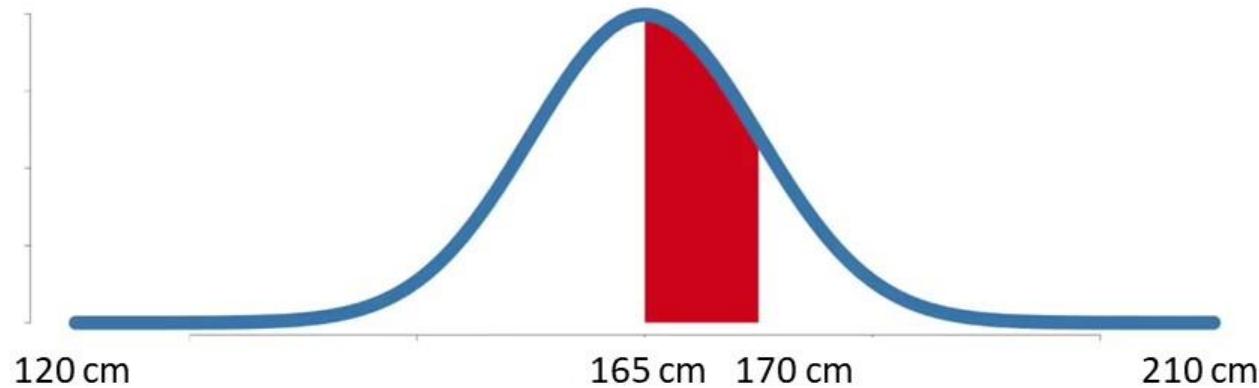


$$p(\text{height between } 165 \text{ and } 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$$

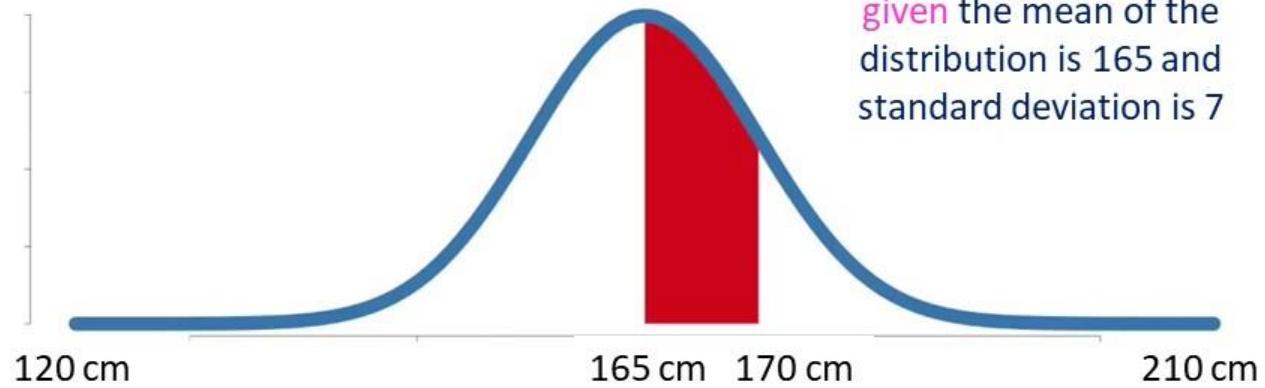


given

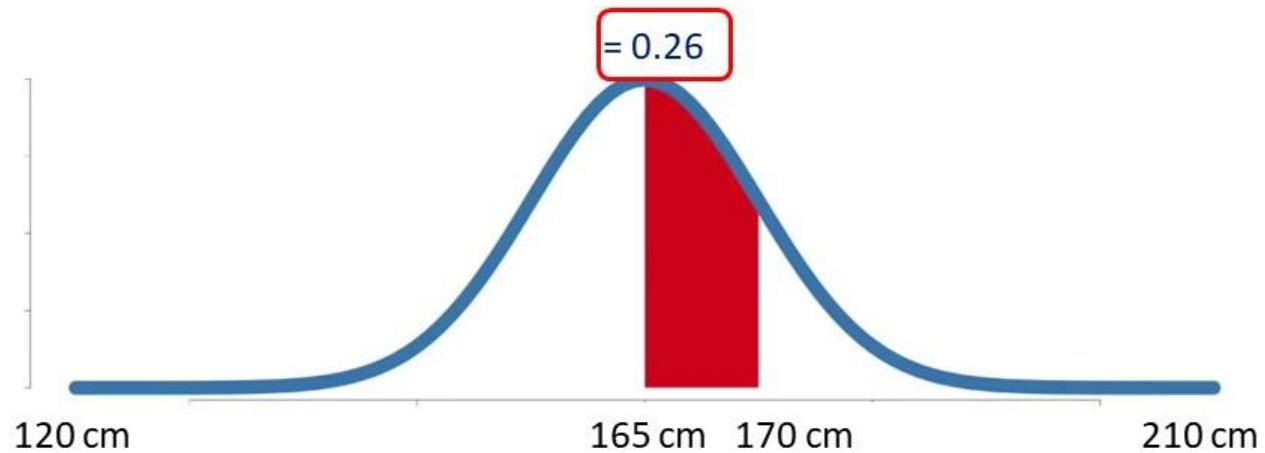
$$p(\text{height between } 165 \text{ and } 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$$



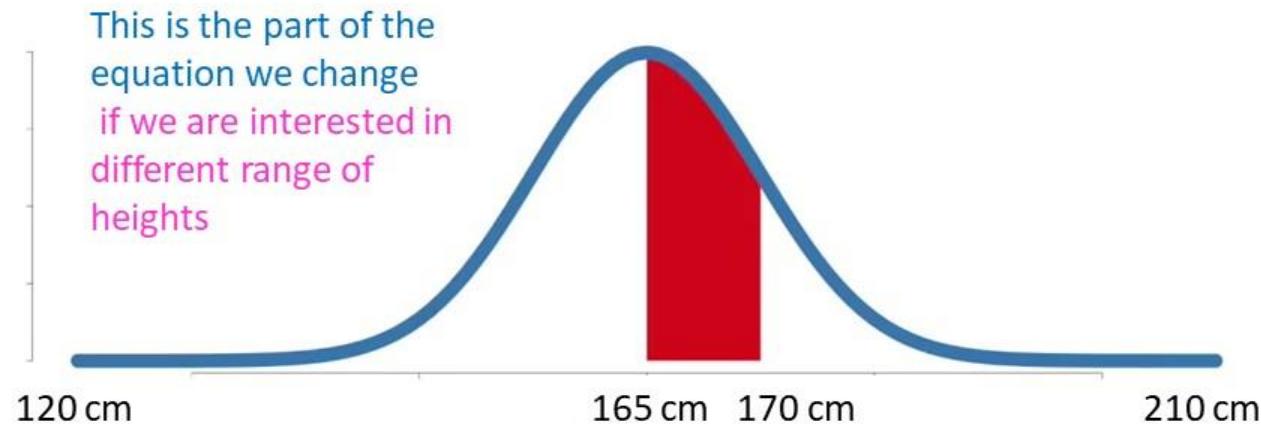
$p(\text{height between } 165 \text{ and } 170 \text{ cm} | \text{ mean} = 165 \text{ and SD} = 7)$



$p(\text{height between } 165 \text{ and } 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$

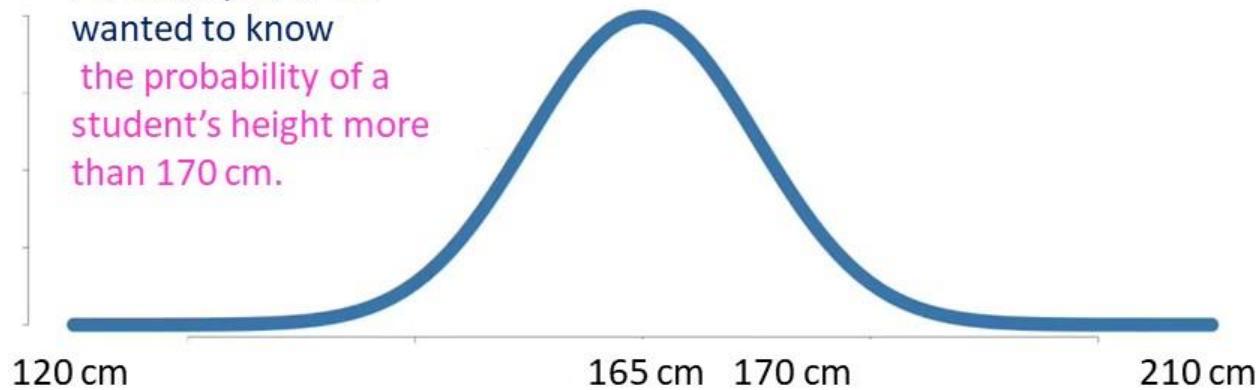


$p(\text{height between } 165 \text{ and } 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$

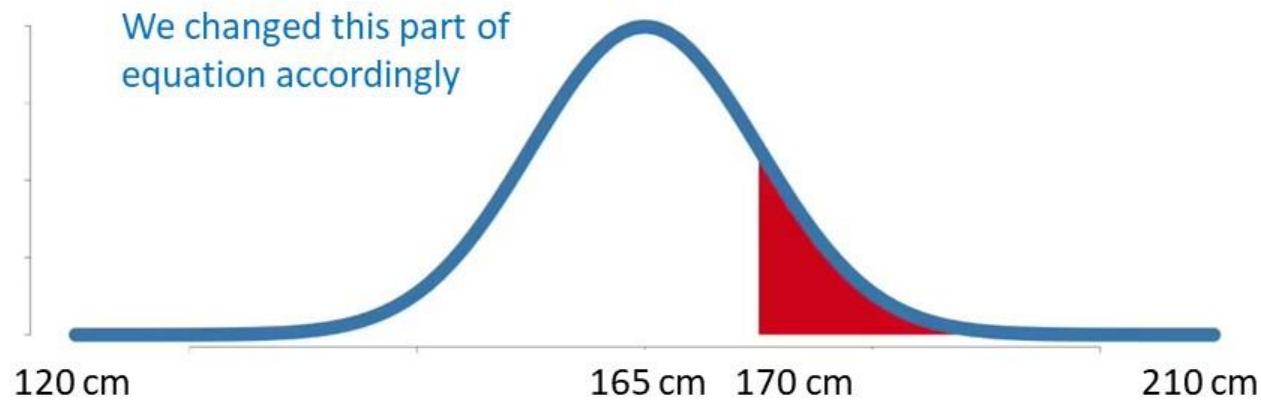


$$p(\text{height} > 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$$

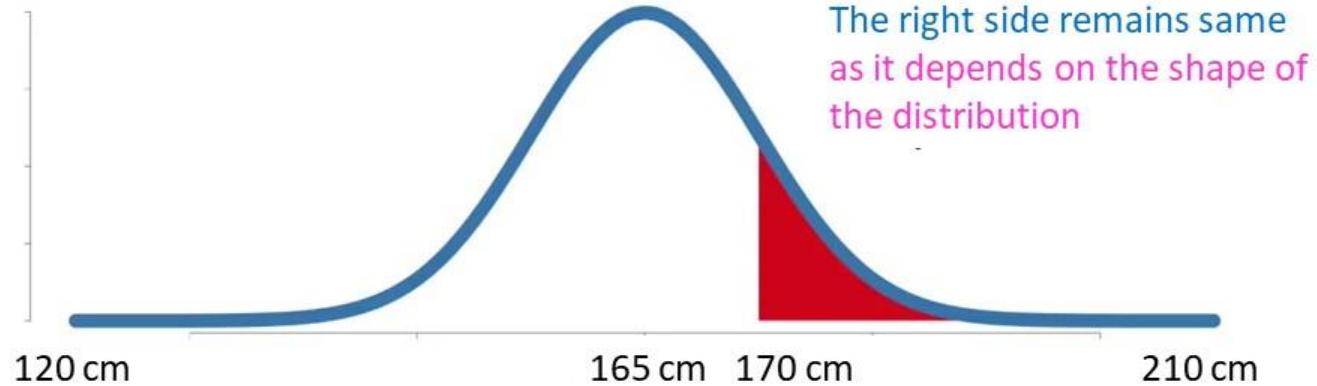
For example, if we wanted to know the probability of a student's height more than 170 cm.



$p(\text{height} > 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$



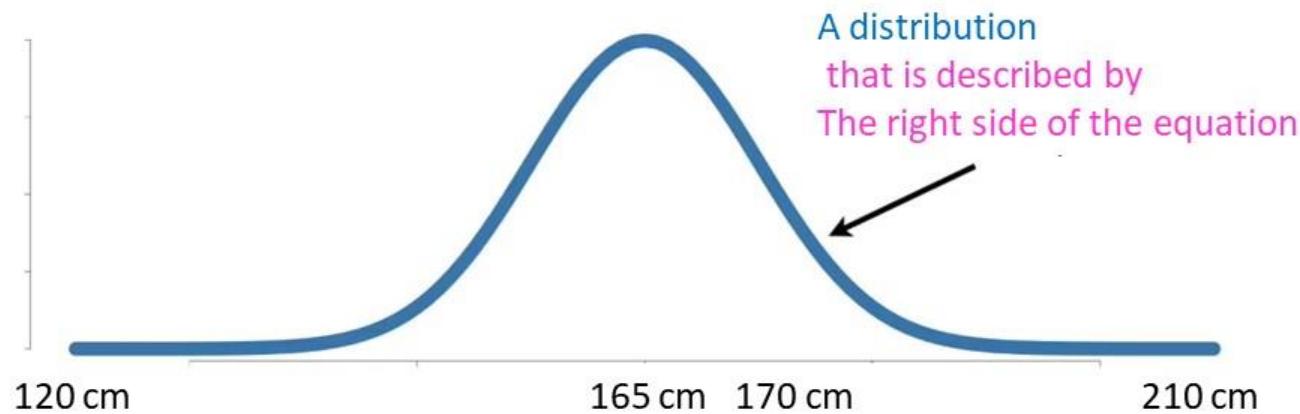
$p(\text{height} > 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$



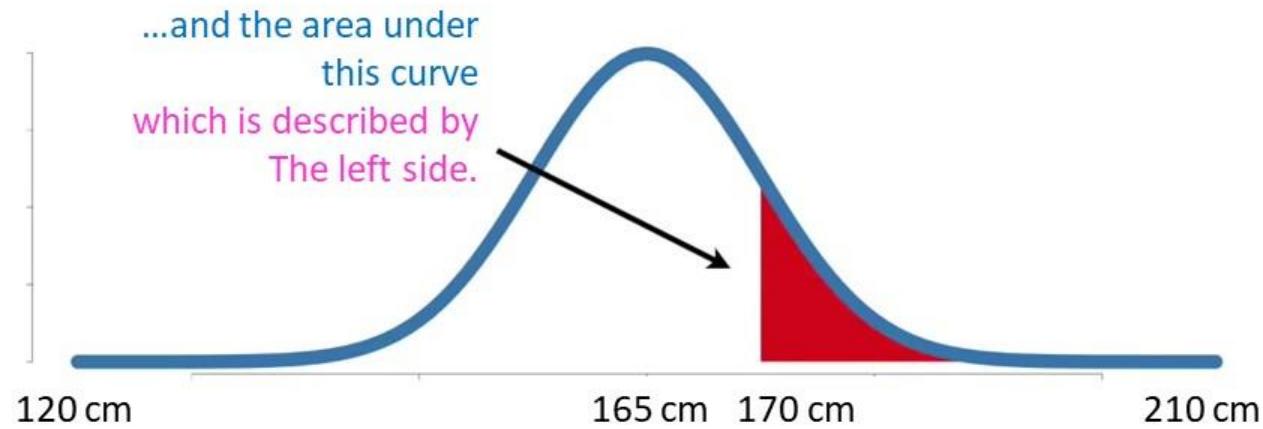
$$p(\text{height} > 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$$



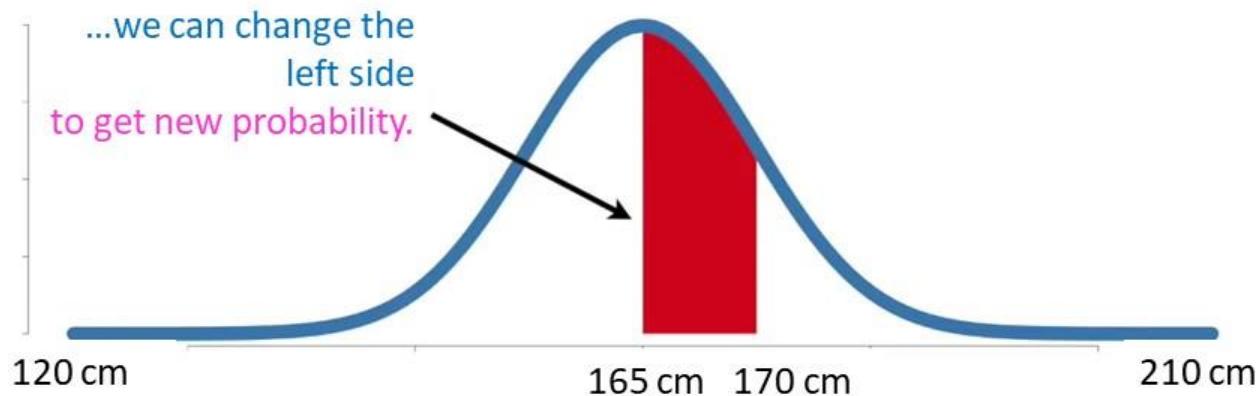
$p(\text{height} > 170 \text{ cm} \mid \text{mean} = 165 \text{ and SD} = 7)$



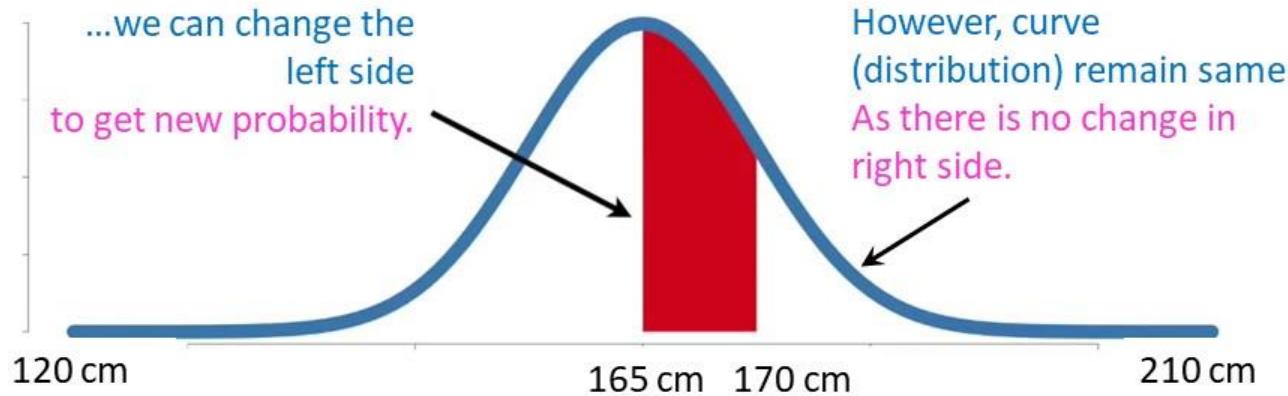
$p(\text{height} > 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$



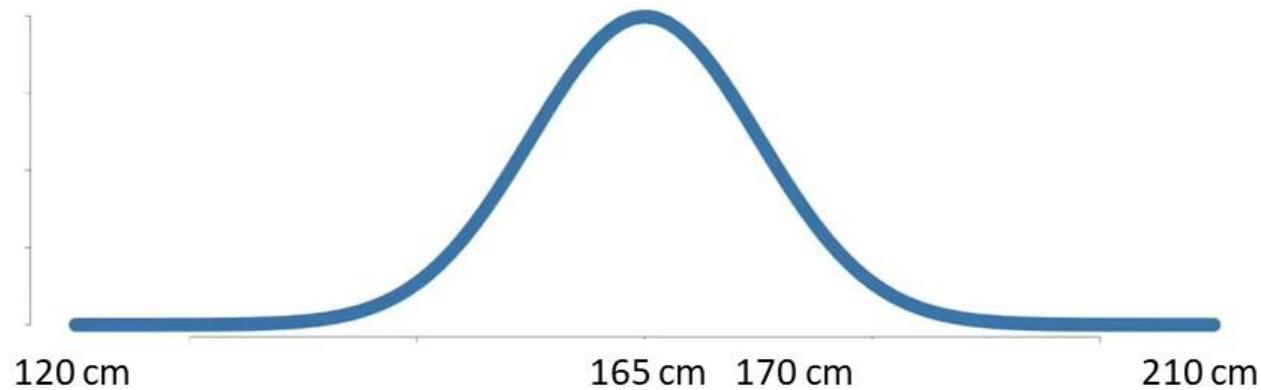
$p(\text{height between } 165 \text{ and } 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$

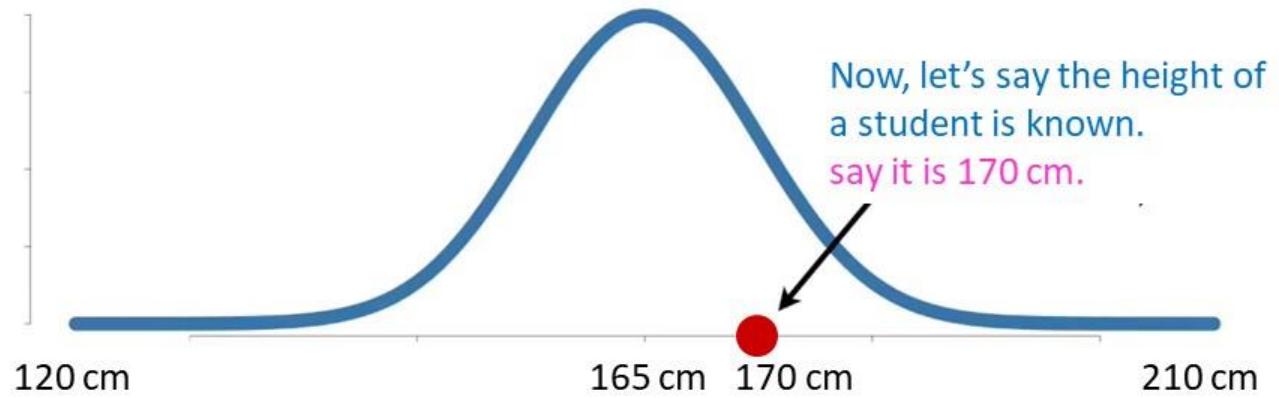


$p(\text{height between } 165 \text{ and } 170 \text{ cm} | \text{mean} = 165 \text{ and SD} = 7)$

