

# CSCI 4022 Spring 2021

## Clustering and Likelihoods

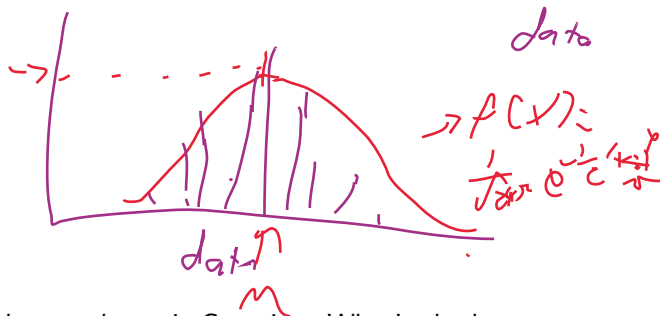
center of data:

$$\bar{X} \approx \mu$$

**Example:** Suppose we have a set of data that we *know* is Gaussian. What's the best way to estimate the true population mean  $\mu$  and  $\sigma^2$ ?

$$\frac{\sum (x_i - \bar{x})^2}{n-1} = s^2 \approx \sigma^2$$

data      pop



# Announcements and To-Dos

## Announcements:

1. HW 2 due Wednesday; HW 3 due next Monday.

↳ Extra OH on Tuesday / tomorrow  
4p-5p

# Clustering Recap

Two methods so far

## Hierarchical

- No initialization
- At each step, we compute *all* distances between pairs of clusters. Then merge the nearest two clusters.
- Once a cluster is formed, it is represented only by its **centroid** (average) or **clustroid** (median/representative point)
- Stop at  $k$  clusters.

## kmeans

- Random initialization of  $k$  clusters
- Loop over points; compute *all* point-to-cluster distances. Then assign each point to its nearest two.
- Loop over clusters; represent each only by its **centroid** (average)
- Stop at “things aren’t changing”

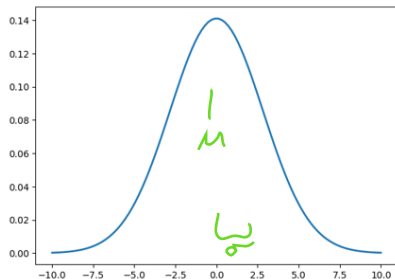
Each method required a notion of distance where it was easy to compare dimension-to-dimension. We can relax this!

# The Gaussian (normal) distribution

You should recall the *normal distribution*, *Gaussian distribution*, or *bell curve*. In one dimension, it's a beautiful little function. (squared negative exponential).

**The normal distribution** has two arguments or *parameters*:

- $\mu$ : The mean or center of the curve. The most important *location*. A single scalar.
- $\sigma^2$ : The variance/width of the curve. The most common measure of *dispersion* or spread. A single positive scalar.



A normal

# The Gaussian (normal) distribution

## Formalism

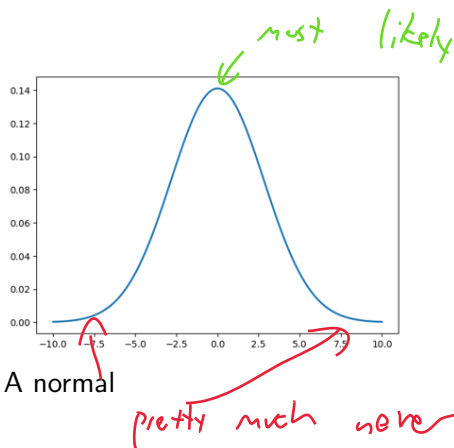
pdf: The *probability density function* of the normal distribution is

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

cdf: The **integral** of the pdf lets us answer probability questions like

$$P(a < X < b) = \int_a^b f(x) dx$$

- The normal is known for assigning *low probabilities* to outliers, or points *far from*  $\mu$



# The 2-D Gaussian

Variance:

$$E[(X - \bar{X})^2]$$

$$E[X]$$

In multiple dimensions, everything gets a bit more complex.

parameters:

$\mu$ : The mean or center of the surface. This is a 2-D  $(x, y)$  tuple, so we may write

$$\mu = (\mu_1, \mu_2) \in \mathbb{R}^2$$

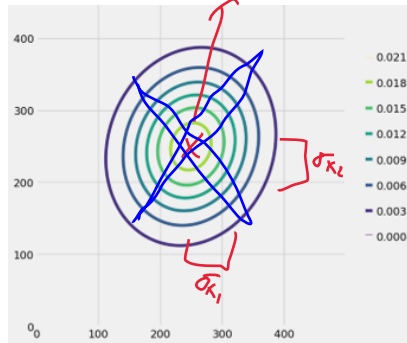
$\Sigma$ : A  $2 \times 2$  covariance matrix whose entries are

$$\begin{bmatrix} \sigma_1^2 & \text{cov}(x_1, x_2) \\ \text{cov}(x_2, x_1) & \sigma_2^2 \end{bmatrix}$$

where  $\sigma_1^2$  is the variance in the **first axis direction**,  $\sigma_2^2$  is the variance in the **second axis direction**, and

$$\text{cov}(x_1, x_2) = E[(X_1 - E[X_1])(X_2 - E[X_2])]$$

Mullen: Likelihood



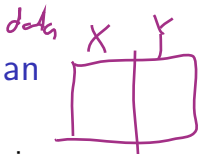
A 2-D normal's contour plot

$$X_1 > E[X_1]$$

$$X_2 > E[X_2]$$



# The 2-D Gaussian



This is actually **5** unique parameters.  
*parameters:*

$\mu_1$ : The location of the center in the axis-1 dimension.

$\mu_2$ : The location of the center in the axis-2 dimension.

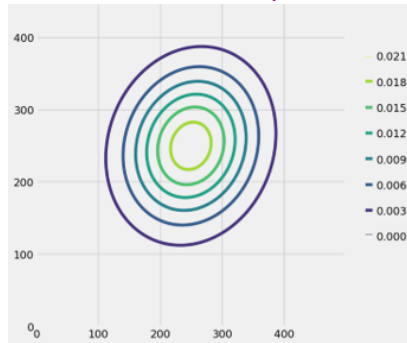
$\sigma_1^2$ : The variance in the axis-1 direction.

$\sigma_2^2$ : The variance in the axis-2 direction.

*cov*:  $cov(x_1, x_2)$  The covariance of the data set (dimension one covariance with dimension 2).

Usage: NP.COV(X1, Y1)

$$\frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x \cdot s_y}$$



A 2-D normal's contour plot

## The 2-D Gaussian: variance and covariance

1. The off-diagonal arguments of  $\Sigma$  lead to *rotations*.

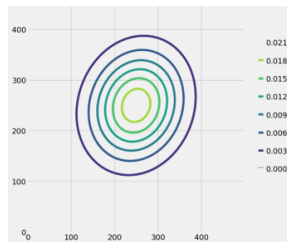
**Results:** if  $cov(x_1, x_2) = 0$ , then the Gaussian is an ellipse that's oriented **vertically/horizontally** (i.e. the semi-major and semi-minor axis are parallel to the coordinate axes.).

2.  $\sigma_1^2$  and  $\sigma_2^2$  determine the widths in their respective coordinate directions

**Results:** if  $cov(x_1, x_2) = 0$  **and**  $\sigma_1 = \sigma_2$ , the contours would be circles... like k-means.

Allowing  $\sigma_1$  and  $\sigma_2$  to vary stretches the circle in the corresponding direction. This gives a new method of clustering: **Gaussian Mixture Models**

$$\begin{bmatrix} \sigma_1^2 & cov(x_1, x_2) \\ cov(x_2, x_1) & \sigma_2^2 \end{bmatrix}$$





## The Gaussian Mixture Model (GMM): 1D Example

Motivating example/cautionary tale: Suppose you go to Chuck E. Cheese's for your niece's birthday party. As you look around, you start to feel rather self-conscious because there seem to be very few people around your age. Needing to pass the time, because you feel so, so awkward, you collect data on everyone's ages.



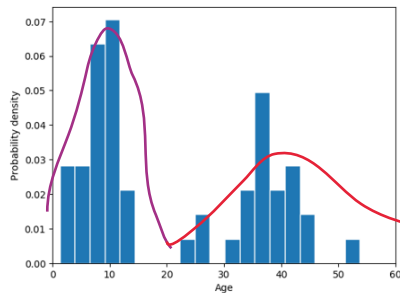
## The Gaussian Mixture Model (GMM): 1D Example

Motivating example/cautionary tale: Suppose you go to Chuck E. Cheese's for your niece's birthday party. As you look around, you start to feel rather self-conscious because there seem to be very few people around your age. Needing to pass the time, because you feel so, so awkward, you collect data on everyone's ages.