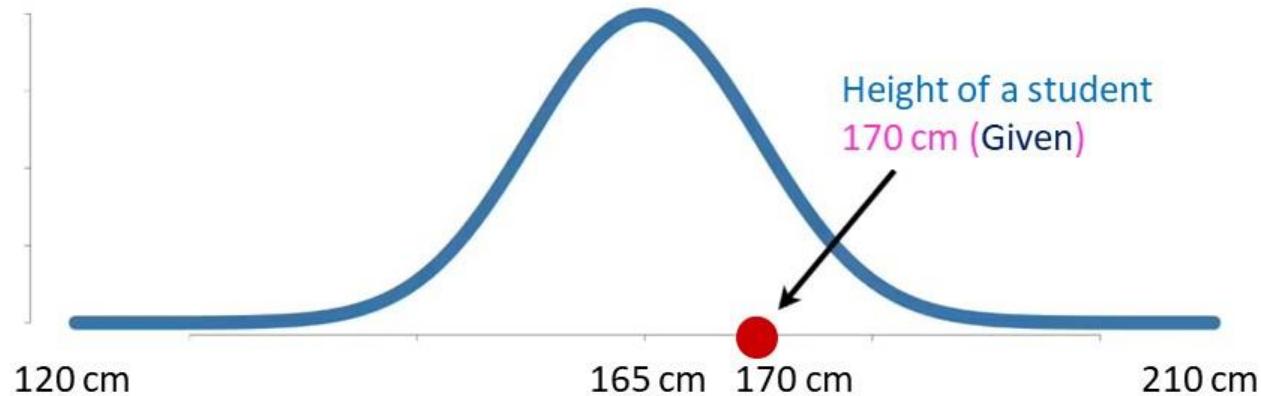


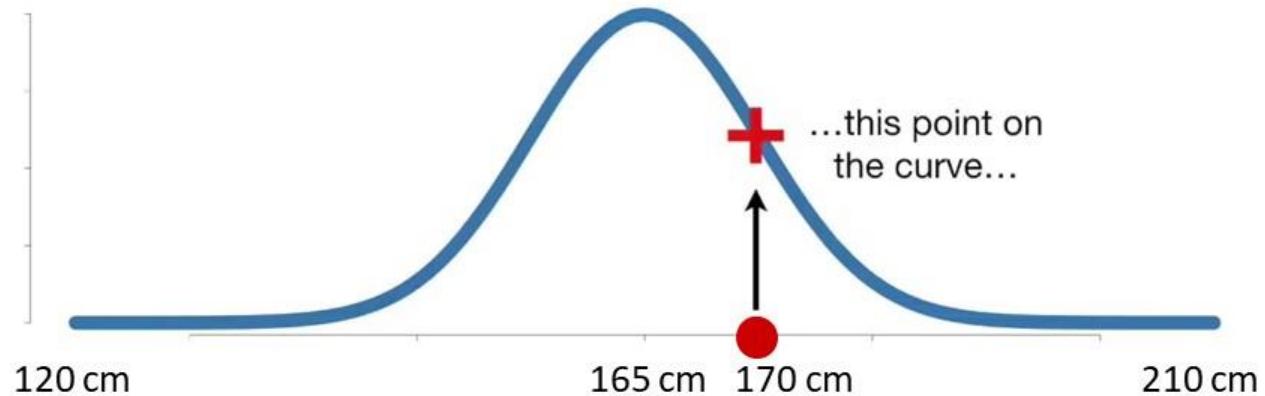
We have discussed about probability,  
Let's discuss about  
*likelihood...*



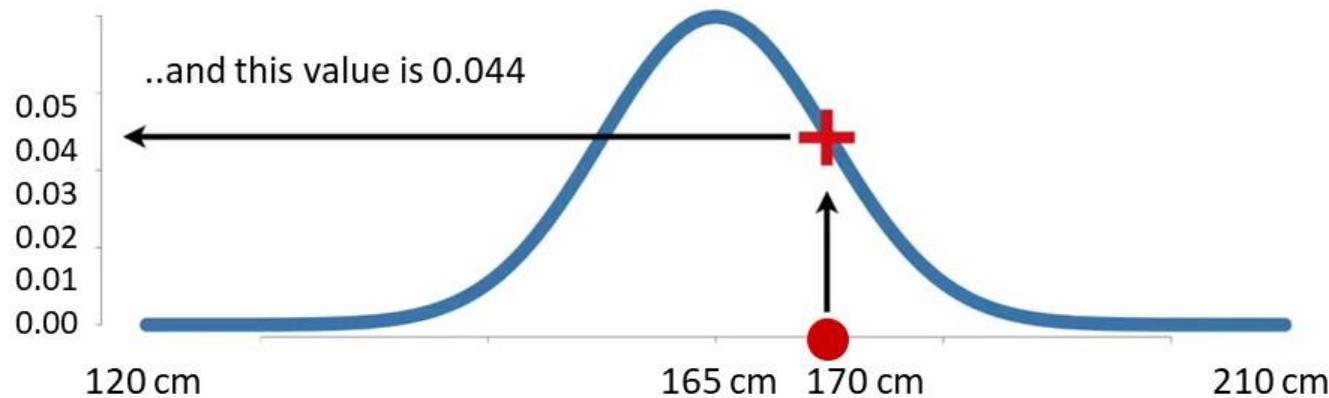


Now, we want to know  
the likelihood of this instance (170 cm)  
belonging to a distribution



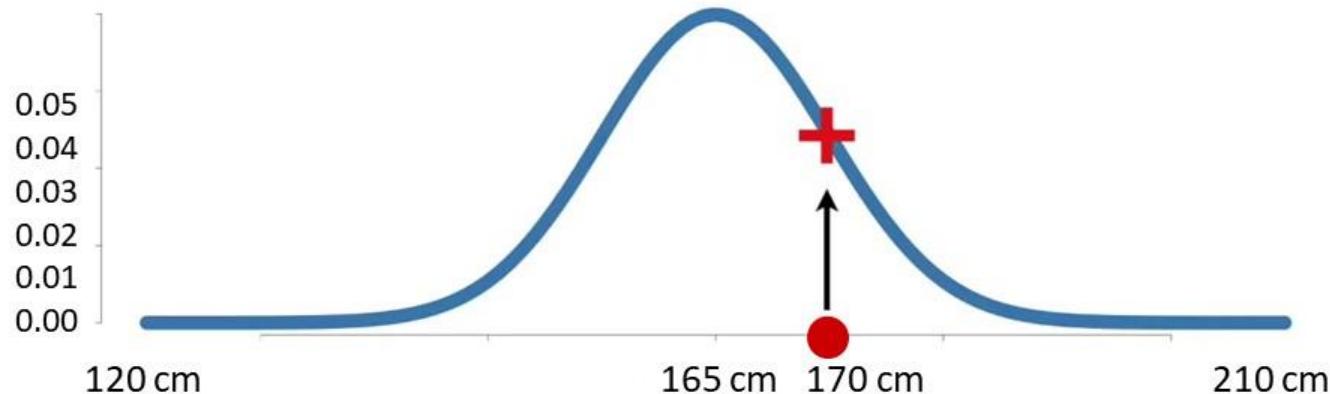


The likelihood of this instance (170 cm)  
belonging to a distribution



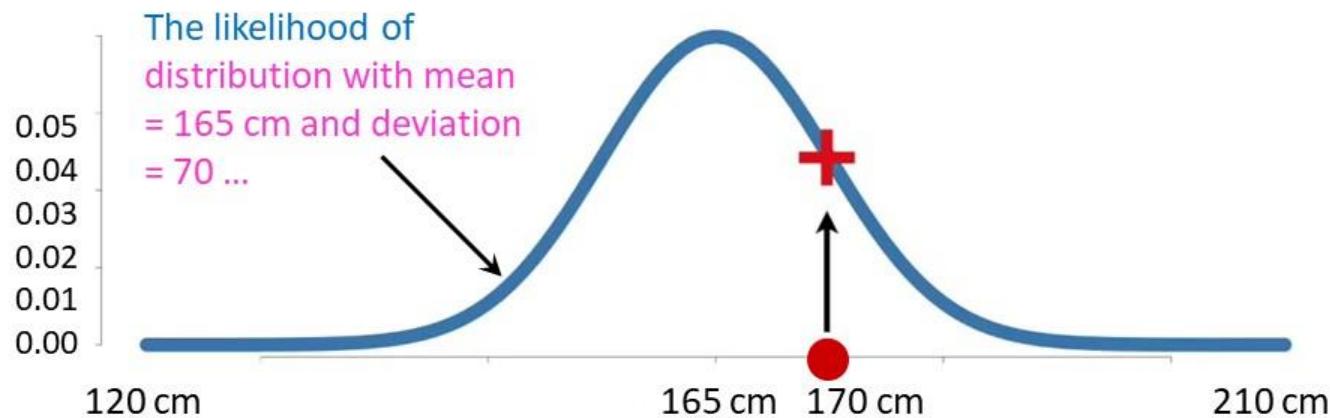
Mathematically expressed as:

$$L(\text{mean} = 165 \text{ and standard deviation} = 7 \mid \text{height of student} = 170 \text{ cm})$$



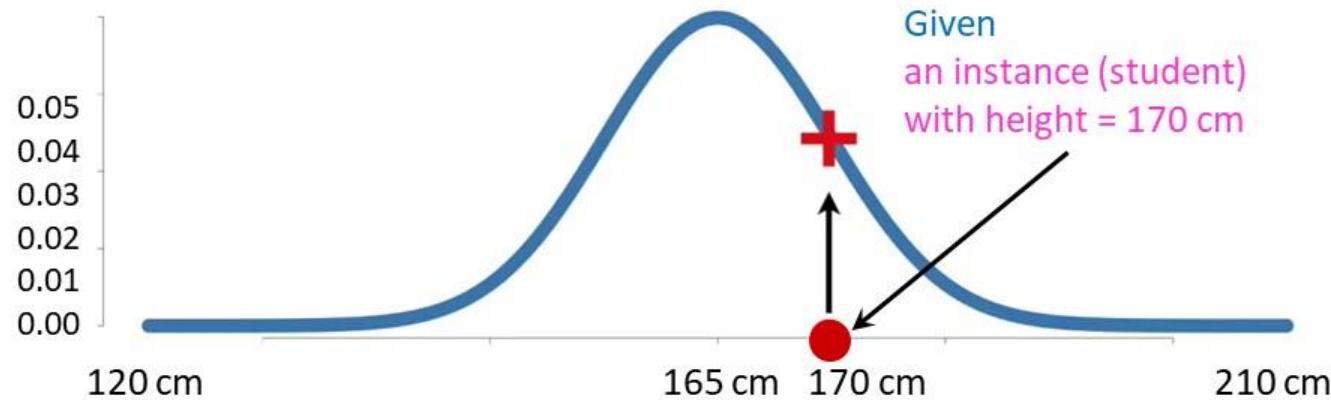
Mathematically expressed as:

$$L(\text{mean} = 165 \text{ and standard deviation} = 7 \mid \text{height of student} = 170 \text{ cm})$$

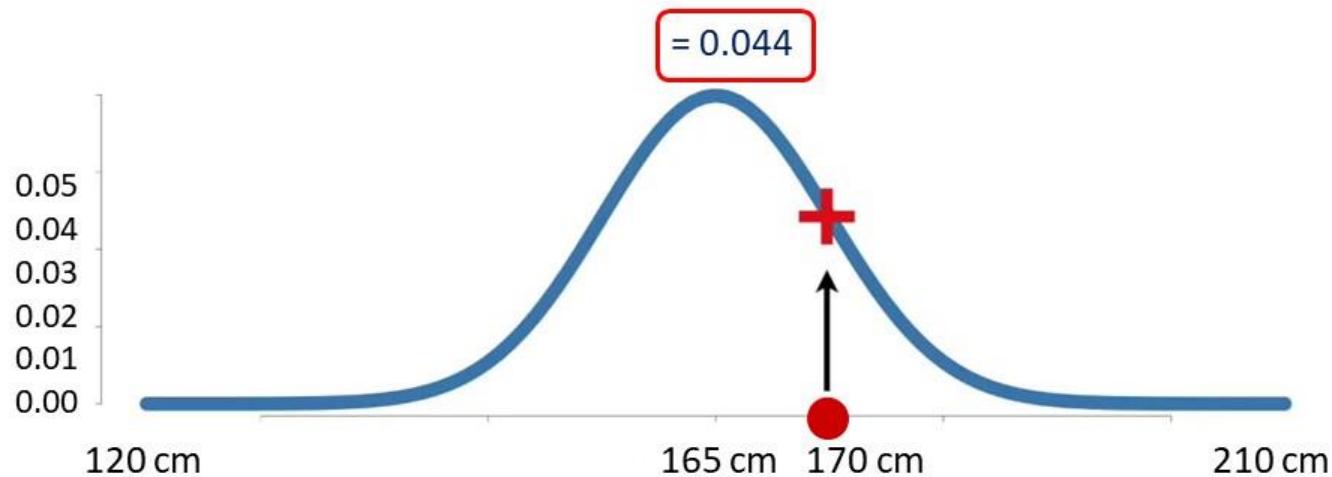


Mathematically expressed as:

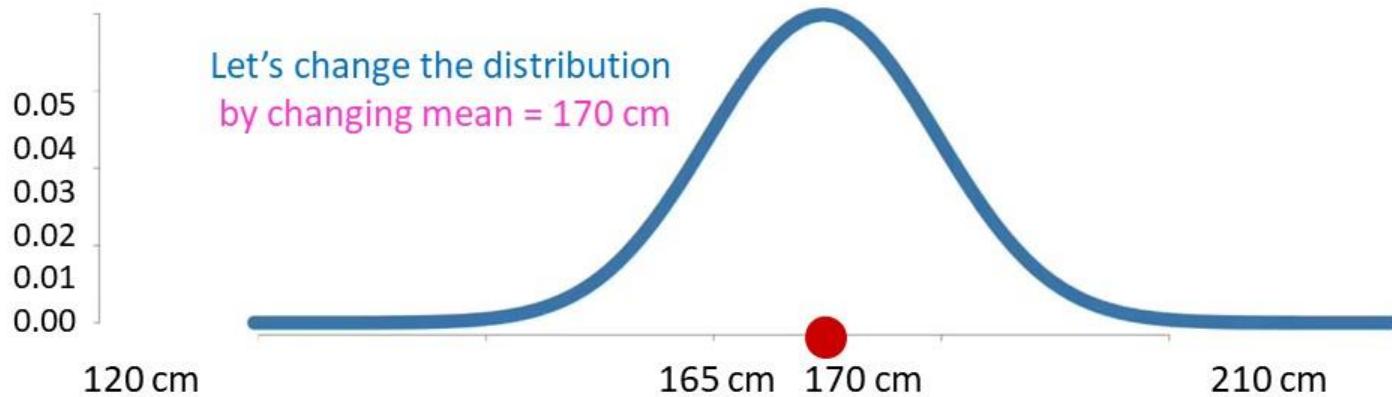
$$L(\text{mean} = 165 \text{ and standard deviation} = 7 \mid \text{height of student} = 170 \text{ cm})$$



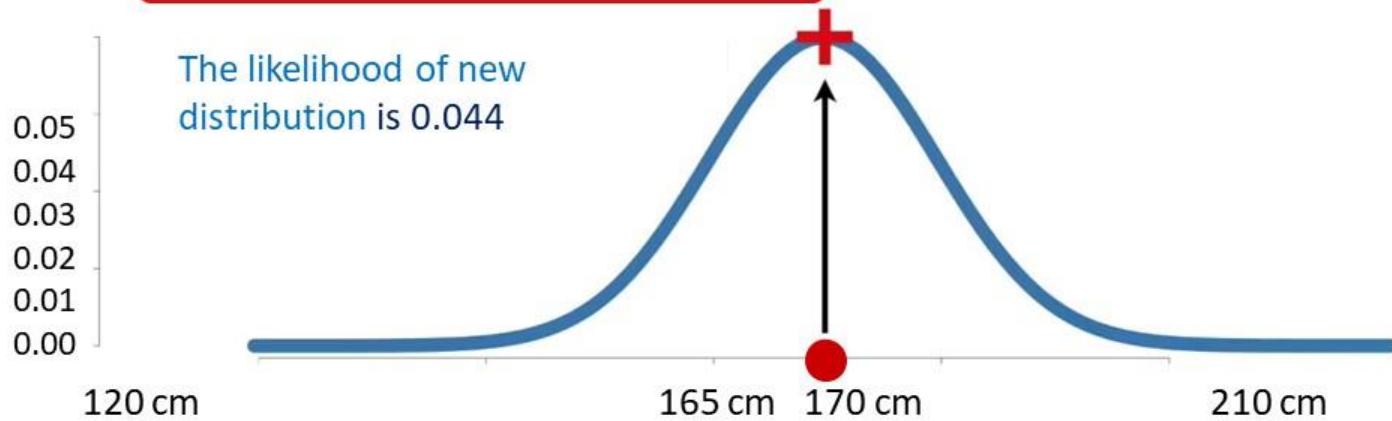
$L(\text{mean} = 165 \text{ and standard deviation} = 7 \mid \text{height of student} = 170 \text{ cm})$



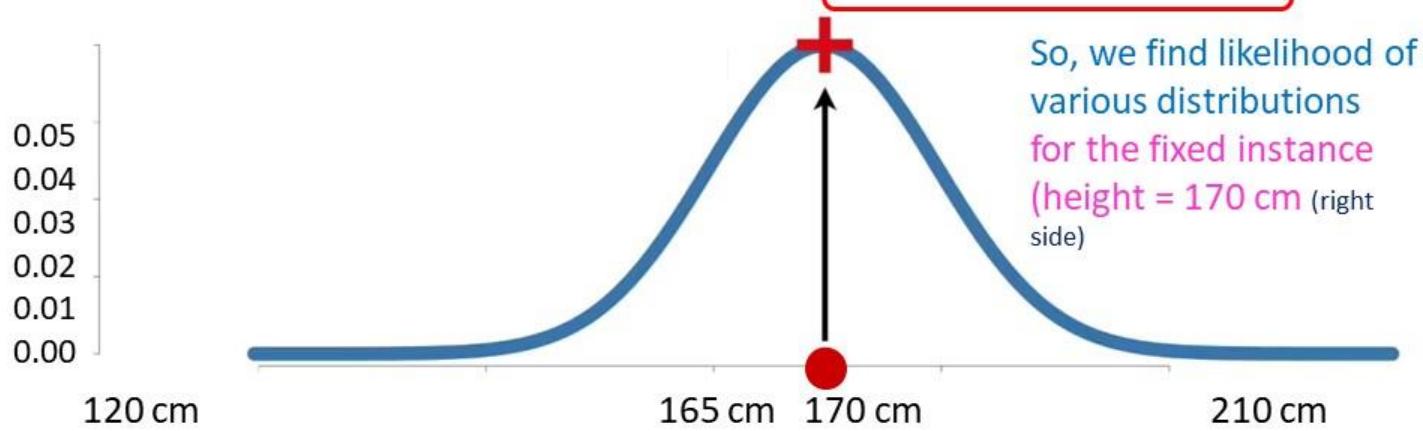
$L(\text{mean} = 170 \text{ and standard deviation} = 7 | \text{height of student} = 170 \text{ cm})$



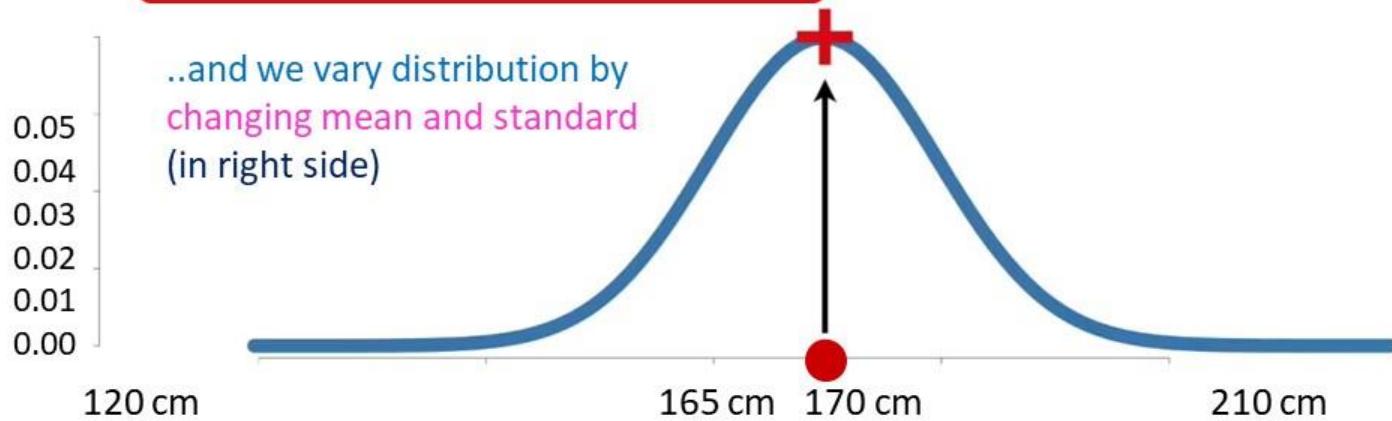
$L(\text{mean} = 170 \text{ and standard deviation} = 7 | \text{height of student} = 170 \text{ cm})$



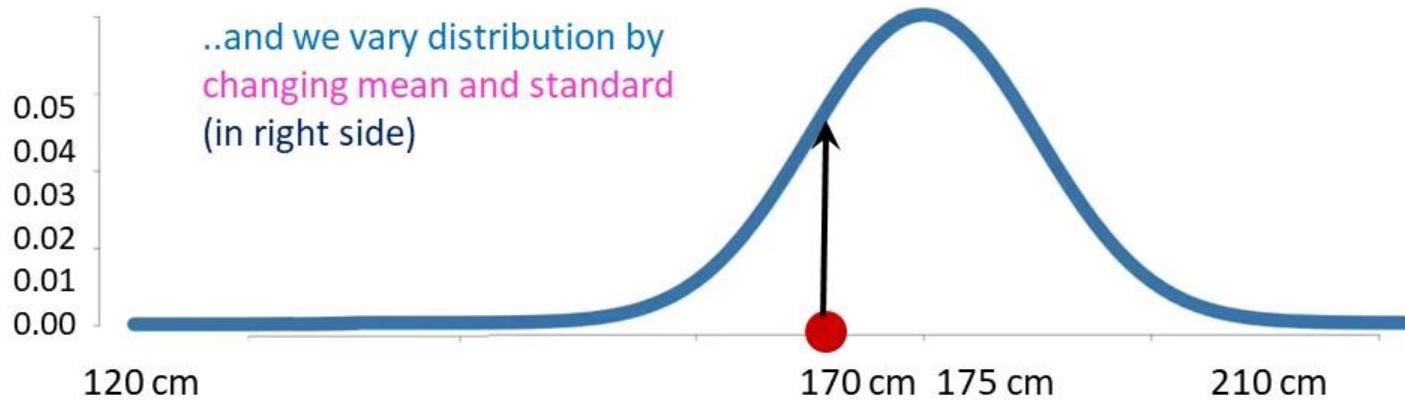
$L(\text{mean} = 170 \text{ and standard deviation} = 7 | \text{height of student} = 170 \text{ cm})$



$L(\text{mean} = 170 \text{ and standard deviation} = 7 | \text{height of student} = 170 \text{ cm})$



$L(\text{mean} = 175 \text{ and standard deviation} = 7 | \text{height of student} = 170 \text{ cm})$

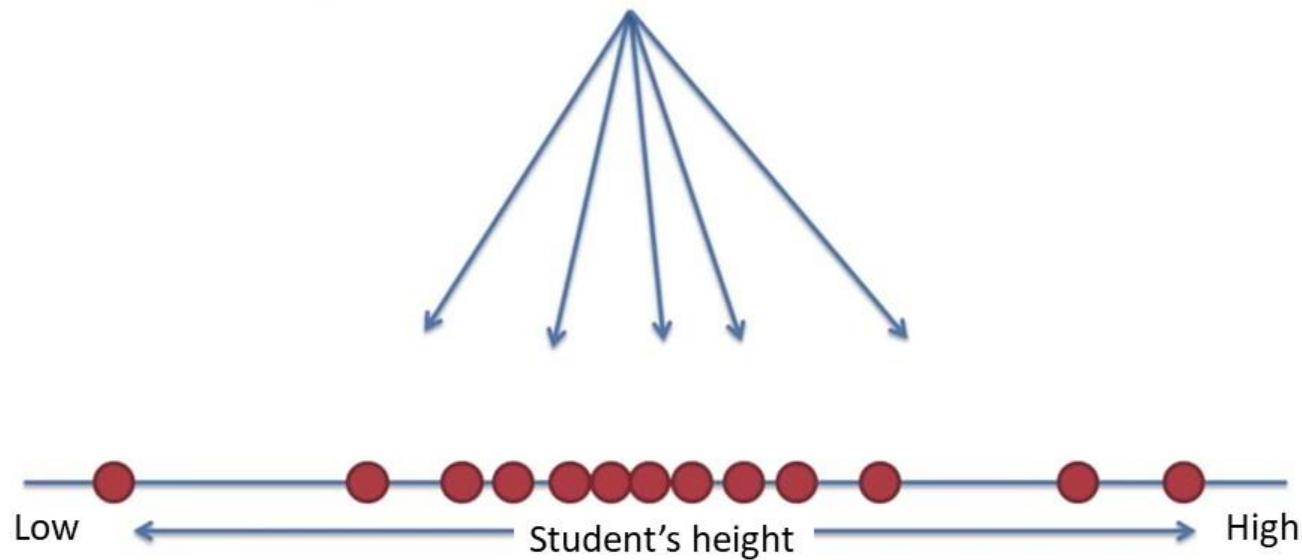


---

# **Maximum Likelihood Estimation**

---

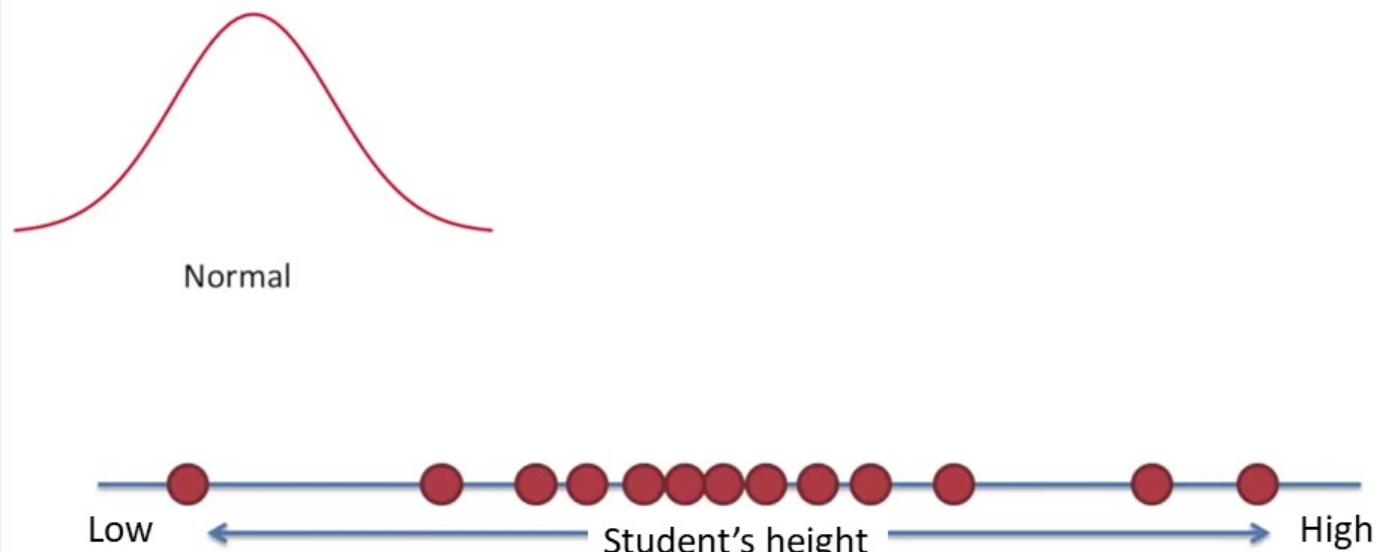
Let's say we measured height of a set of students



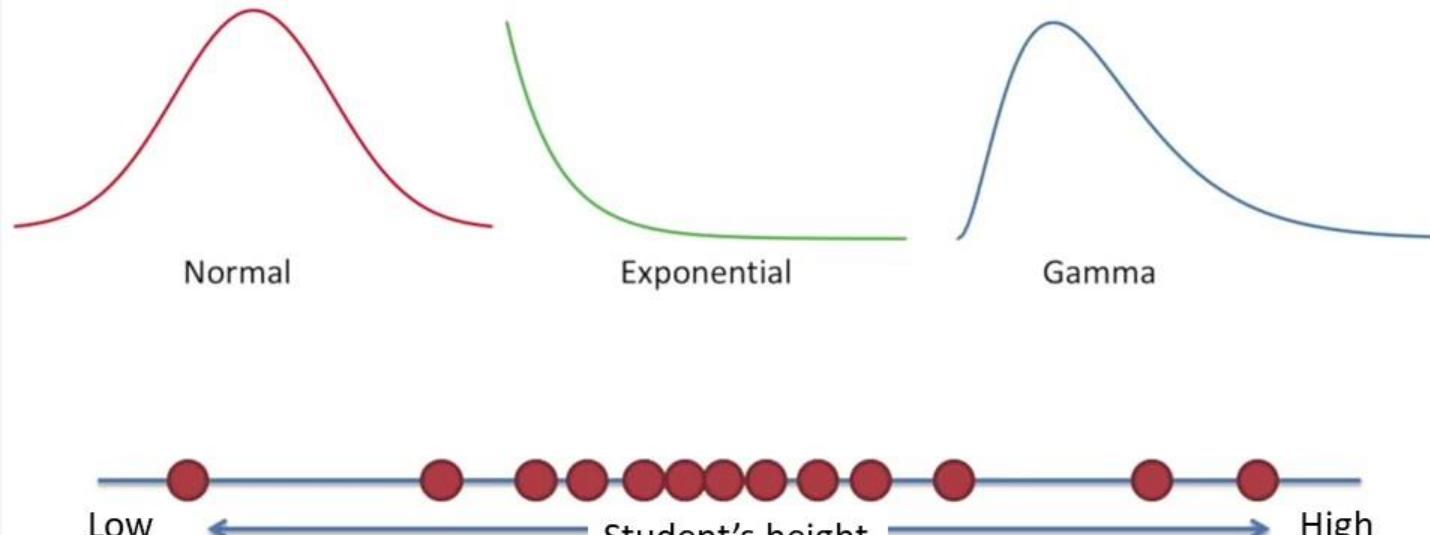
The goal of maximum likelihood is to find the optimal way to fit a distribution to the data



There are lots of different types of distributions  
for different types of data...



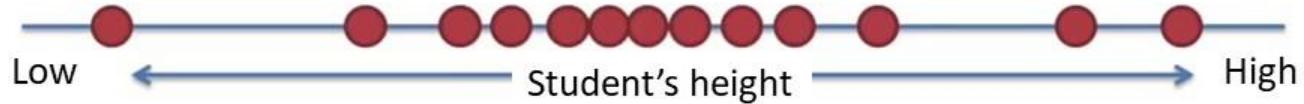
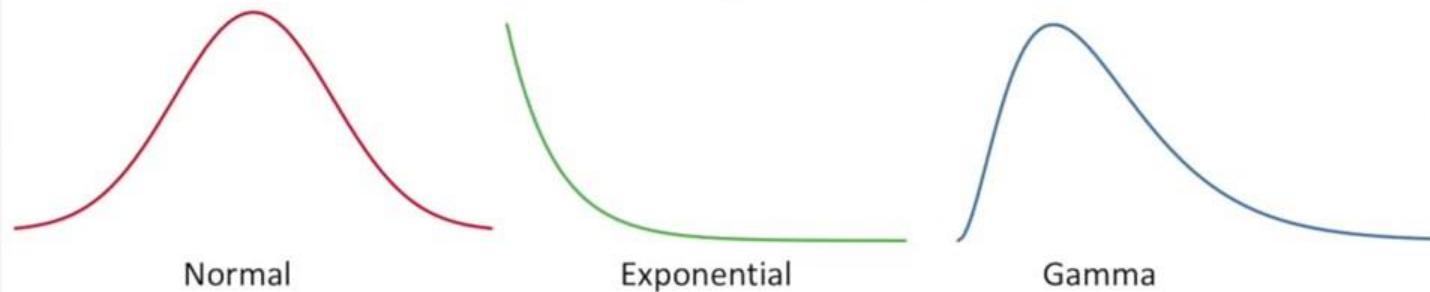
There are lots of different types of distributions  
for different types of data...



The reason we want to fit a distribution to the data is

Distribution is much more than data

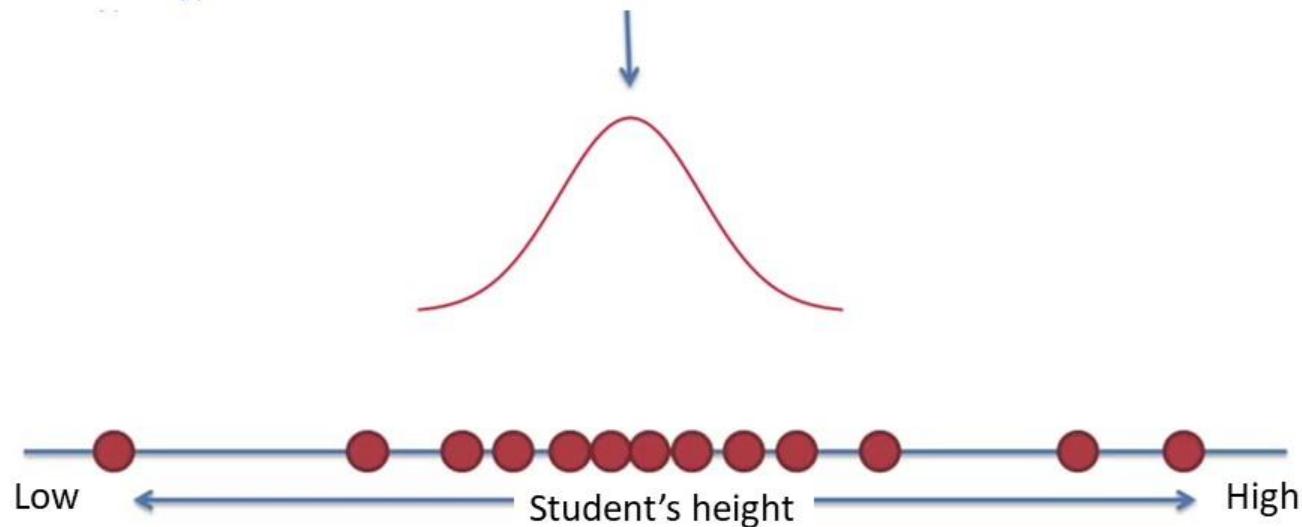
As distribution represent whole pollution,  
while data is a sample of the population



In this case, we may assume that heights are normally distributed

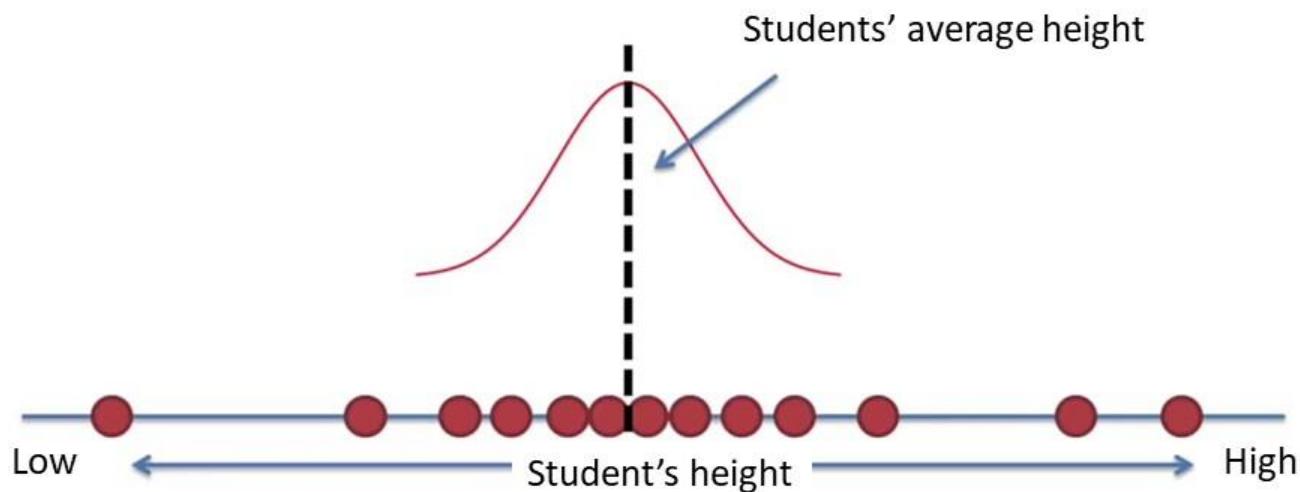


That means, it has came from  
this type of distribution



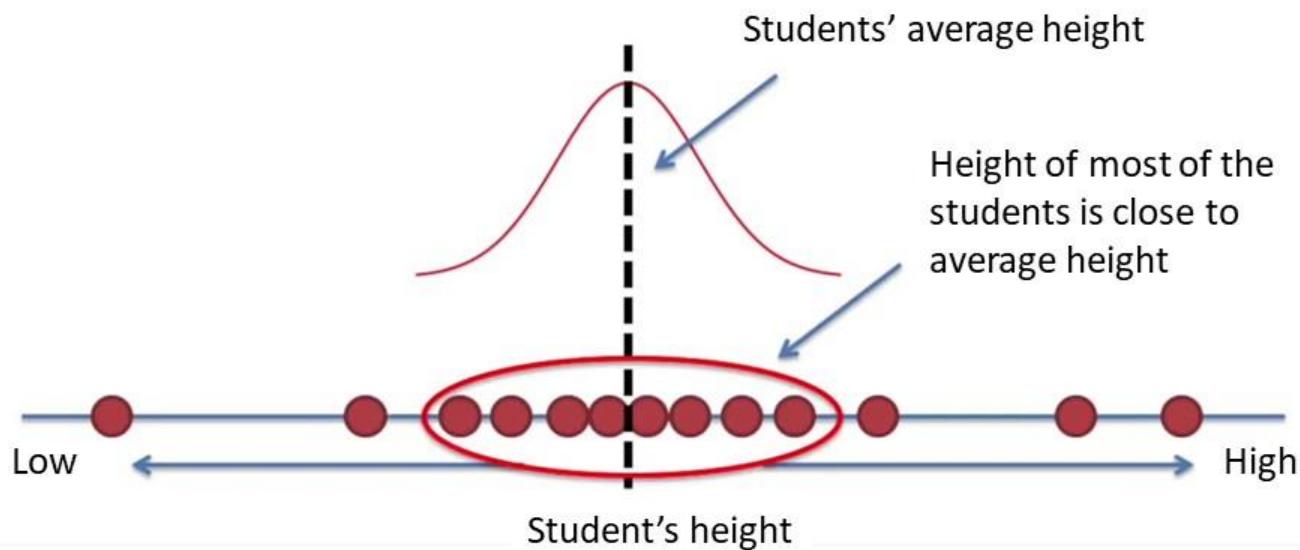
“Normally distributed” means:

- 1) We expect most of the students have height close to mean (average)



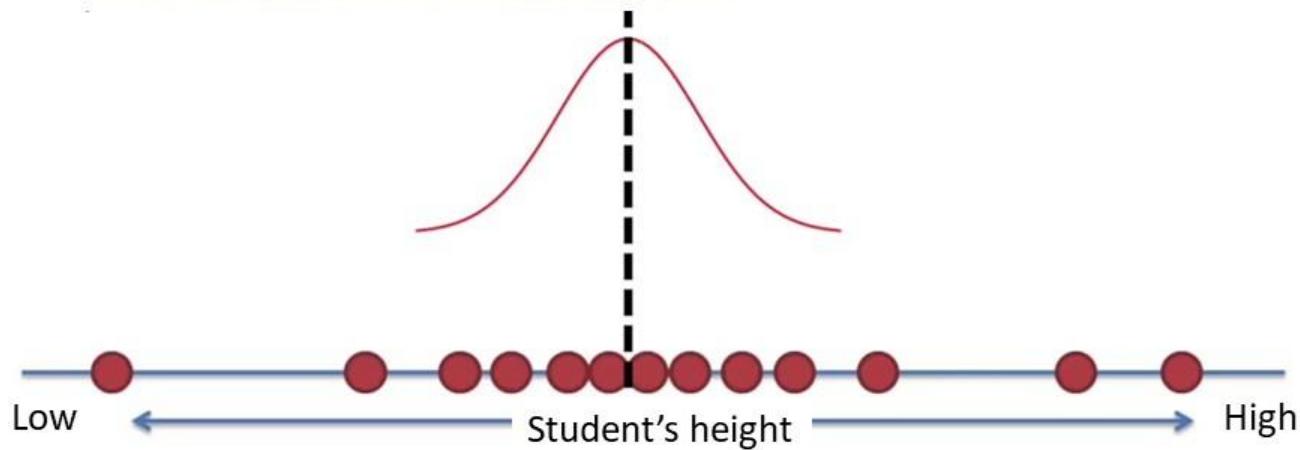
“Normally distributed” means:

- 1) We expect most of the students have height close to mean (average)



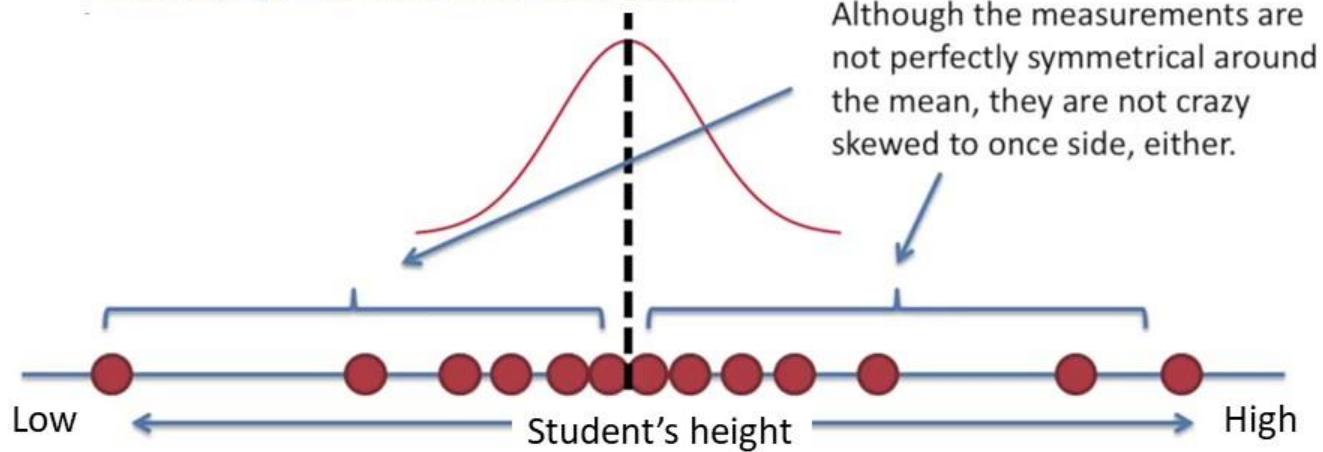
“Normally distributed” means:

- 1) We expect most of the students have height (measurements) close to mean (average)
- 2) We expect the measurements to be relatively symmetrical around the mean

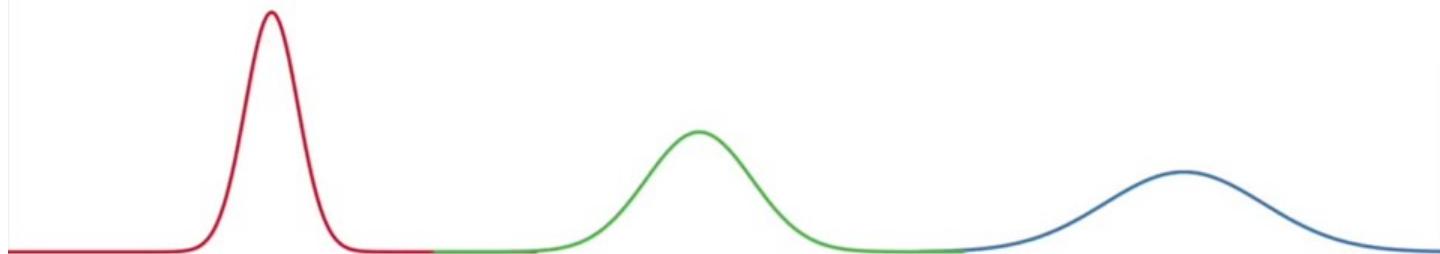


“Normally distributed” means:

- 1) We expect most of the students have height (measurements) close to mean (average)
- 2) We expect the measurements to be relatively symmetrical around the mean

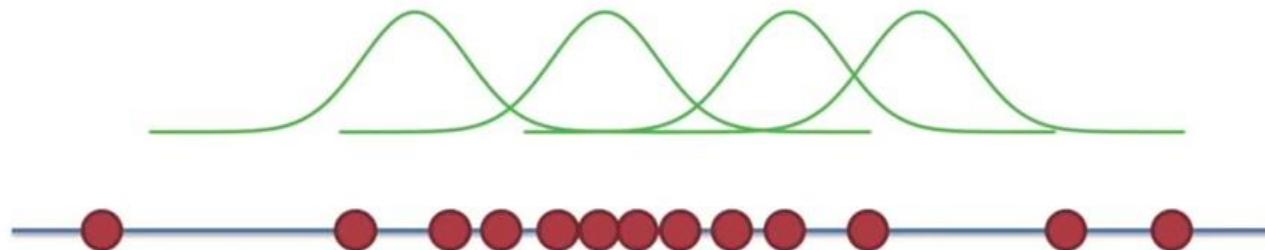


Normal distribution may have different kinds of shapes and sizes...

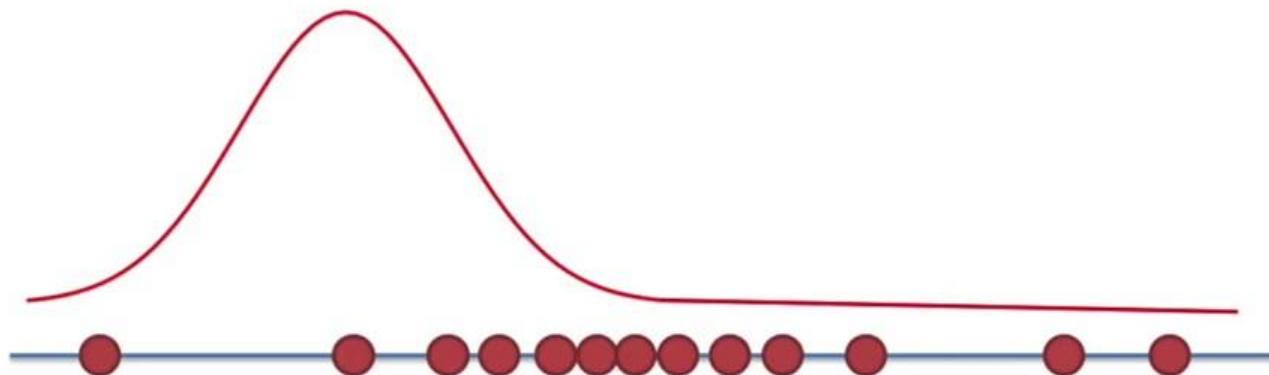


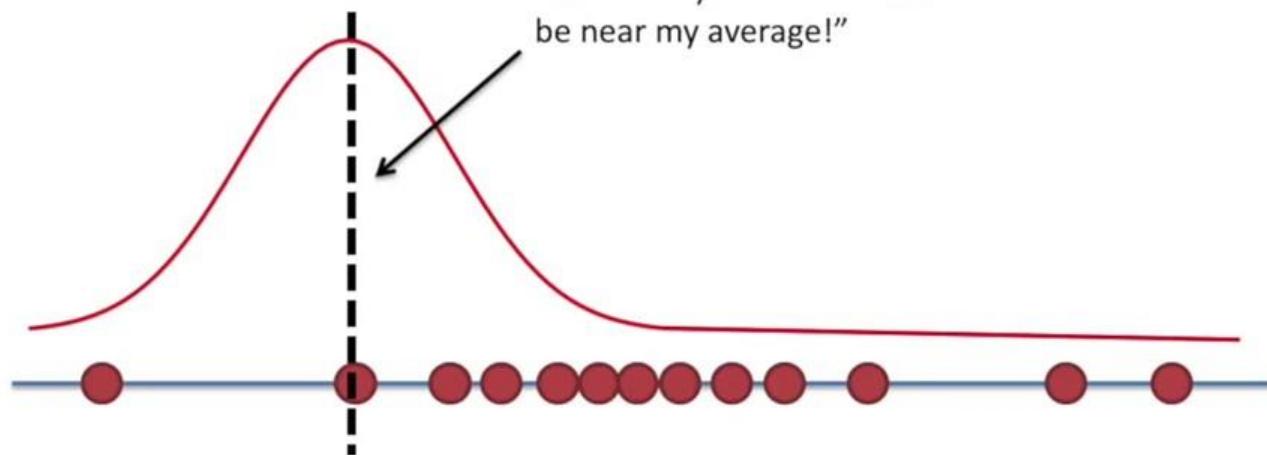
Once we decide the shape of the distribution, we need to figure out its center (position) In this case, we may assume that heights are

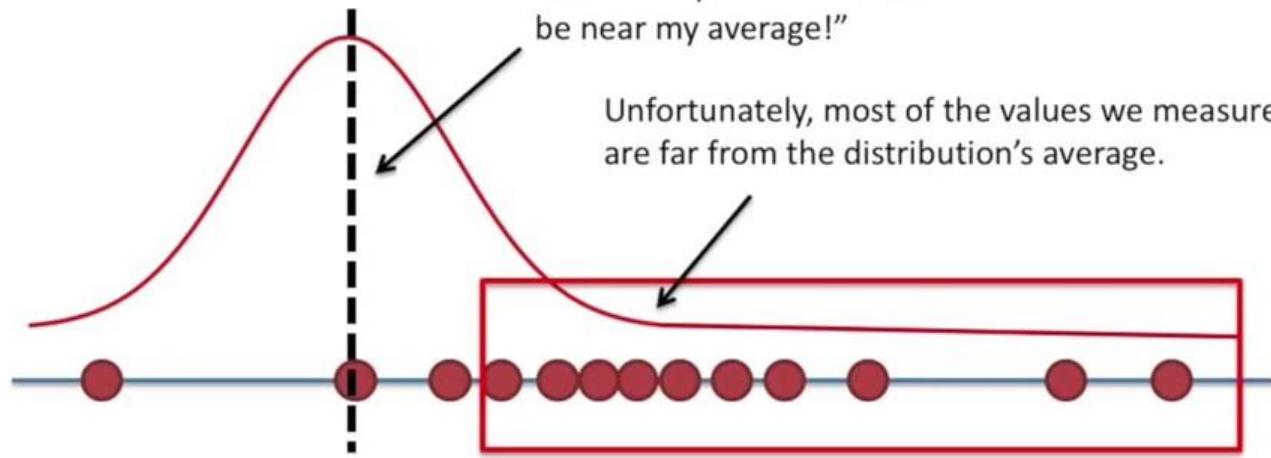
Is one location is “better” than another?