You are then kicked out of and permanently banned from Chuck E. Cheese's because you were accosting children and asking about their ages...

... Fair enough.

Now that your afternoon is freed up, you plot up a histogram of the age data you so creepily and painstakingly collected. It looks like this:

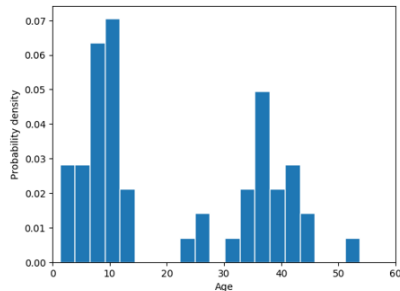


# The Gaussian Mixture Model (GMM): 1D Example

**The task:** How can we best model this distribution?

It's clearly **bimodal**, so a simple normal Gaussian is not appropriate. Instead, we view it as possibly the combination of **two** distributions:

1. The **kids'** ages seem like they might be reasonably Gaussian, aside from the pesky fact of non-negativity on ages.
2. The **parents'** ages might also be modeled by a *different* Gaussian distribution
3. This would make the overall distribution of patrons' ages a *Gaussian mixture model*



# The Gaussian Mixture Model (GMM): 1D Example

**The whole model:**

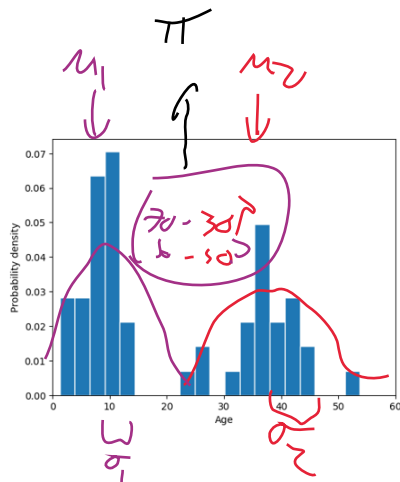
**Kids' ages:** a normal  $X_1 \sim N(\mu_1, \sigma_1^2)$

**Parents' ages:** a normal  $X_2 \sim N(\mu_2, \sigma_2^2)$

**Who's who:** a patron could be either of these at some *mixing proportion*. Define the *Bernoulli* random variable  $\Delta$  where the probability that a patron is an adult ( $\Delta = 1$ ) by  $\pi$ . Then:

**All patrons' ages:**

$$\hat{X} = (1 - \Delta) \cdot X_1 + \Delta \cdot X_2$$



# The Gaussian Mixture Model (GMM): 1D Example

$$\Delta = 1 \text{ w/ Prob } \pi$$

**The whole model:**

$$X = (1 - \Delta) \cdot X_1 + \Delta \cdot X_2$$

Let's unpack, since this is secretly just adding up all the possible ways we can observe a specific age:

$$\underbrace{X}_{\text{Prob of specific age}} = \underbrace{(1 - \Delta)}_{\substack{\text{Prob person is adult} \\ \text{if } \Delta=0 \\ \text{ok}}} \cdot \underbrace{X_1}_{\text{Prob an "adult" is that age}} + \underbrace{\Delta}_{\substack{\text{Prob person is child} \\ \text{if } \Delta=1 \\ \text{ok}}} \cdot \underbrace{X_2}_{\text{Prob a "child" is that age}}$$

Handwritten notes:  $\Delta=0$  (purple),  $\Delta=1$  (red), and arrows pointing to the corresponding terms in the equation.

## The Gaussian Mixture Model (GMM): 1D Example

**The whole model:**  $X = (1 - \Delta) \cdot X_1 + \Delta \cdot X_2$

**Definition:** The GMM is a *generative* model, since it specifies the probabilities for new data points.

1. Sample or simulate a  $\Delta$  with a coin flip or `NP.RANDOM.CHOICE`
2. Based on  $\Delta$ , sample a random normal from:
  - 2.1  $X_1$  as a  $N(\mu_1, \sigma_1^2)$  if  $\Delta = 0$  **OR**
  - 2.2  $X_2$  as a  $N(\mu_2, \sigma_2^2)$  if  $\Delta = 1$