Before we get too technical, let's just pick any old normal distribution and see how well it fits the data.

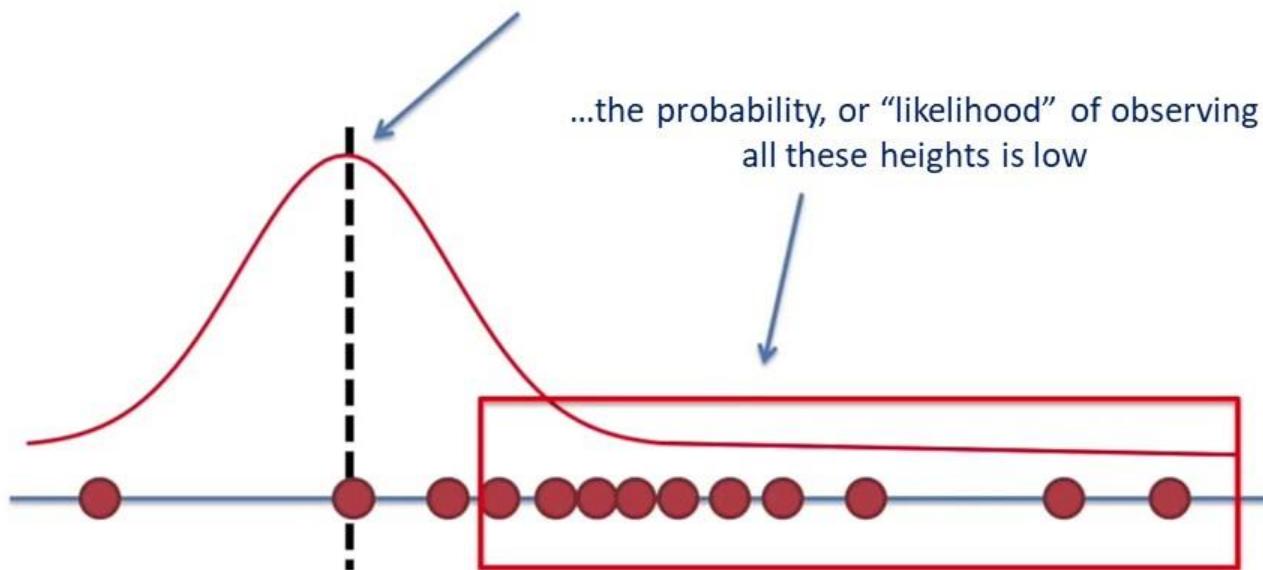




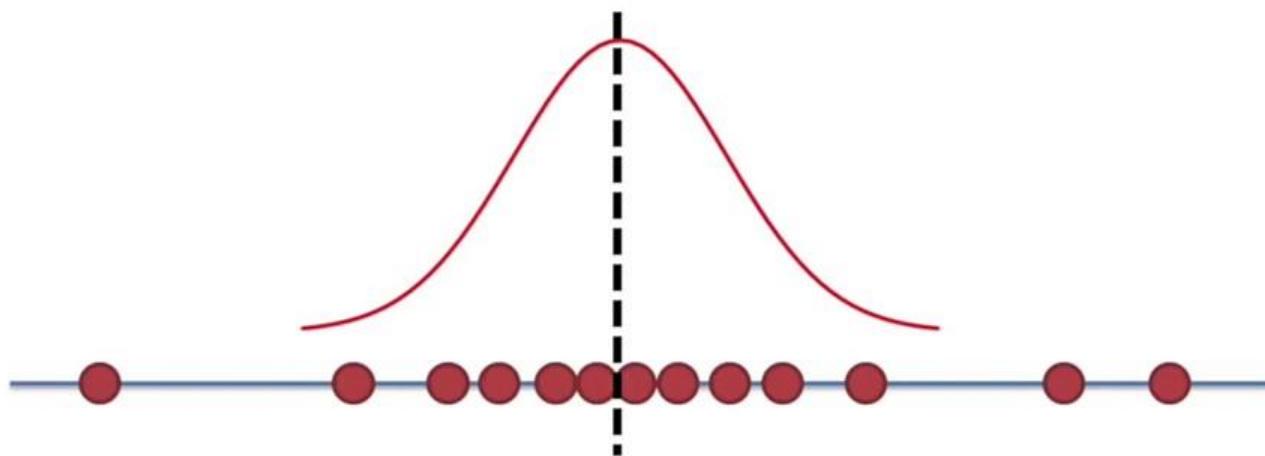


According to a normal distribution with a mean value over here...

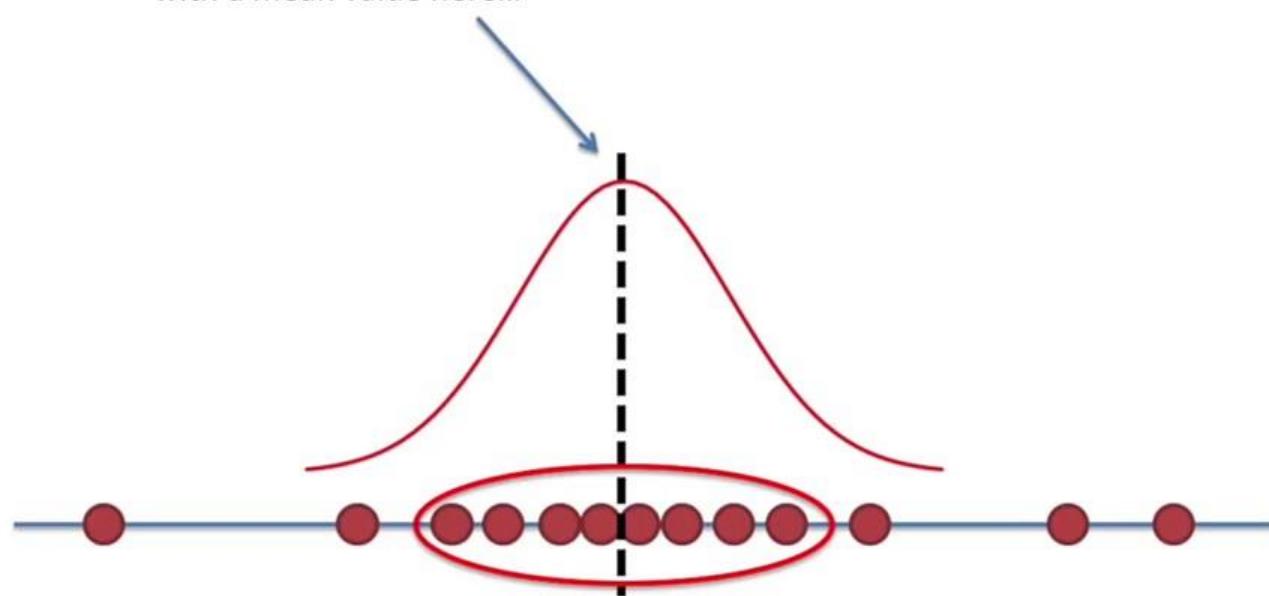
...the probability, or “likelihood” of observing all these heights is low



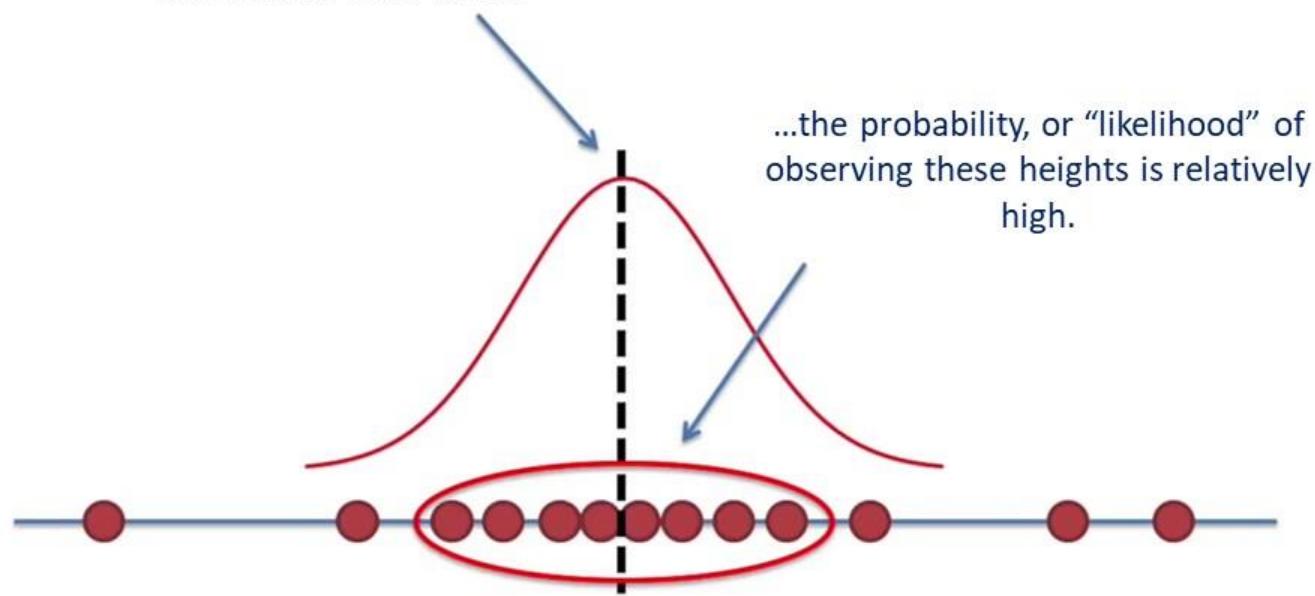
What if we shift the normal distribution over, so that its mean was the same as the average weight?



According to a normal distribution  
with a mean value here...

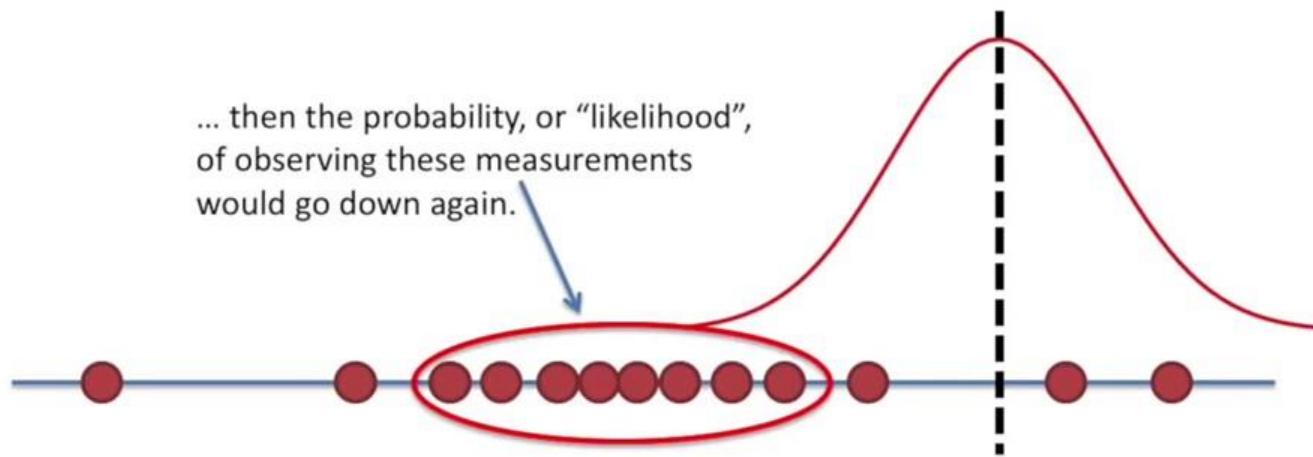


According to a normal distribution  
with a mean value here...



If we kept shifting the normal distribution over...

... then the probability, or “likelihood”,  
of observing these measurements  
would go down again.

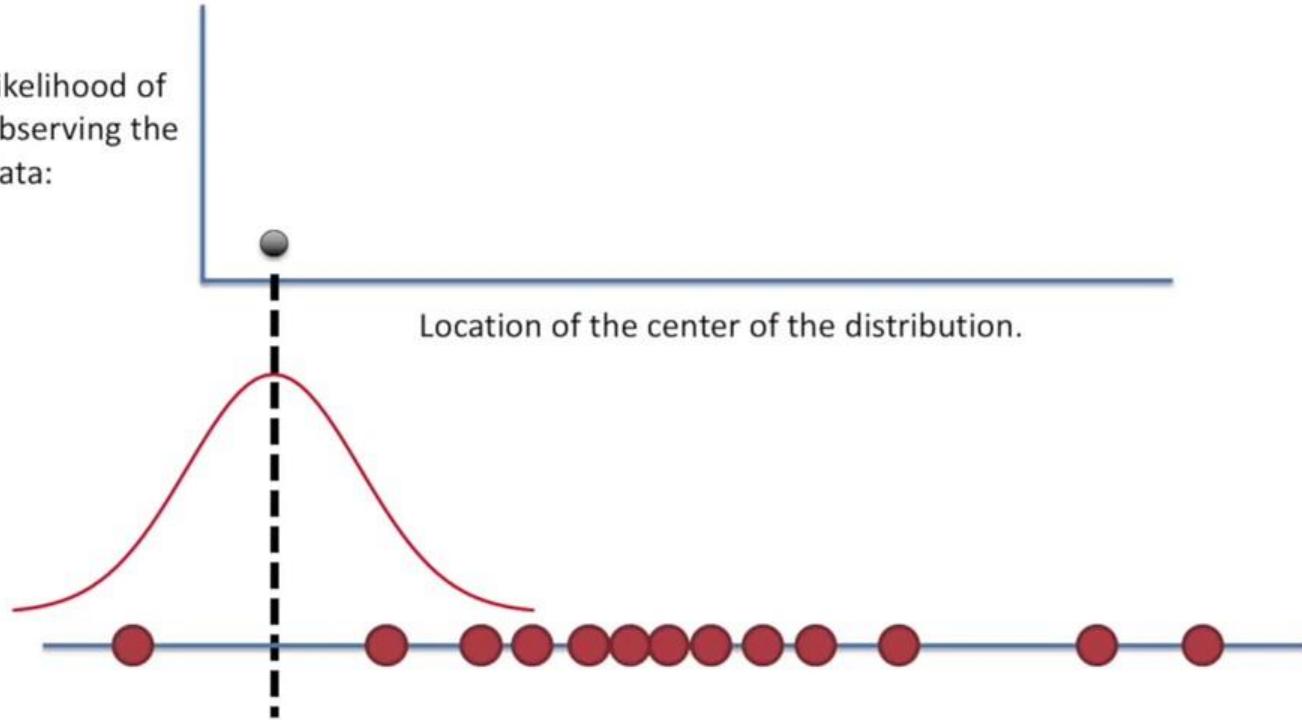


Likelihood of observing the data:

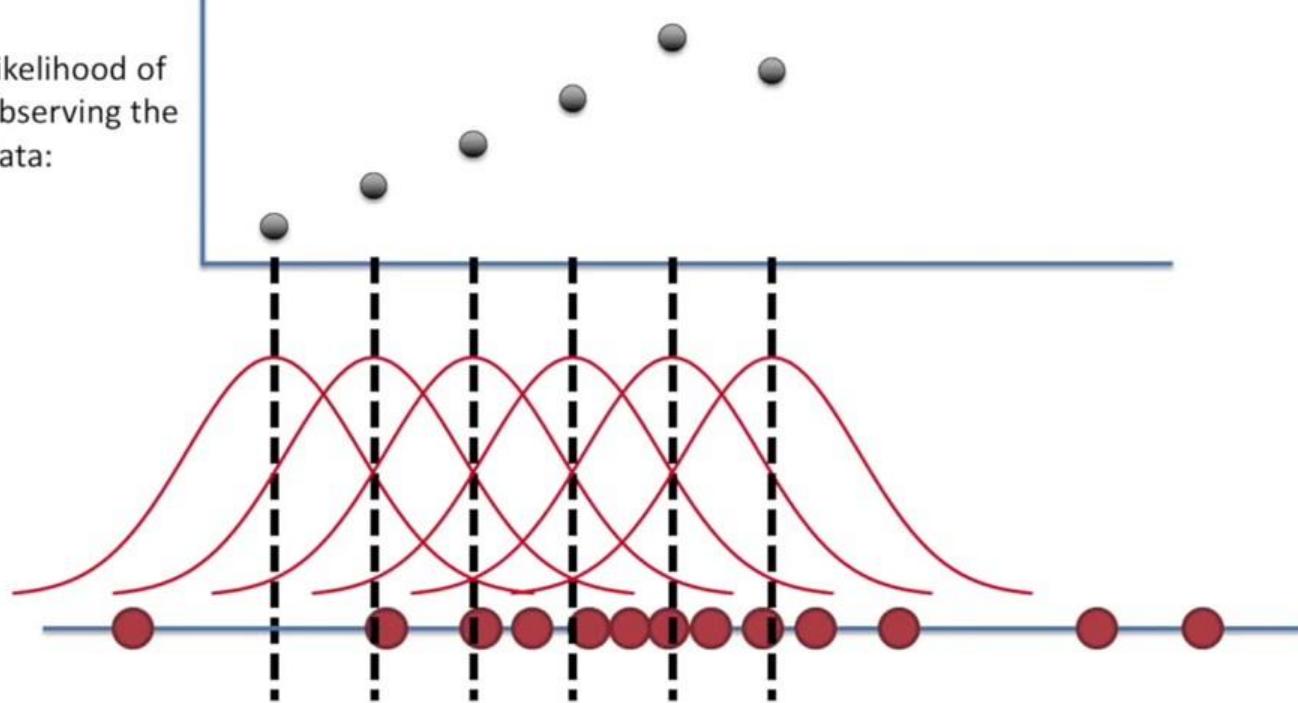
Location of the center of the distribution.



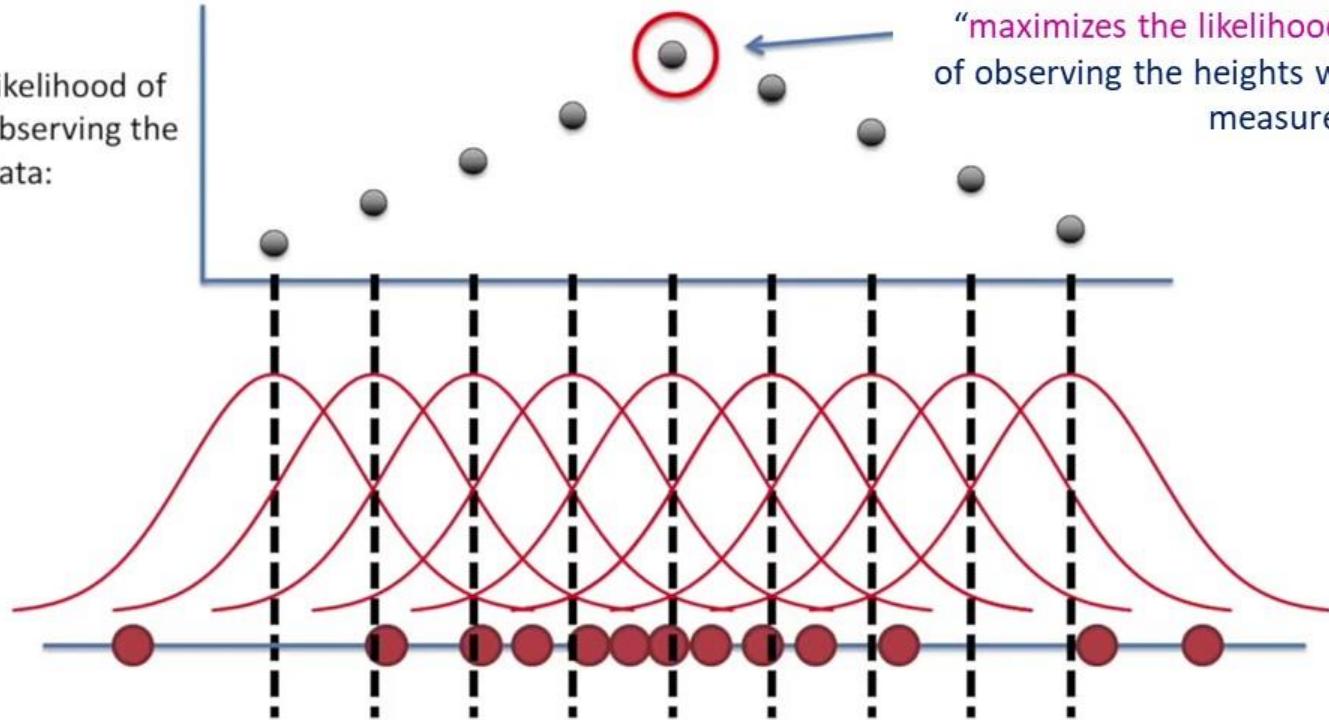
Likelihood of observing the data:

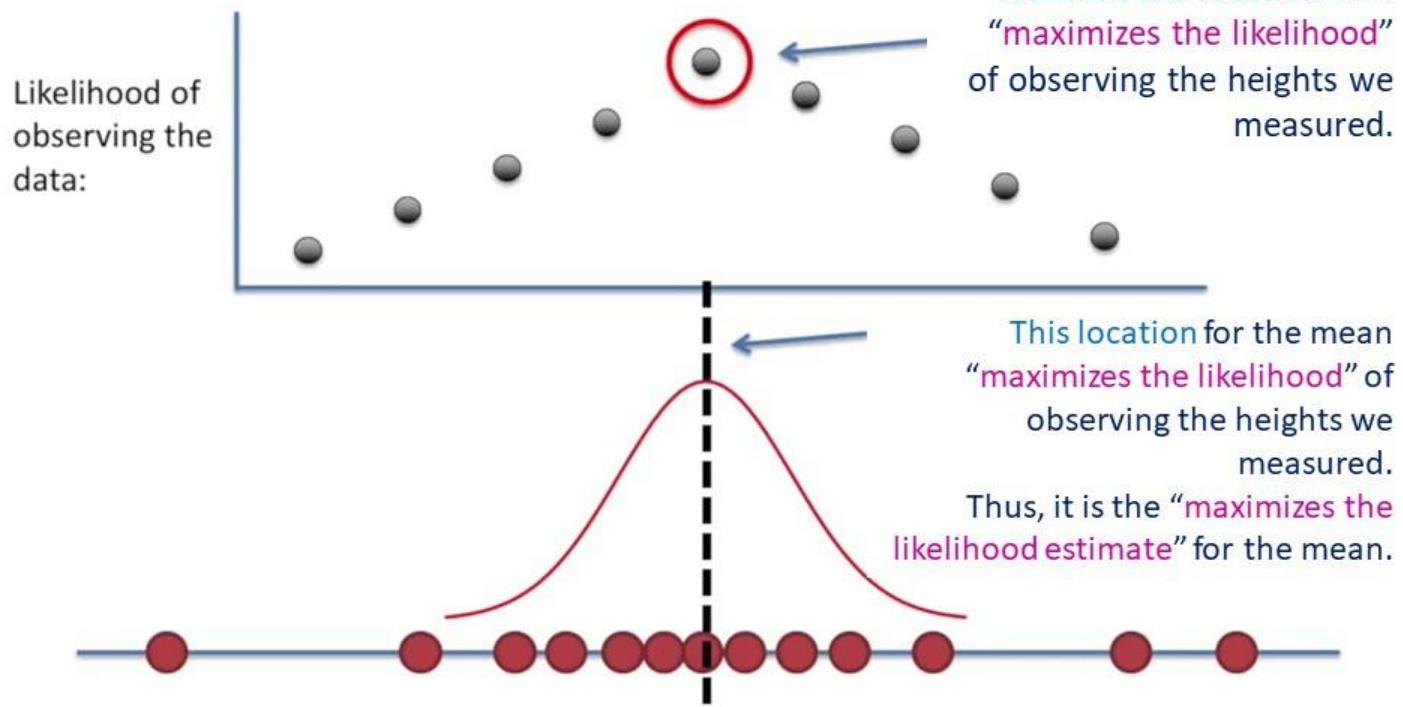


Likelihood of observing the data:

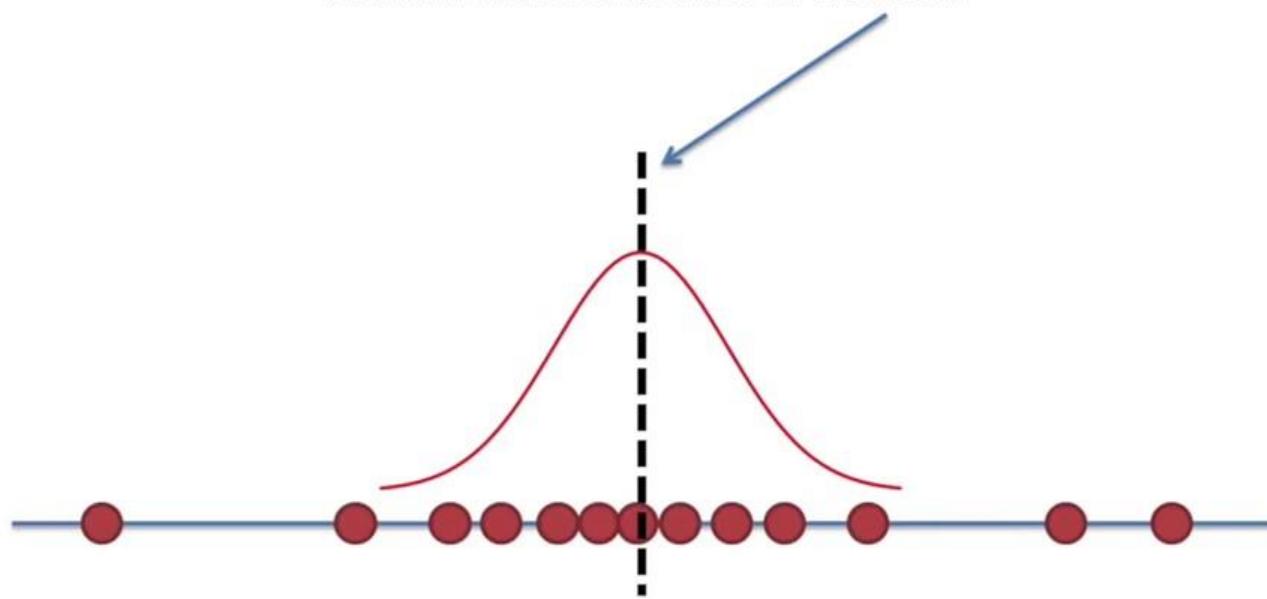


Likelihood of observing the data:

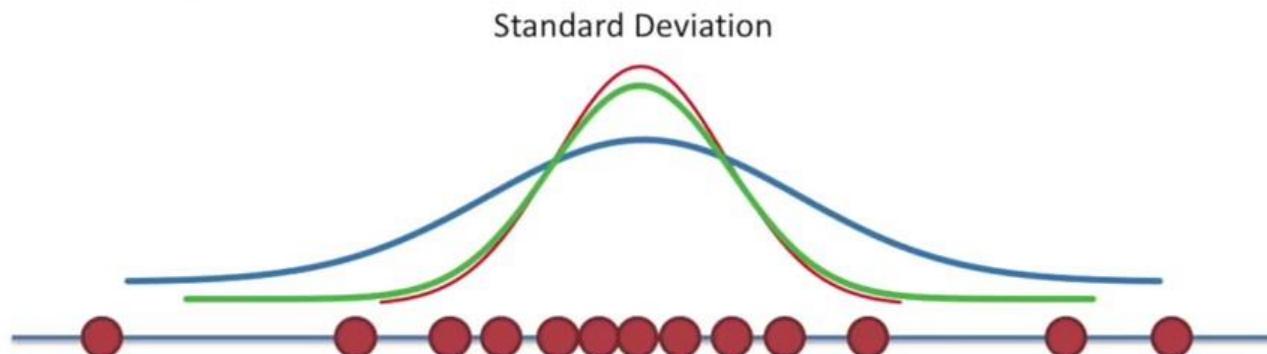




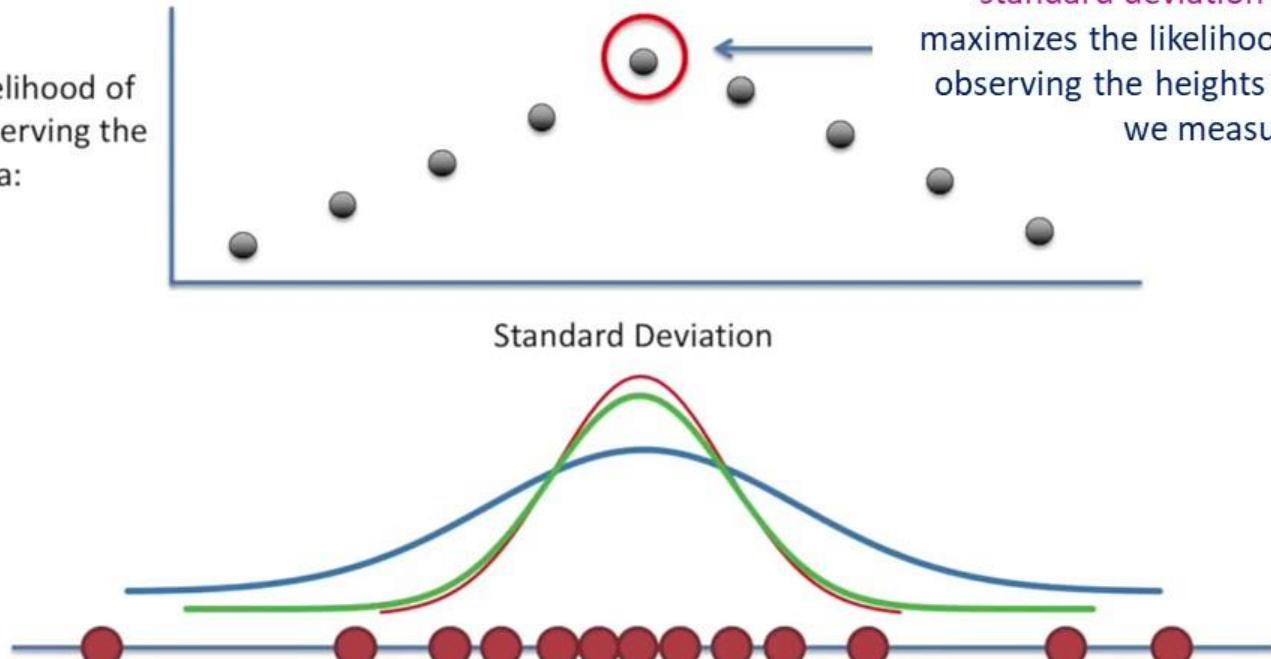
Great! Now we have figured out the maximum likelihood estimate for the mean!



Likelihood of observing the data:

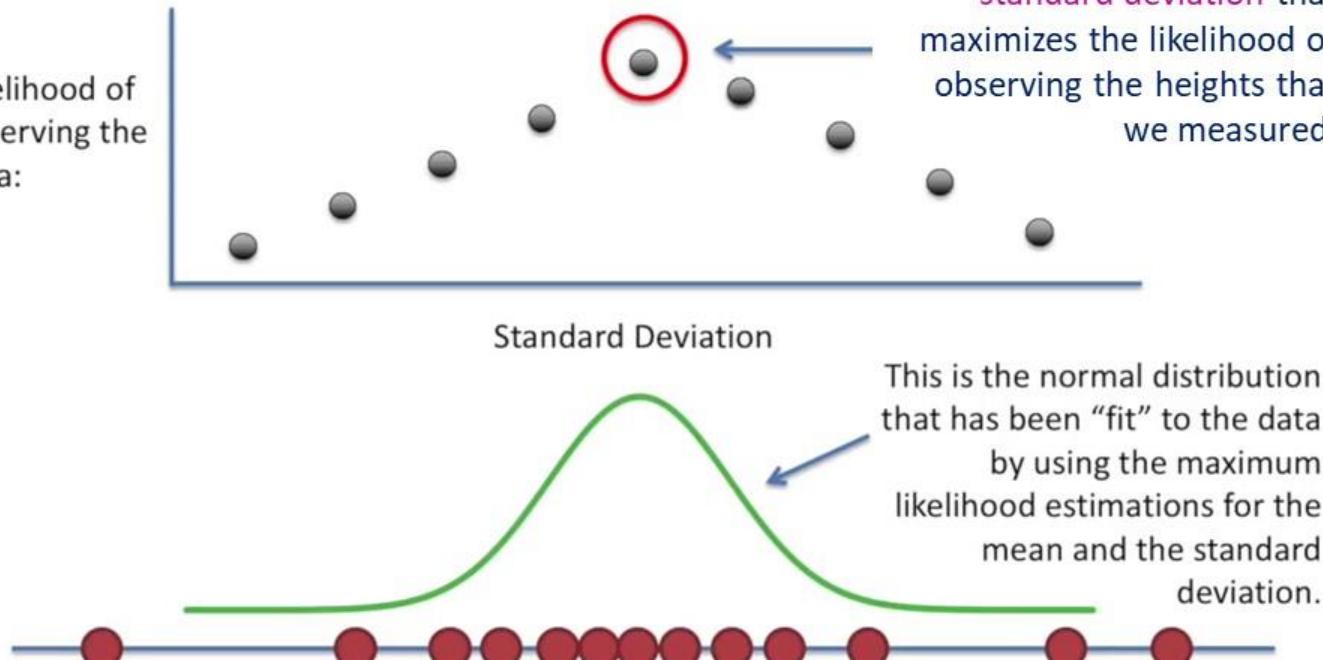


Likelihood of observing the data:

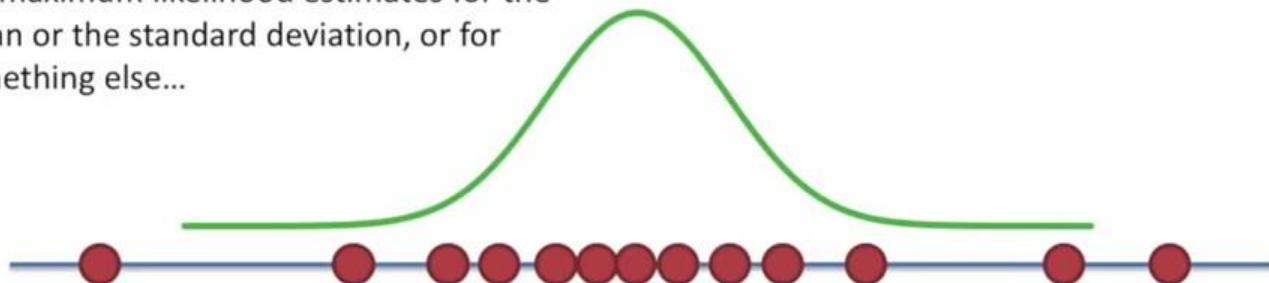


Now we've found the **standard deviation** that maximizes the likelihood of observing the heights that we measured.

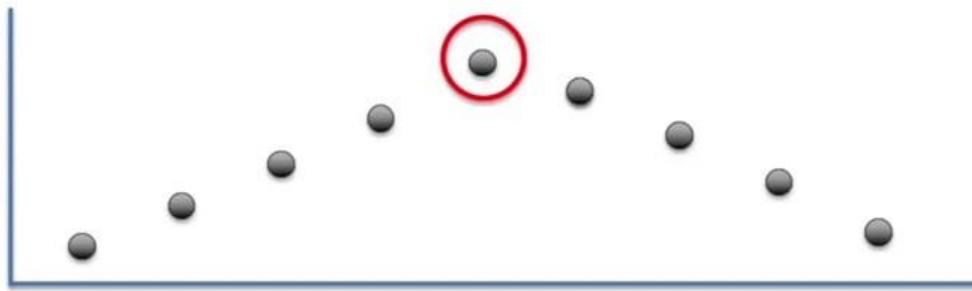
Likelihood of observing the data:



Now when someone says that they have the maximum likelihood estimates for the mean or the standard deviation, or for something else...

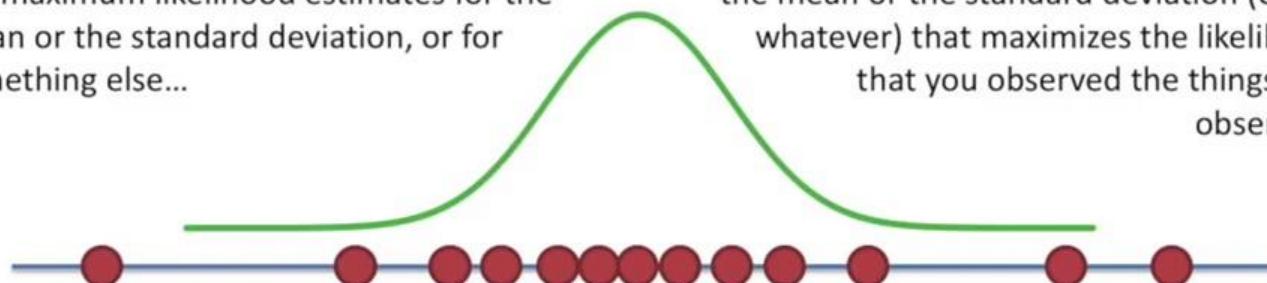


Likelihood of observing the data:



Now when someone says that they have the maximum likelihood estimates for the mean or the standard deviation, or for something else...

... you know that they found the value for the mean or the standard deviation (or for whatever) that maximizes the likelihood that you observed the things you observed.



## Curve Fitting: Bayesian Approach

---

We capture our assumptions about  $\mathbf{w}$ , before observing the data, in the form of a prior probability distribution  $p(\mathbf{w})$ .

The effect of the observed data  $D = \{t_1, \dots, t_N\}$  is expressed through the conditional probability  $p(D|\mathbf{w})$ .

Bayes' theorem, which takes the form

$$p(\mathbf{w}|D) = p(D|\mathbf{w}) p(\mathbf{w}) / p(D)$$

Thus, allows us to evaluate the uncertainty in  $\mathbf{w}$  after we have observed  $D$  in the form of the posterior probability  $p(\mathbf{w}|D)$ .

---

## Curve Fitting: Bayesian Approach

---

The quantity  $p(D|w)$  on the right-hand side of Bayes' theorem is evaluated for the observed data set  $D$

It can be viewed as a **function** of the parameter vector  $w$ , in which case it is called the **likelihood function**.

It expresses **how probable** the observed data set is **for different settings** of the parameter vector  $w$ .

**Note** that the likelihood is not a probability distribution over  $w$ , and its integral with respect to  $w$  does not (necessarily) equal one.

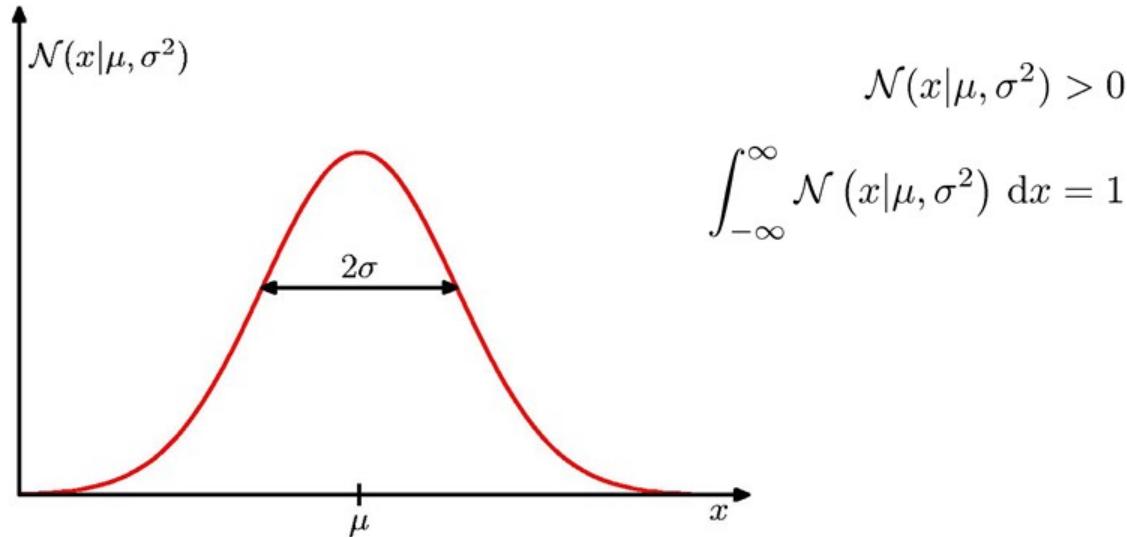
$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

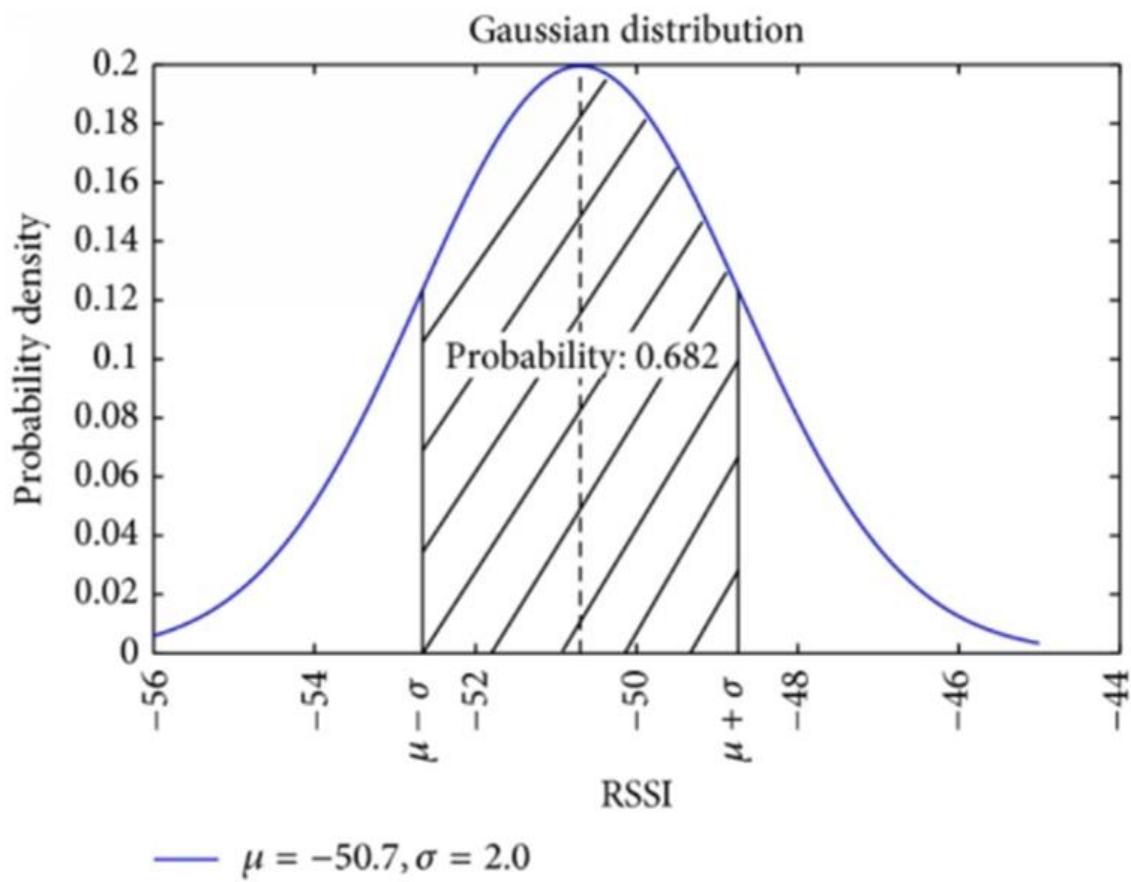
---

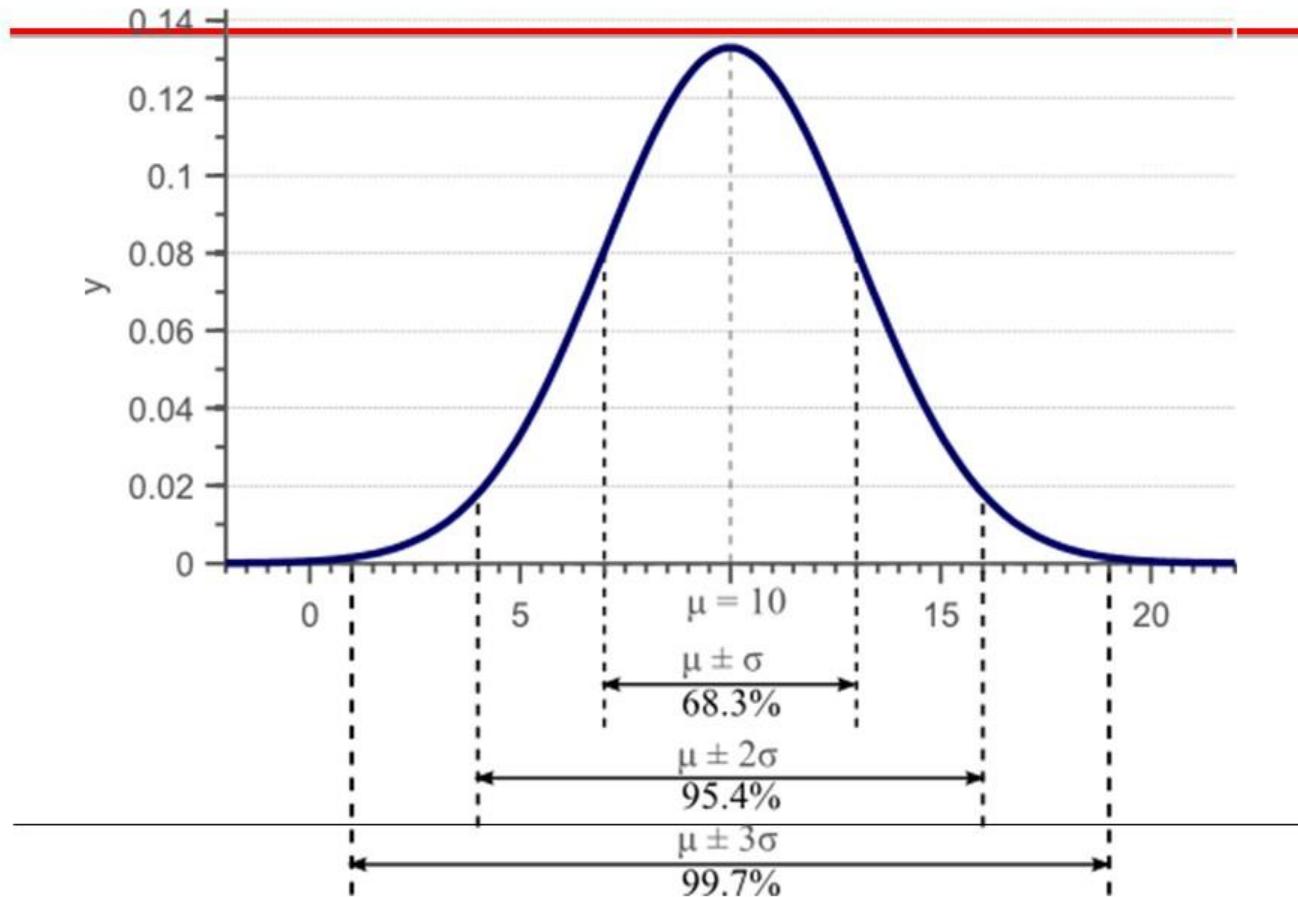
# The Gaussian Distribution

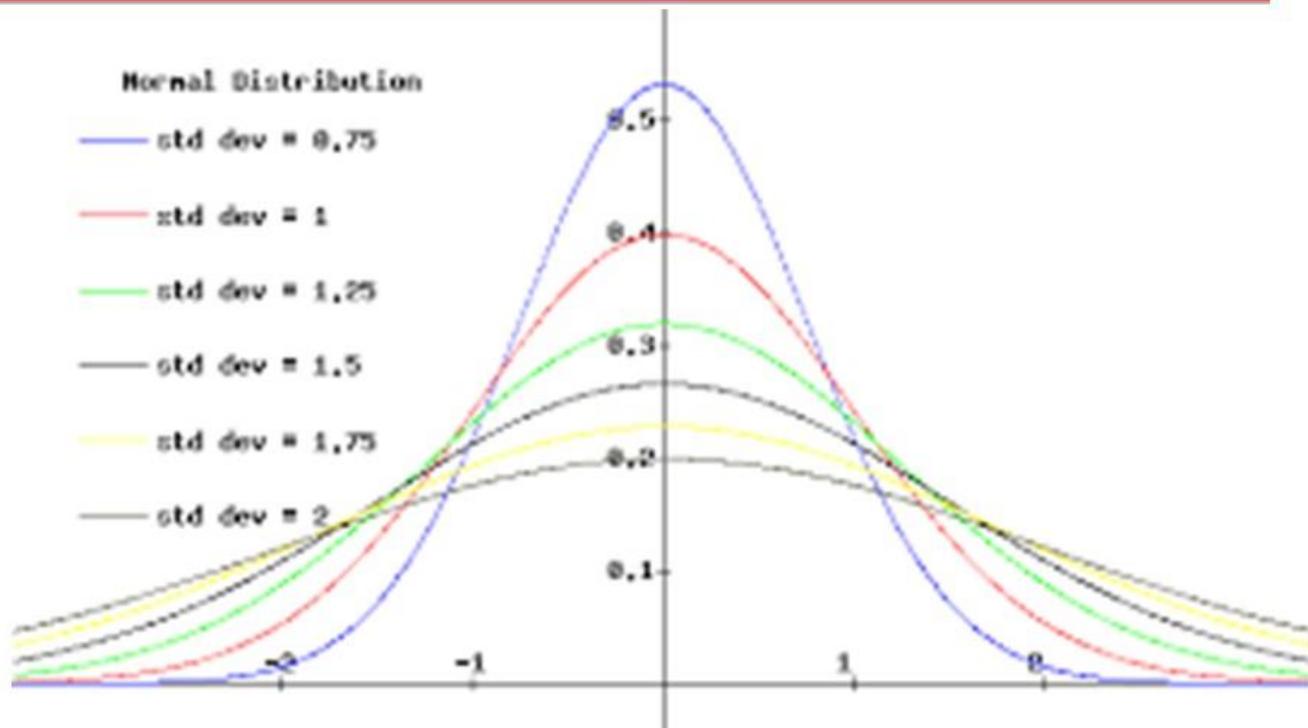
---

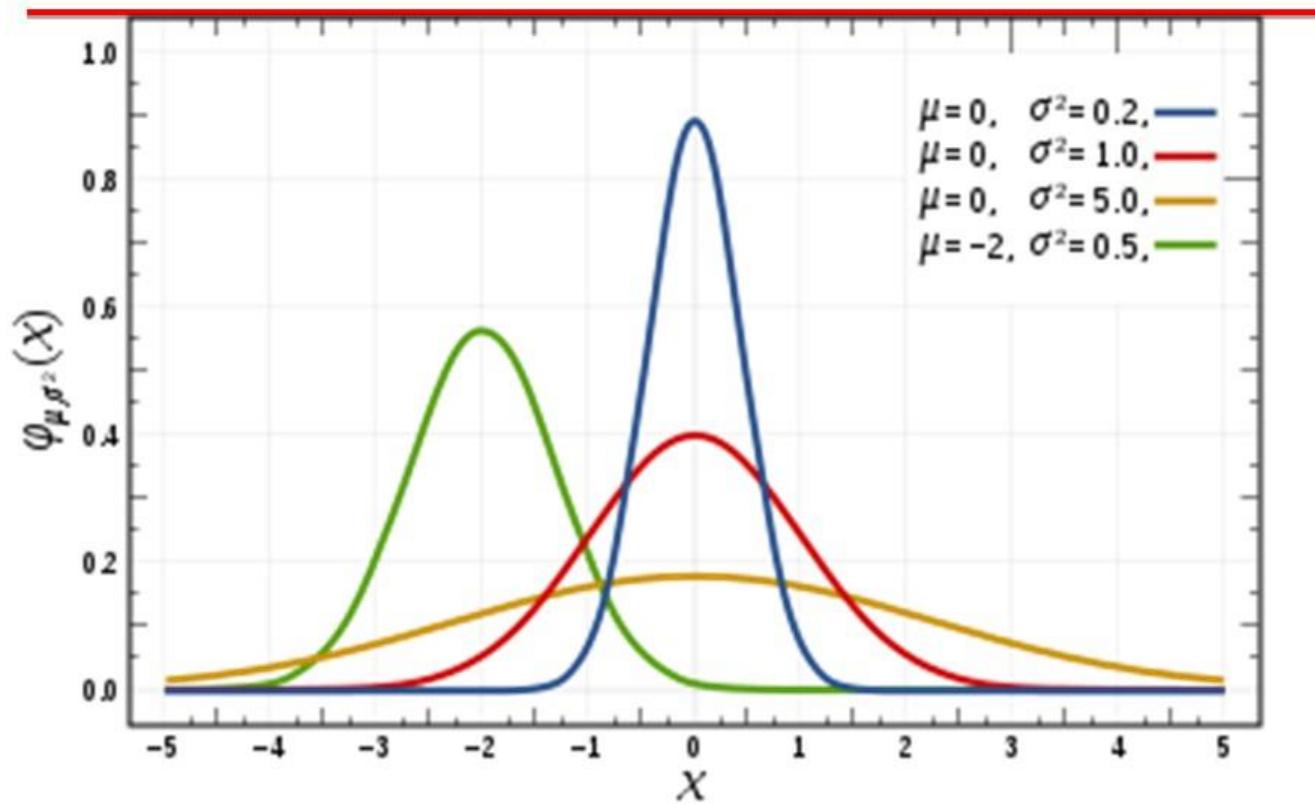
$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\}$$











# Gaussian Mean and Variance

---

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, dx = \mu$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 \, dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

# The Multivariate Gaussian

---

$$N(\mathbf{X}|\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu)\right\}$$

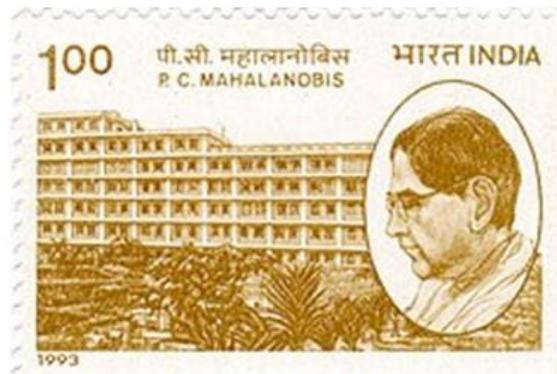
$\mathbf{X}$  is  $n$  dimensional vector

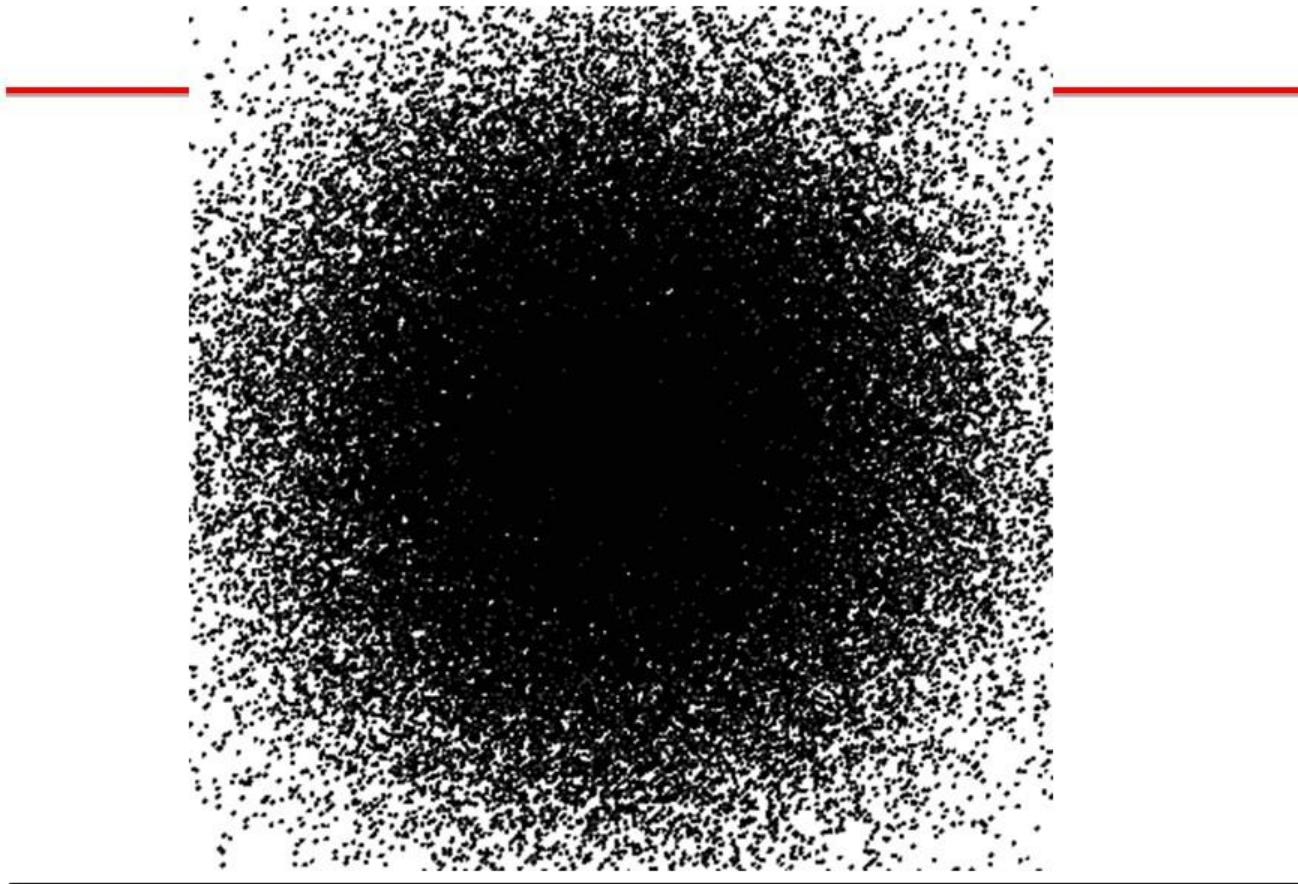
$\mu$  is mean vector of dimension  $n$

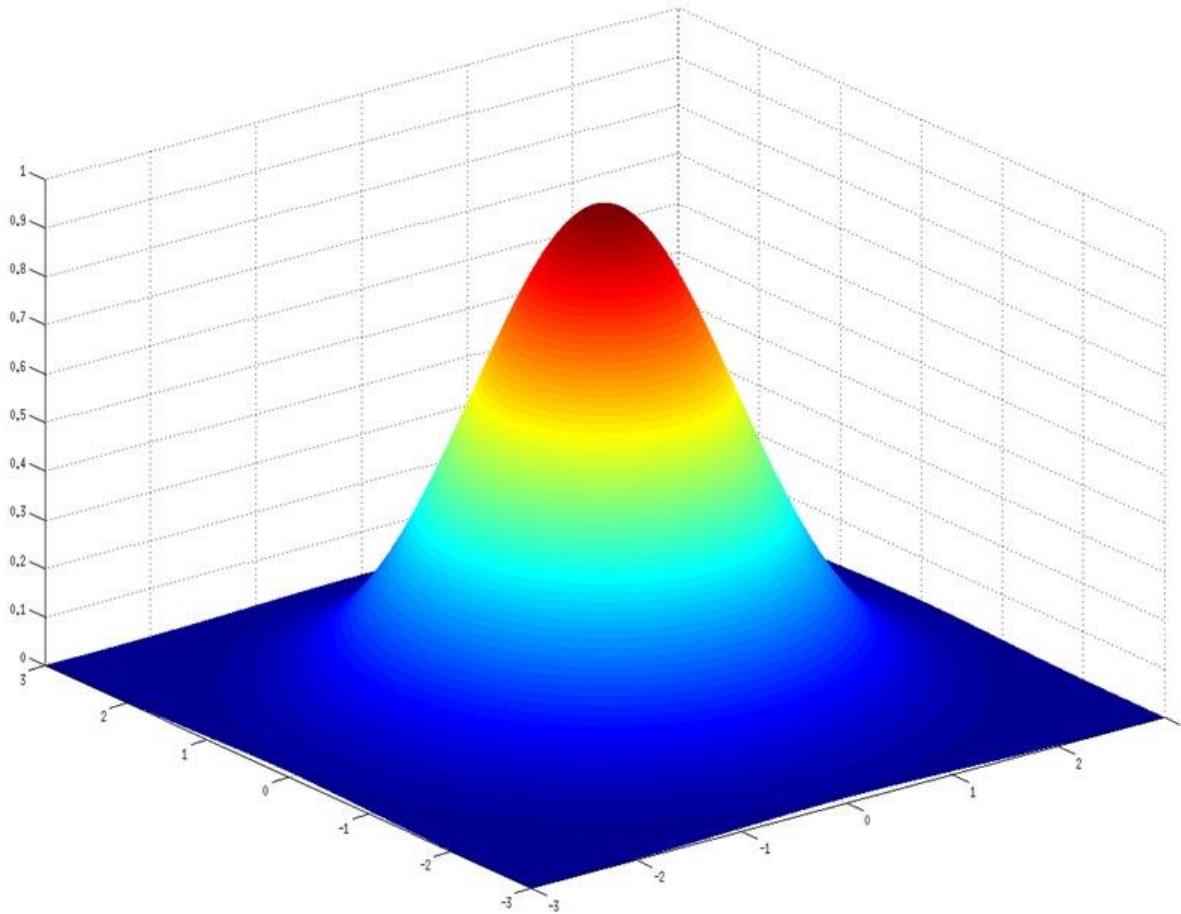
$\Sigma$  is  $n \times n$  covariance matrix

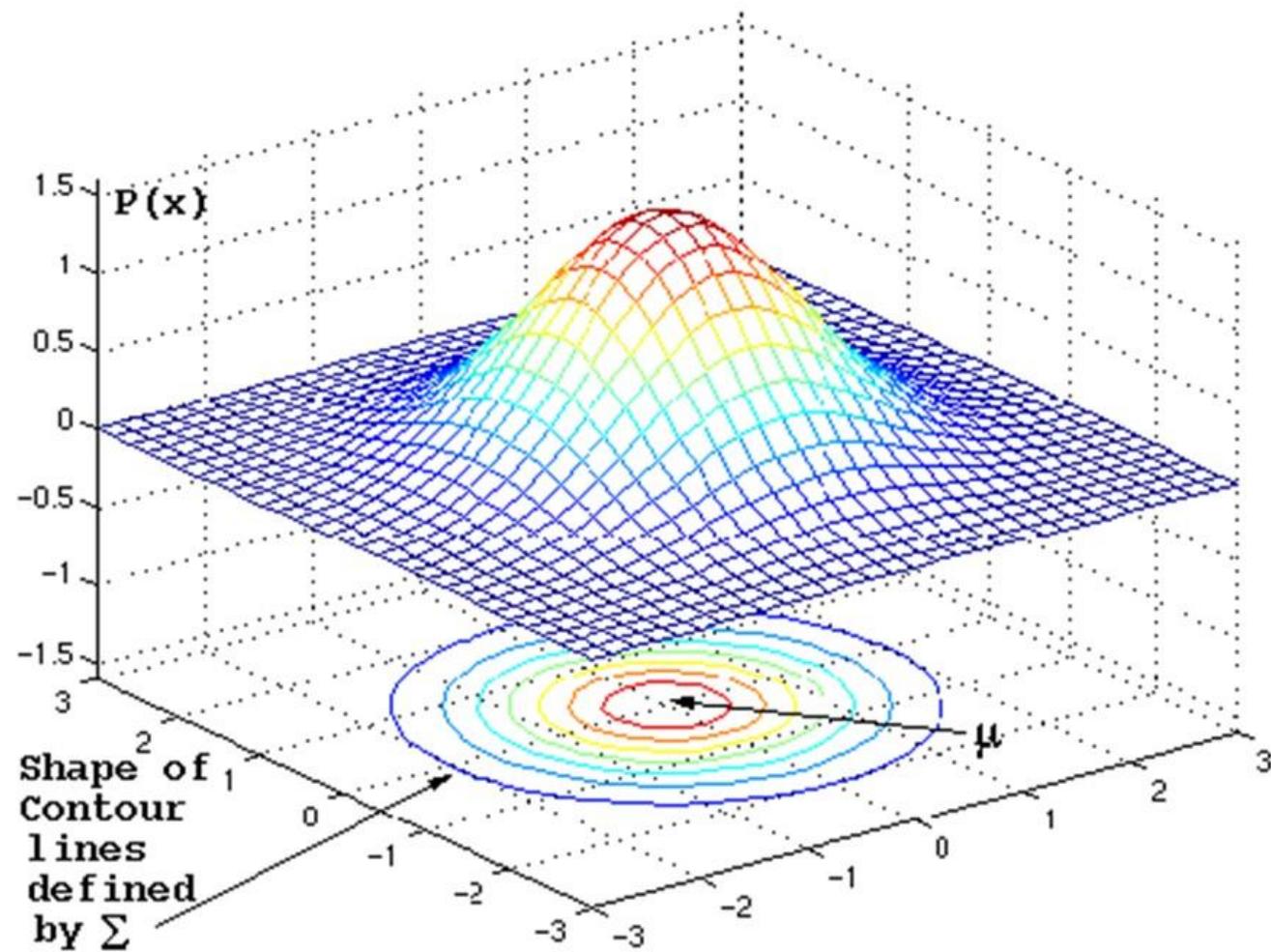
Note: 1.  $\Sigma$  is a symmetric matrix

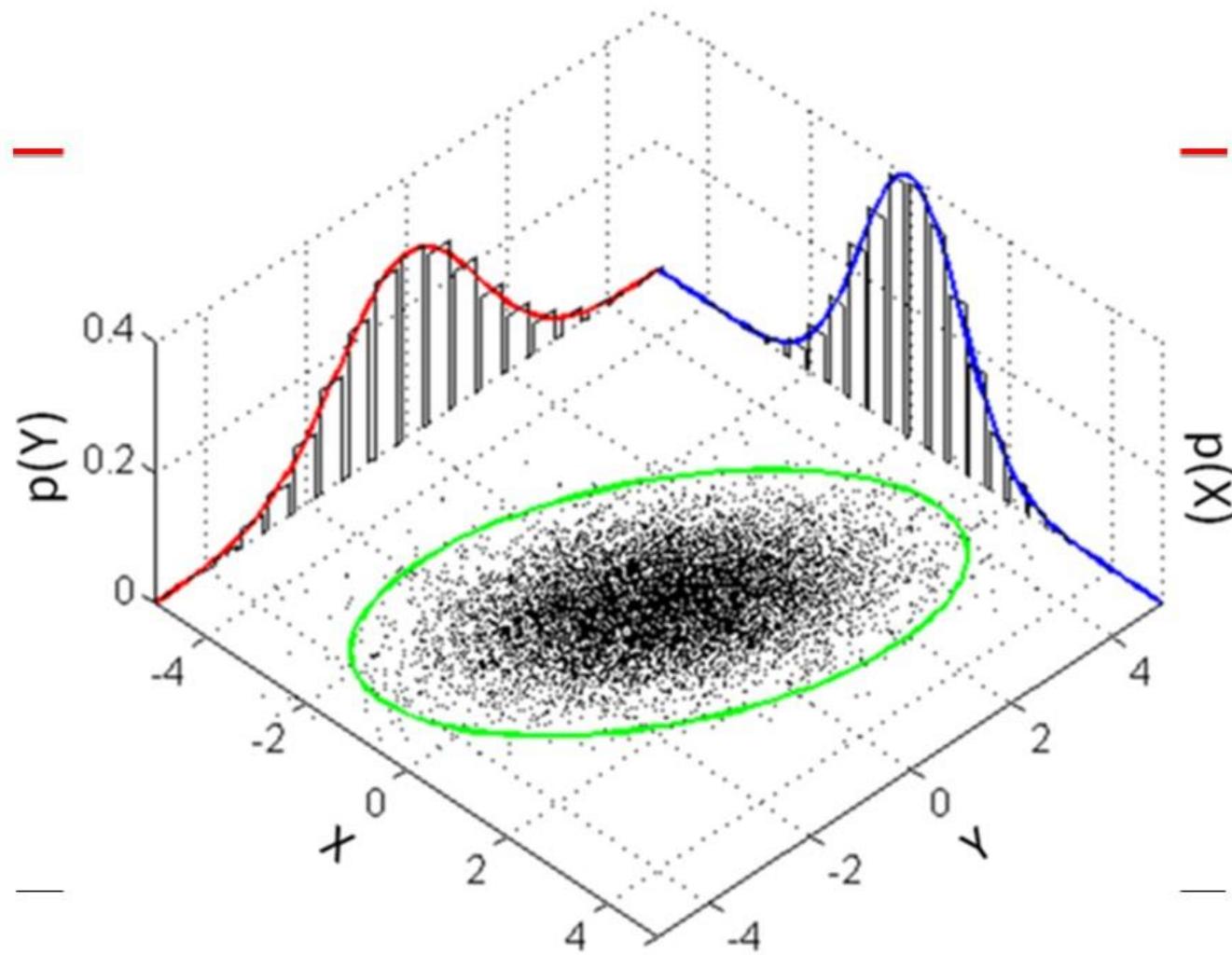
2.  $(\mathbf{X} - \mu)^T \Sigma^{-1} (\mathbf{X} - \mu)$  is a scalar quantity ( $1 \times n \times n \times n \times n \times 1 = 1 \times 1 = 1$ ).  
It is termed as Mahalanobis distance (squared).





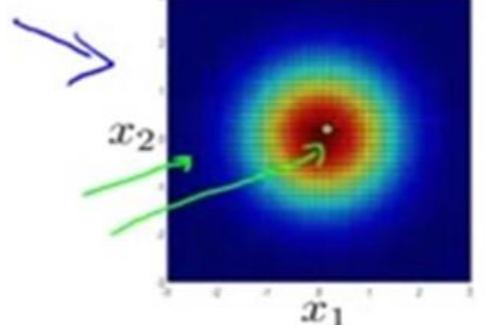
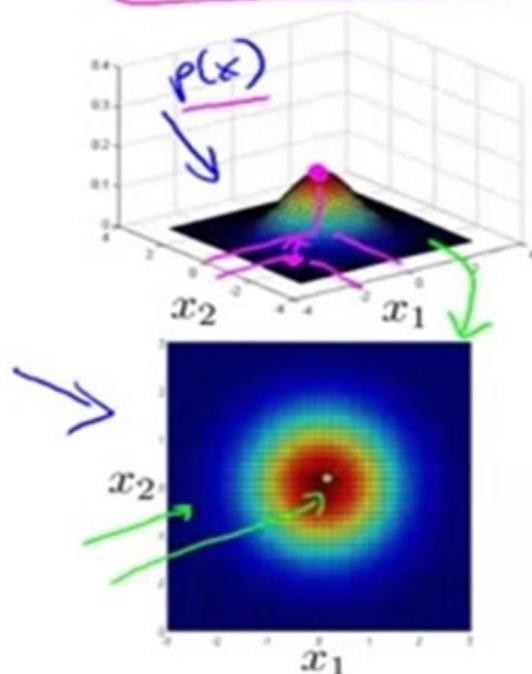






## Multivariate Gaussian (Normal) examples

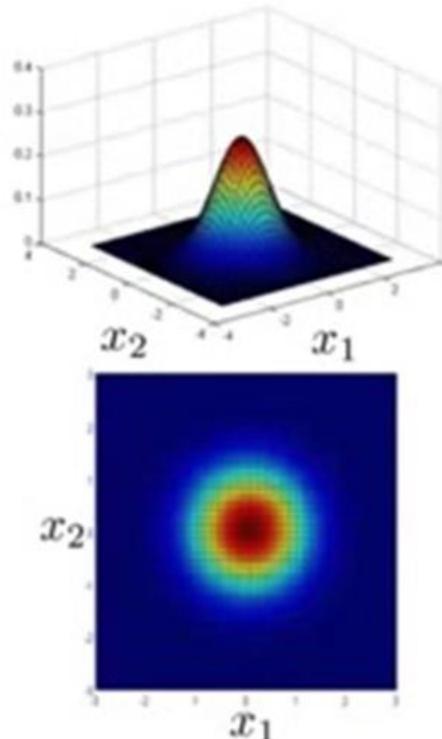
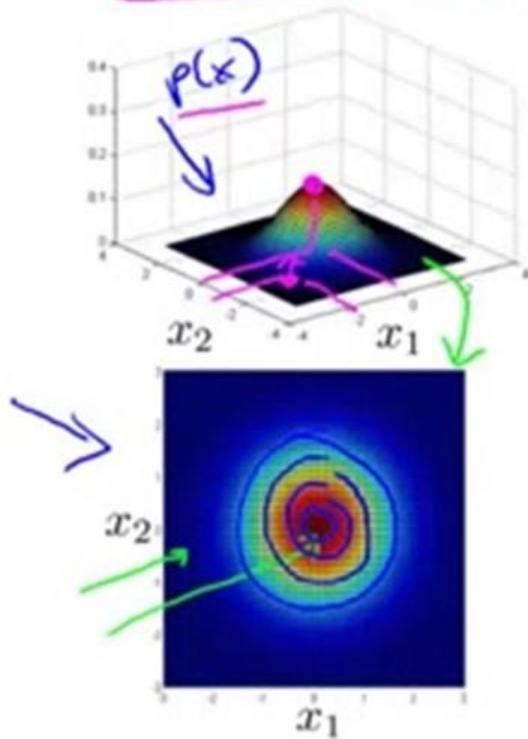
$$\rightarrow \boxed{\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}$$



## Multivariate Gaussian (Normal) examples

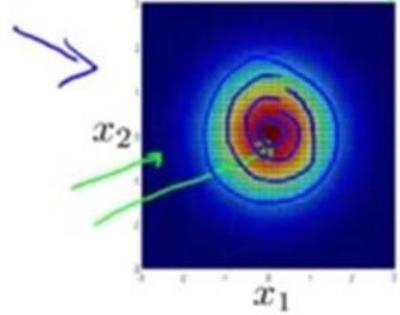
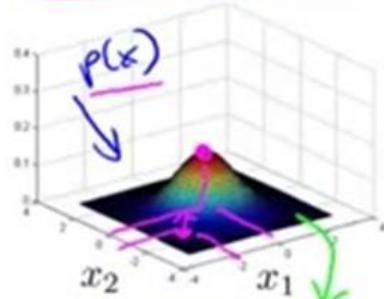
$$\rightarrow \boxed{\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}$$

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

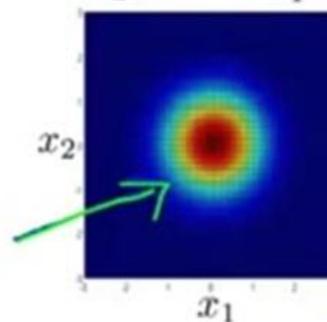
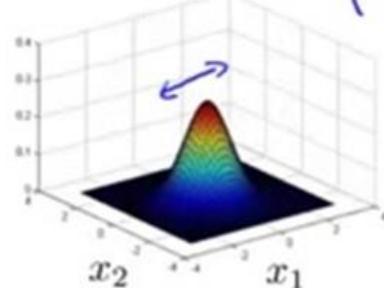


## Multivariate Gaussian (Normal) examples

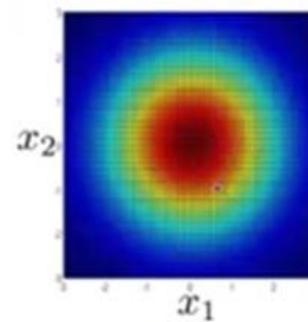
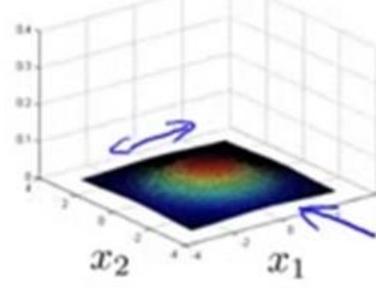
$$\rightarrow \mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



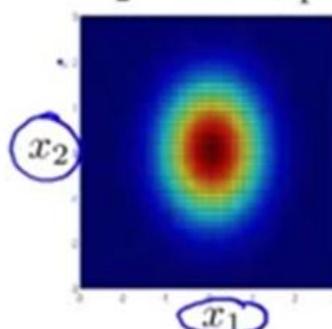
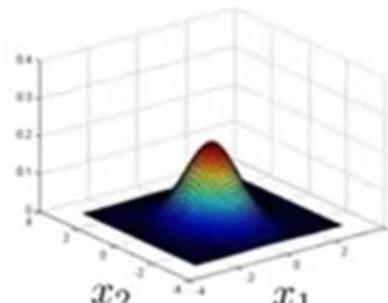
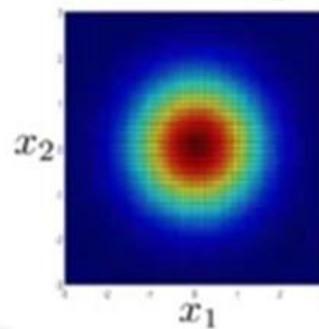
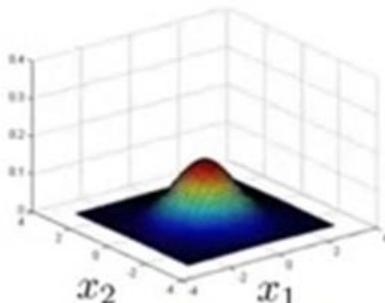
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



## Multivariate Gaussian (Normal) examples

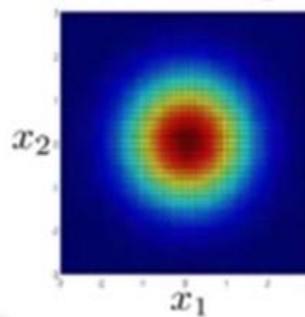
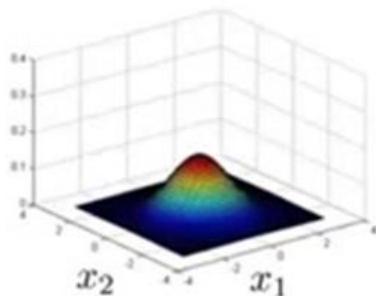
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

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

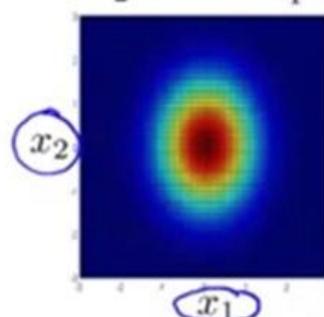
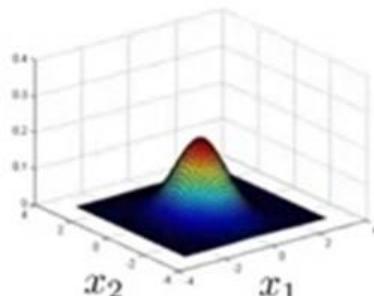


## Multivariate Gaussian (Normal) examples

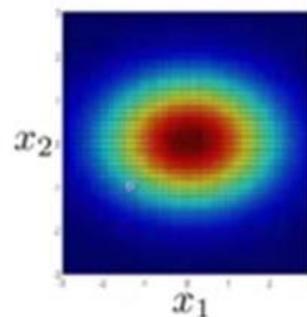
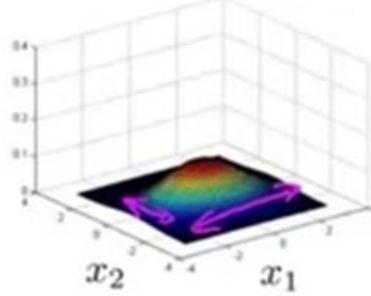
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



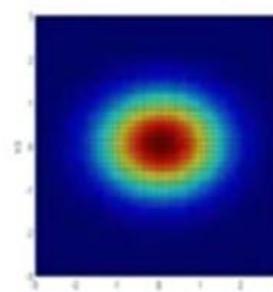
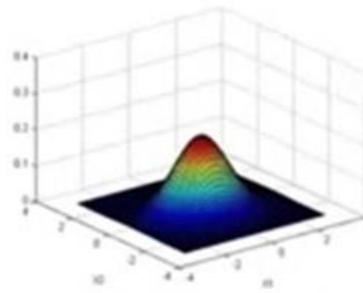
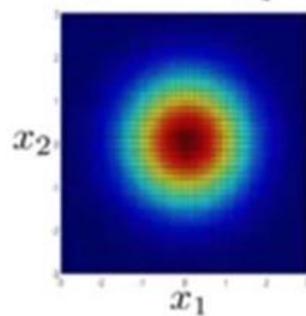
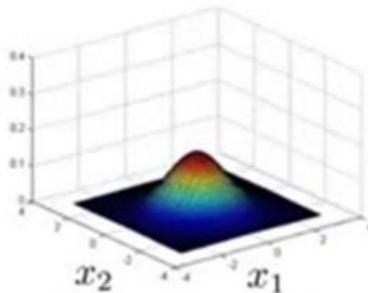
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix}$$



## Multivariate Gaussian (Normal) examples

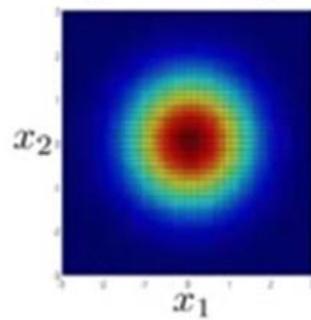
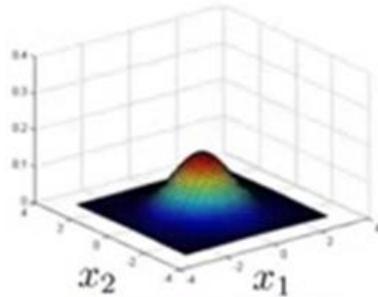
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

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

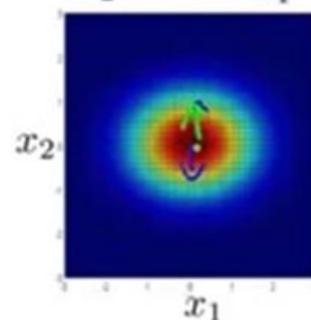
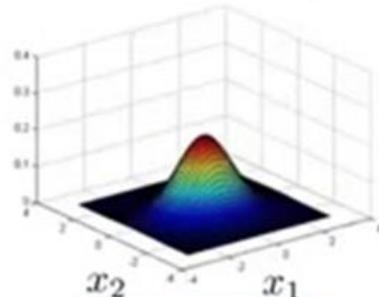


## Multivariate Gaussian (Normal) examples

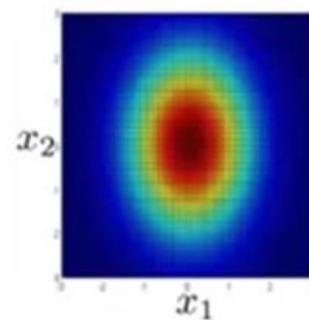
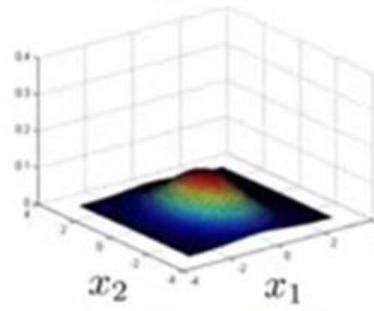
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



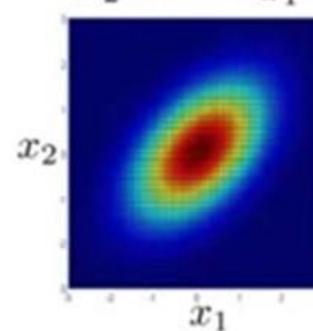
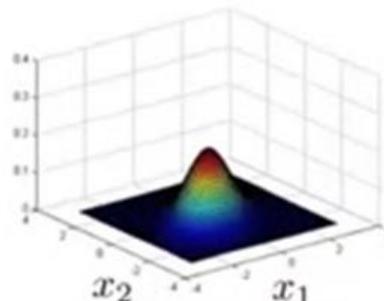
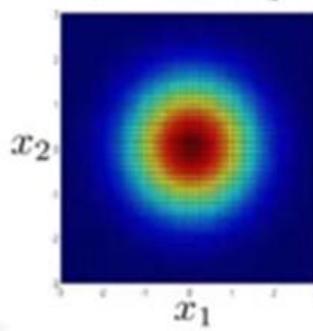
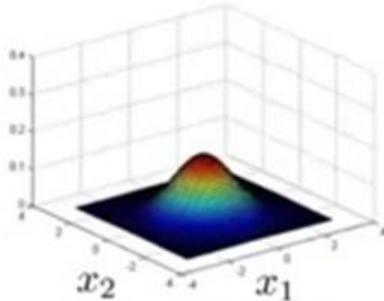
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$$



## Multivariate Gaussian (Normal) examples

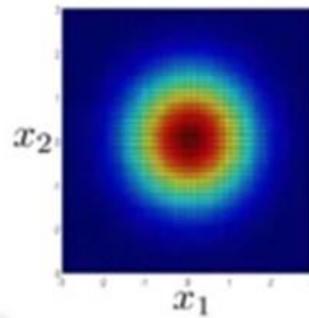
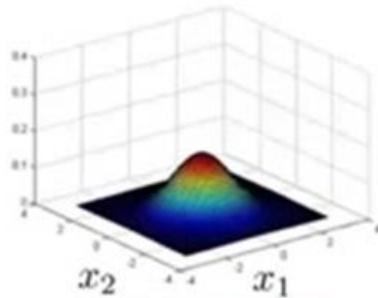
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

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

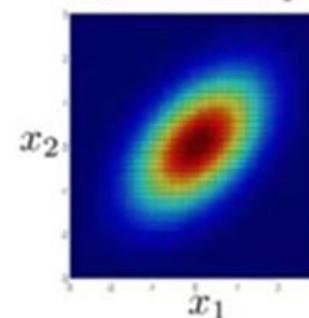
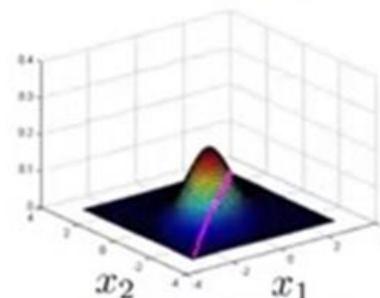


## Multivariate Gaussian (Normal) examples

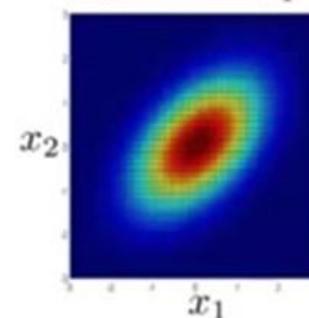
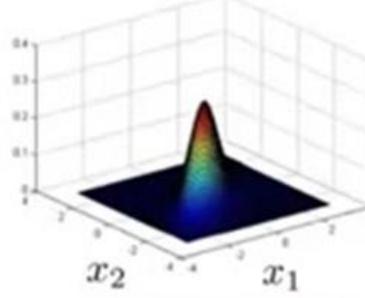
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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

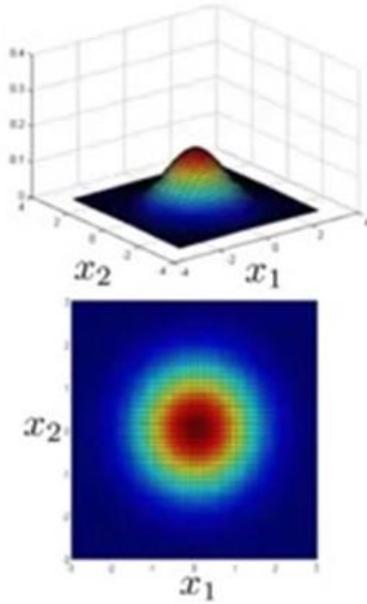


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

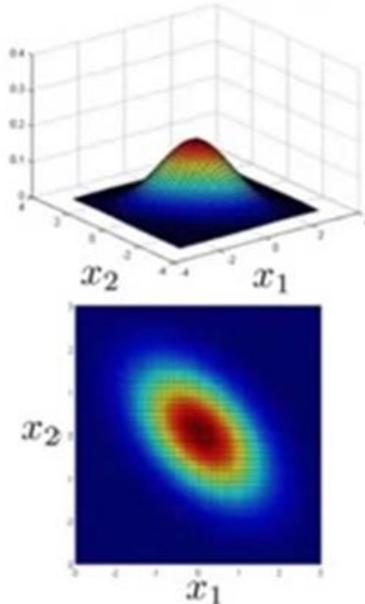


## Multivariate Gaussian (Normal) examples

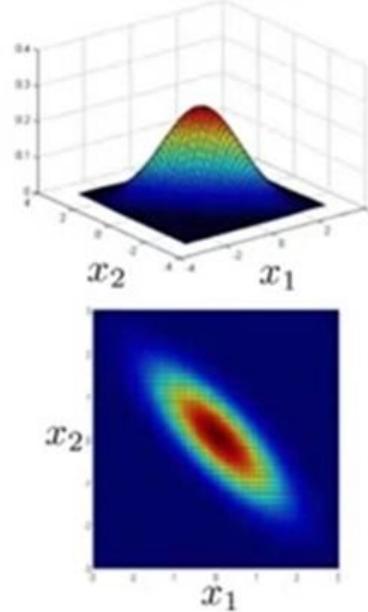
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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

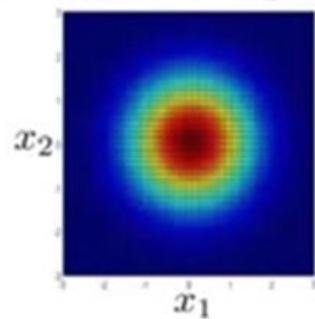
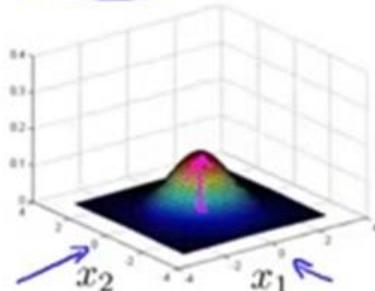


$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$

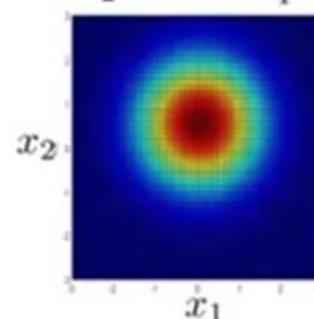
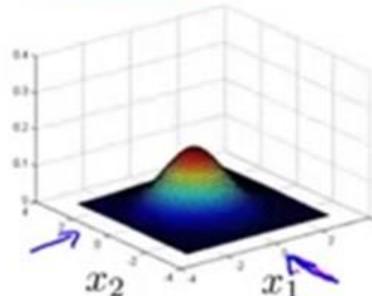


## Multivariate Gaussian (Normal) examples

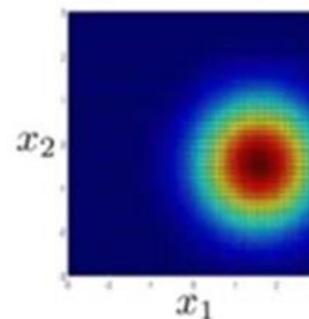
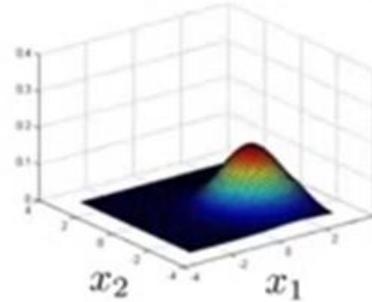
$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



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



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



# Geometry of the Gaussian

$$\Delta^2 = (x - \mu) \Sigma^{-1} (x - \mu)^T$$

Oval shows constant  $\Delta^2$  value...

Write  $\Sigma$  in terms of  
eigenvectors...

$$\Sigma = U \Lambda U^T$$

$$\Sigma = \begin{bmatrix} u_1 & u_2 \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \end{bmatrix} \begin{bmatrix} u_1 & \rightarrow \\ u_2 & \rightarrow \end{bmatrix}$$

Then...

$$y_i = u_i^T (x - \mu)$$

$$\Delta^2 = \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2}$$

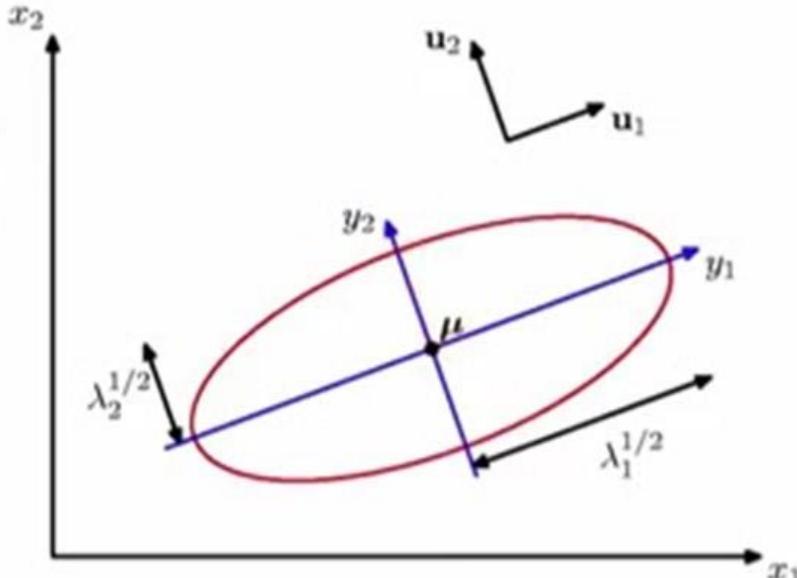
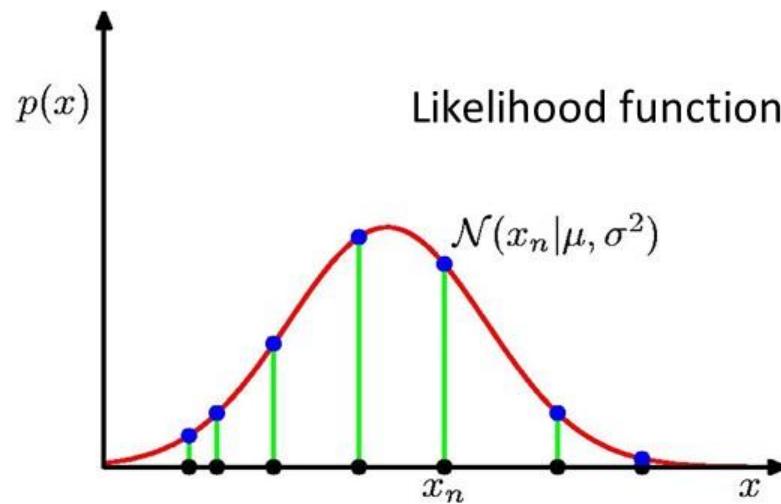


Figure from Chris Bishop's slides



# Gaussian Parameter Estimation

---



$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$

# Maximum (Log) Likelihood

---

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \quad \sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

---

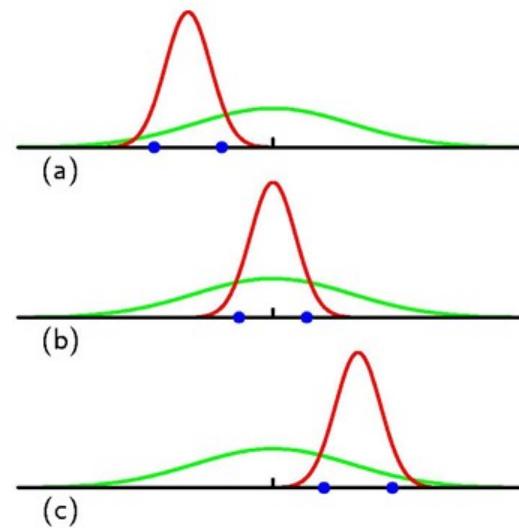
# Properties of $\mu_{\text{ML}}$ and $\sigma_{\text{ML}}^2$

---

$$\mathbb{E}[\mu_{\text{ML}}] = \mu$$

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \left(\frac{N-1}{N}\right) \sigma^2$$

$$\begin{aligned}\tilde{\sigma}^2 &= \frac{N}{N-1} \sigma_{\text{ML}}^2 \\ &= \frac{1}{N-1} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2\end{aligned}$$



# Curve Fitting Re-visited

---

The goal in the curve fitting problem is to be able to  
make predictions for the target variable  $t$   
given some new value of the input variable  $x$   
on the basis of a set of training data comprising  $N$  input values  
 $X = \{x_1, \dots, x_N\}$  and their corresponding target values  $t = \{t_1, \dots, t_N\}$ .

We can express our uncertainty over the value of the target variable using a probability distribution.

For this purpose, we shall assume that, given the value of  $x$ , the corresponding value of  $t$  has a Gaussian distribution with a mean equal to the value  $y(x, w)$  of the polynomial curve given by (1.1).

Thus we have: 
$$p(t|x, w, \beta) = N(t|y(x, w), \beta^{-1})$$

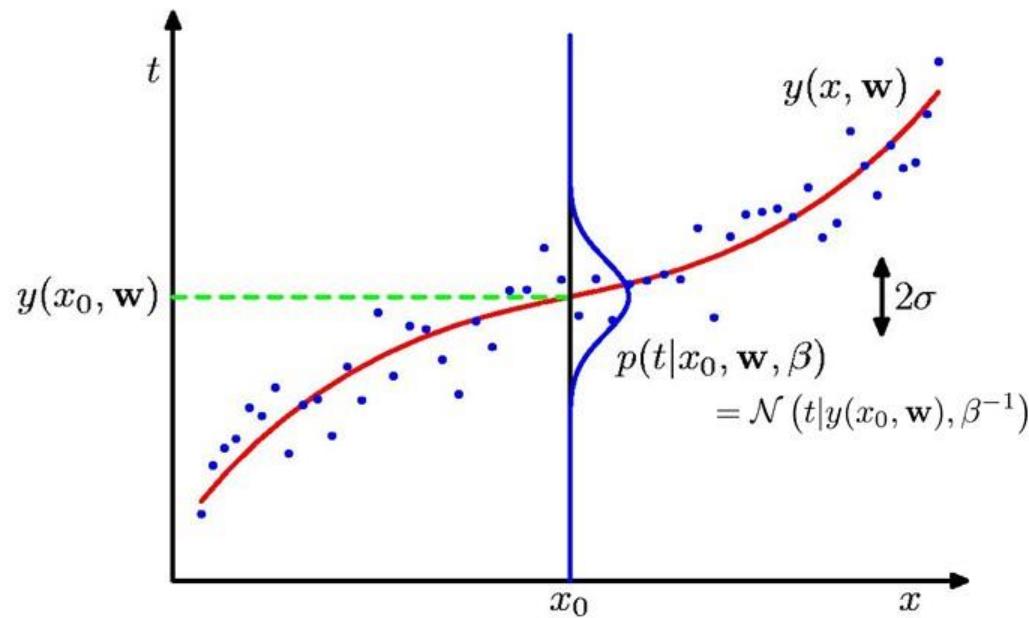
Where  $\beta = 1/\sigma^2$  is defined as precision parameter

---

Assume that data (target values) are drawn from above distribution independently.

## Curve Fitting Re-visited...

---



# Maximum Likelihood

---

We can use the training data  $\{\mathbf{x}, \mathbf{t}\}$  to determine the values of the unknown parameters  $\mathbf{w}$  and  $\beta$  by maximum likelihood. The likelihood function is given by

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1})$$
$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \underbrace{\sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

For obtaining  $\mathbf{w}_{ML}$ , maximize above equation with respect to  $\mathbf{w}$ .

we can omit the last two terms as they do not depend on  $\mathbf{w}$ .

Also, we note that scaling the log likelihood by a positive constant coefficient does not alter the location of the maximum with respect to  $\mathbf{w}$ , and so we can replace the coefficient  $\beta/2$  with  $1/2$ .

Finally, instead of maximizing the log likelihood, we can equivalently minimize the negative log likelihood.

---

# Maximum Likelihood

---

We therefore see that maximizing likelihood is equivalent, so far as determining  $\mathbf{w}$  is concerned, to minimizing the *sum-of-squares error function* defined by (1.2).

Thus the sum-of-squares error function has arisen as a consequence of maximizing likelihood under the assumption of a Gaussian noise distribution.

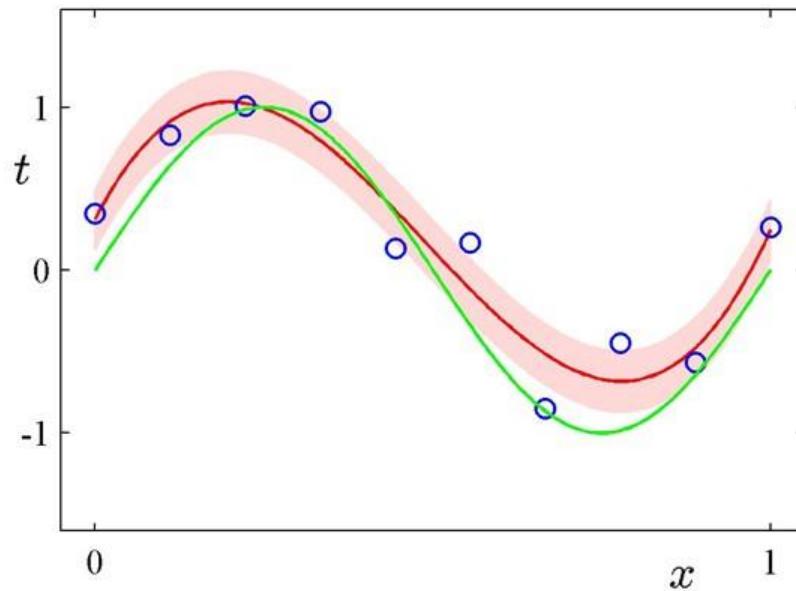
We can also use maximum likelihood to determine the precision parameter  $\theta$  of the Gaussian conditional distribution. Maximizing likelihood with respect to  $\theta$  gives

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$

# Predictive Distribution

---

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$



# MAP: A Step towards Bayes

---

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

Determine  $\mathbf{w}_{\text{MAP}}$  by minimizing regularized sum-of-squares error,  $\tilde{E}(\mathbf{w})$ .

---

# Bayesian Curve Fitting

---

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w} = \mathcal{N}(t|m(x), s^2(x))$$

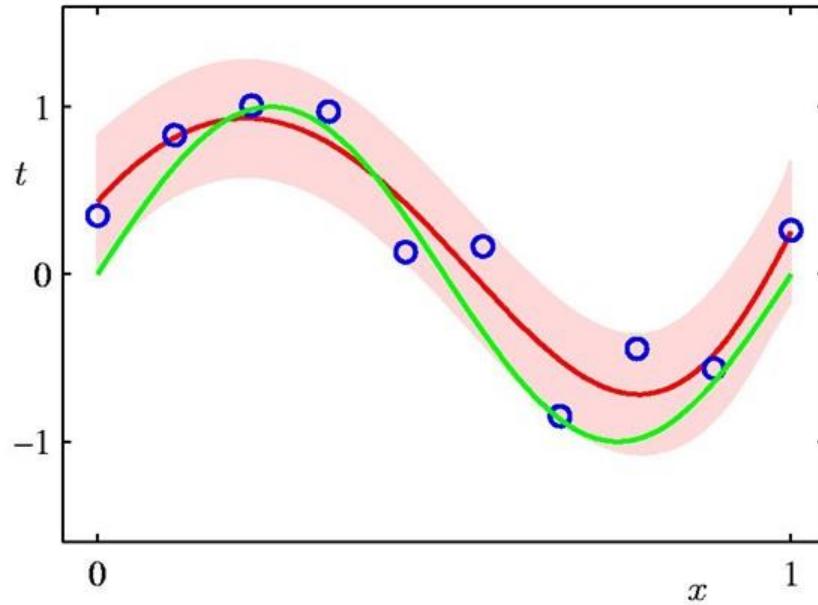
$$m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S} \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \quad s^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S} \boldsymbol{\phi}(x)$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \quad \boldsymbol{\phi}(x_n) = (x_n^0, \dots, x_n^M)^T$$

# Bayesian Predictive Distribution

---

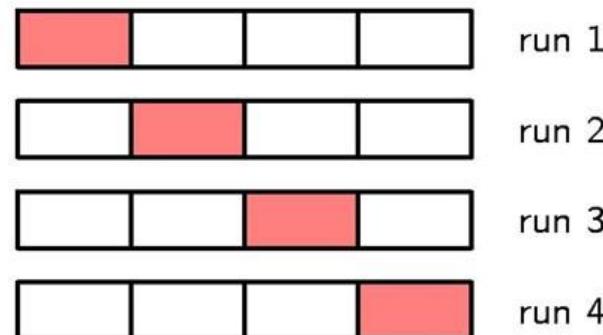
$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$



# Model Selection

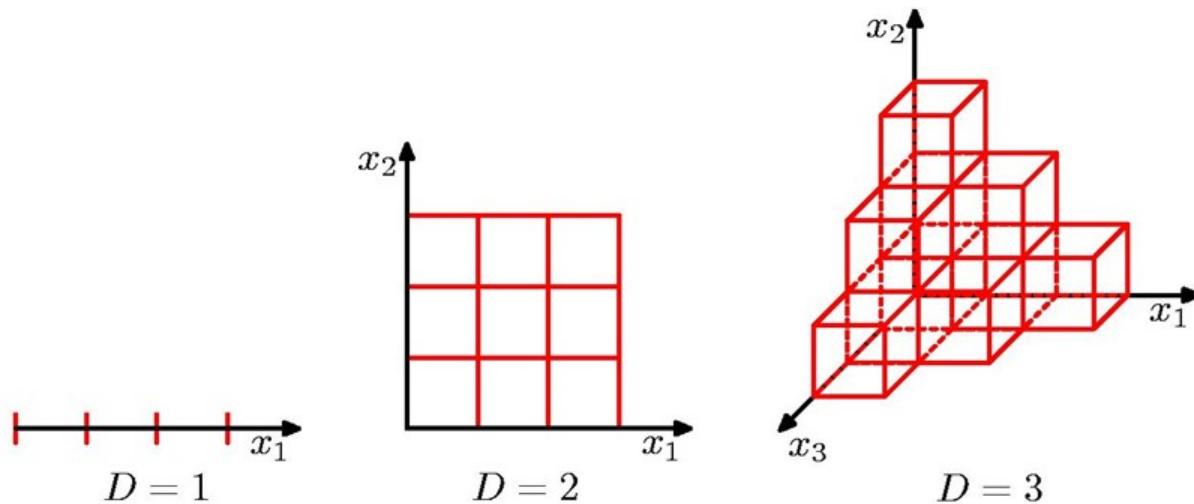
---

## Cross-Validation



# Curse of Dimensionality

---



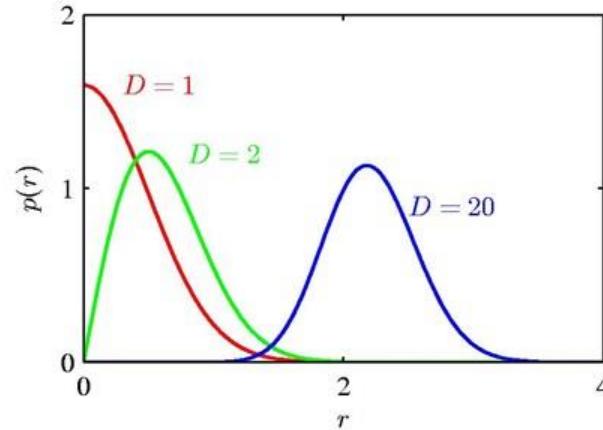
# Curse of Dimensionality

---

Polynomial curve fitting, M = 3

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

Gaussian Densities in  
higher dimensions



# Decision Theory

---

Inference step

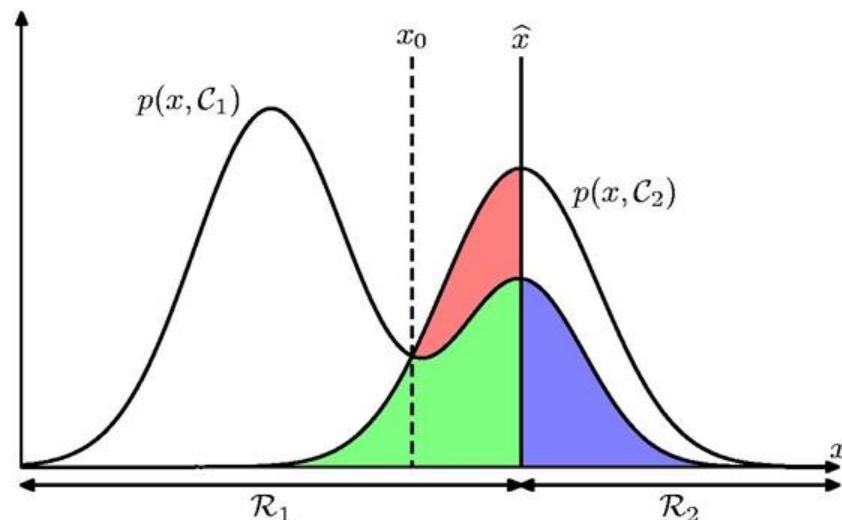
Determine either  $p(t|x)$  or  $p(x,t)$ .

Decision step

For given  $x$ , determine optimal  $t$ .

# Minimum Misclassification Rate

---



$$\begin{aligned} p(\text{mistake}) &= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1) \\ &= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}. \end{aligned}$$

---

# Minimum Expected Loss

---

Example: classify medical images as ‘cancer’ or ‘normal’

		Decision	
		cancer	normal
Truth	cancer	0	1000
	normal	1	0

# Minimum Expected Loss

---

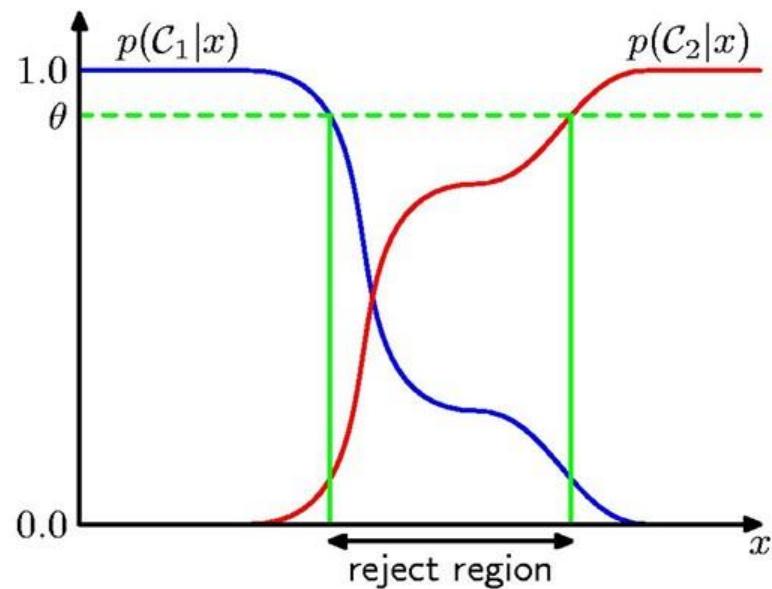
$$\mathbb{E}[L] = \sum_k \sum_j \int_{\mathcal{R}_j} L_{kj} p(\mathbf{x}, \mathcal{C}_k) d\mathbf{x}$$

Regions  $\mathcal{R}_j$  are chosen to minimize

$$\mathbb{E}[L] = \sum_k L_{kj} p(\mathcal{C}_k | \mathbf{x})$$

## Reject Option

---



# Why Separate Inference and Decision?

---

- Minimizing risk (loss matrix may change over time)
  - Reject option
  - Unbalanced class priors
  - Combining models
-

# Decision Theory for Regression

---

Inference step

Determine  $p(\mathbf{x}, t)$ .

Decision step

For given  $\mathbf{x}$ , make optimal prediction,  $y(\mathbf{x})$ , for  $t$ .

Loss function:  $\mathbb{E}[L] = \iint L(t, y(\mathbf{x}))p(\mathbf{x}, t) d\mathbf{x} dt$

---

# The Squared Loss Function

---

$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$$\begin{aligned}\{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2\end{aligned}$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

$$y(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}]$$

# Generative vs Discriminative

---

Generative approach:

Model  $p(t, \mathbf{x}) = p(\mathbf{x}|t)p(t)$

Use Bayes' theorem  $p(t|\mathbf{x}) = \frac{p(\mathbf{x}|t)p(t)}{p(\mathbf{x})}$

Discriminative approach:

Model  $p(t|\mathbf{x})$  directly

# Entropy

---

$$H[x] = - \sum_x p(x) \log_2 p(x)$$

Important quantity in

- coding theory
  - statistical physics
  - machine learning
-

# Entropy

---

Coding theory:  $x$  discrete with 8 possible states; how many bits to transmit the state of  $x$ ?

All states equally likely

$$H[x] = -8 \times \frac{1}{8} \log_2 \frac{1}{8} = 3 \text{ bits.}$$

# Entropy

---

$x$	a	b	c	d	e	f	g	h
$p(x)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$	$\frac{1}{64}$
code	0	10	110	1110	111100	111101	111110	111111

$$\begin{aligned} H[x] &= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{8} \log_2 \frac{1}{8} - \frac{1}{16} \log_2 \frac{1}{16} - \frac{4}{64} \log_2 \frac{1}{64} \\ &= 2 \text{ bits} \end{aligned}$$

$$\begin{aligned} \text{average code length} &= \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + 4 \times \frac{1}{64} \times 6 \\ &= 2 \text{ bits} \end{aligned}$$

---

# Entropy

---

In how many ways can N identical objects be allocated M bins?

$$W = \frac{N!}{\prod_i n_i!}$$

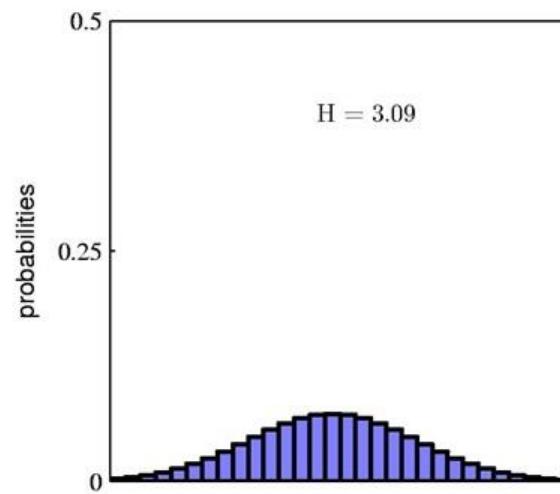
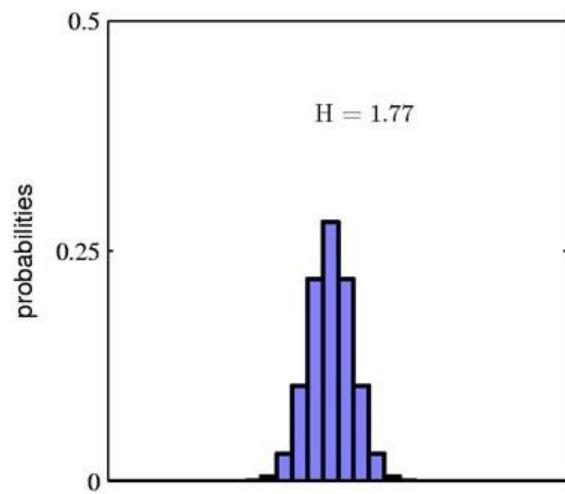
$$H = \frac{1}{N} \ln W \simeq - \lim_{N \rightarrow \infty} \sum_i \left( \frac{n_i}{N} \right) \ln \left( \frac{n_i}{N} \right) = - \sum_i p_i \ln p_i$$

Entropy maximized when  $\forall i : p_i = \frac{1}{M}$

---

# Entropy

---



# Differential Entropy

---

Put bins of width  $\epsilon$  along the real line

$$\lim_{\Delta \rightarrow 0} \left\{ - \sum_i p(x_i) \Delta \ln p(x_i) \right\} = - \int p(x) \ln p(x) dx$$

Differential entropy maximized (for fixed  $\sigma^2$ ) when

$$p(x) = \mathcal{N}(x|\mu, \sigma^2)$$

in which case

$$H[x] = \frac{1}{2} \{1 + \ln(2\pi\sigma^2)\}.$$

# Conditional Entropy

---

$$H[y|x] = - \iint p(y, x) \ln p(y|x) dy dx$$

$$H[x, y] = H[y|x] + H[x]$$

---

# The Kullback-Leibler Divergence

---

$$\begin{aligned}\text{KL}(p\|q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}) d\mathbf{x} - \left( - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \right) \\ &= - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x}\end{aligned}$$

$$\text{KL}(p\|q) \simeq \frac{1}{N} \sum_{n=1}^N \{-\ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \ln p(\mathbf{x}_n)\}$$

$$\text{KL}(p\|q) \geq 0 \quad \text{KL}(p\|q) \neq \text{KL}(q\|p)$$

---

# Mutual Information

---

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) \| p(\mathbf{x})p(\mathbf{y})) \\ &= - \iint p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x} d\mathbf{y} \end{aligned}$$

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