Our task is sometimes to *generate*, but first we have to *estimate* the underlying parameters used in the model. To use the model, we have **5** things to estimate or choose.

$$\Theta = (\pi, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)$$

# The Gaussian Mixture Model (GMM): 1D Example

**The whole model:**  $X = (1 - \Delta) \cdot X_1 + \Delta \cdot X_2$  We need to estimate:

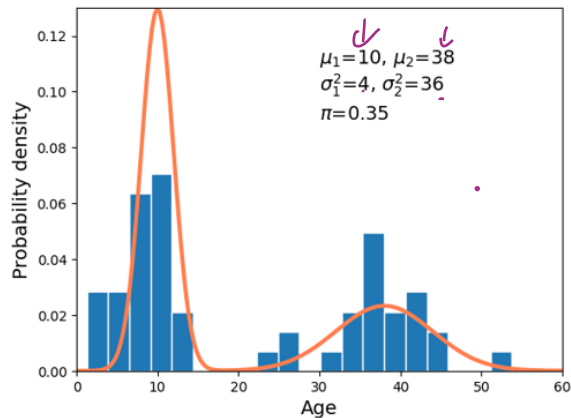
$$\Theta = (\pi, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)$$

Assuming we actually *know* all 5 parameters, we can simply write down the full probability density function for our process. If we denote  $\phi(x|\mu, \sigma^2)$  as the normal with mean  $\mu$  and variance  $\sigma^2$ , the model is now

$$f(x|\Theta) = (1 - \pi)\phi(x|\mu_1, \sigma_1^2) + \pi\phi(x|\mu_2, \sigma_2^2)$$

## GMM Example: Varying Theta

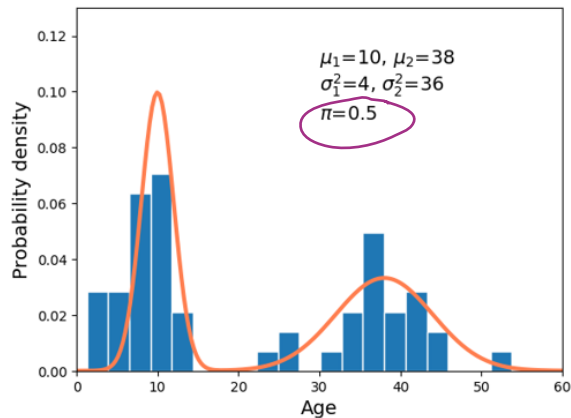
Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .





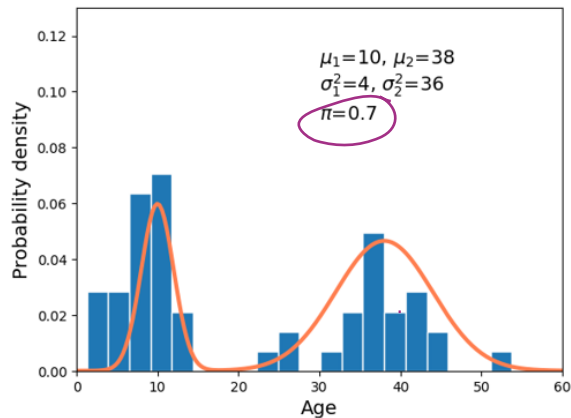
## GMM Example: Varying Theta

Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .



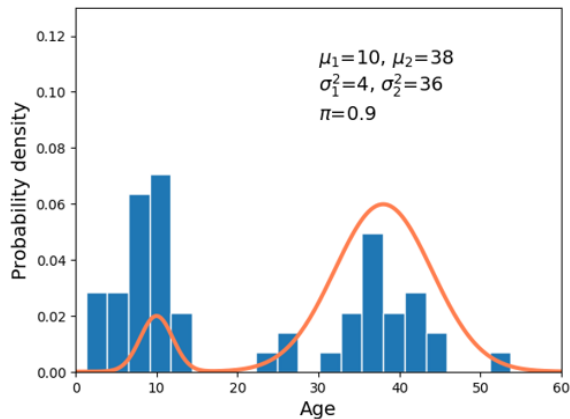
## GMM Example: Varying Theta

Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .



## GMM Example: Varying Theta

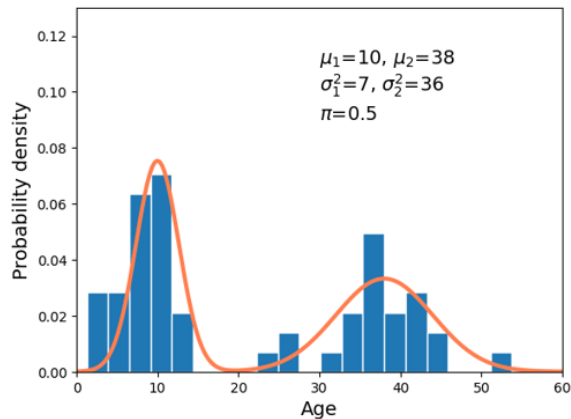
Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .



This is too many adults: let's go back to  $\pi = .5...$

## GMM Example: Varying Theta

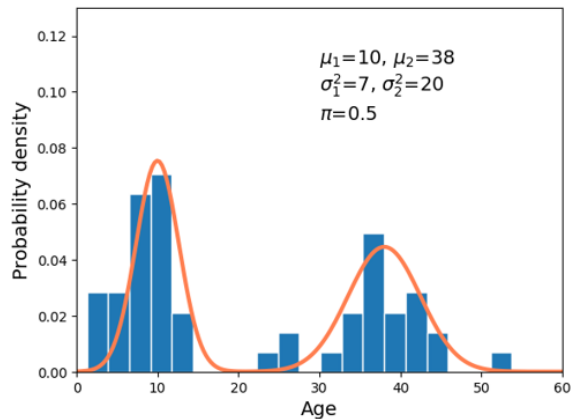
Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .



Changing  $\sigma_1$

## GMM Example: Varying Theta

Here are some pdfs, depending on different choices of the parameter set  $\Theta$ .



Changing  $\sigma_2$

## GMM Example: Underview

**The whole model:**  $X = (1 - \Delta) \cdot X_1 + \Delta \cdot X_2$  We need to estimate:

$$\Theta = (\pi, \mu_1, \sigma_1^2, \mu_2, \sigma_2^2)$$

1. It would be **tedious** and hand-wavy to manually try to estimate those parameters, even in one dimension.
2. Note that in 2D, each variance was a  $2 \times 2$  covariance matrix.

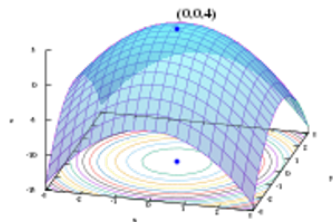
This problem grows in size *quickly*: 2 Gaussians in 2D is now ~~8~~<sup>n</sup> unknowns (4 per Gaussian plus a mixture probability).

$$\mu = (x_1, x_2)$$

$$\Sigma = \begin{pmatrix} \sigma_{x_1}^2 & \sigma_{x_1 x_2} \\ \sigma_{x_1 x_2} & \sigma_{x_2}^2 \end{pmatrix} \Rightarrow 5 \text{ terms per Gaussian}$$

# Optimization

We're going to zoom out for a bit from Gaussians and talk more generally about *estimating parameters*.



- ▶ We are often tasked with asking what the “best” values for a set of parameters is.
  - ▶ In regression: the  $\beta$  values.
  - ▶ In GMMs: the mixture probabilities and component variables
- ▶ The data scientist has to define **best**!
  - ▶ Least Squares?
  - ▶ Calculus for local extrema?
  - ▶ Minimized Loss? How do we define loss? Can we use cross-validation?
  - ▶ Fastest algorithm that's “almost” best?

# Probability Models

In the Gaussian mixture model, we're able to write down a **model**. We could answer exact questions like:

1. What are the exact probabilities associated with new data points? For example, suppose I got a new observation. If I knew the *parameters* ( $\Theta$ ), what would be the  $P(a < Y < b)$  for a new observation  $Y$ ?  $= \int_a^b f(x) dx$
2. What is the *joint distribution* of multiple new data points? Assume we get multiple new *independent* actualizations (draws/observations) from the model. They should **all** follow the exact probabilities  $P(a < Y_i < b)$ !

This is the crux of likelihood theory, which describes a method for estimating terms when we have a probability density function.



## pdfs

**Definition:** The Probability density function (pdf) of a random variable  $X$  is the function that describes the probability distribution of its outcomes. For (1D) continuous distributions, we add up *intervals* of outcomes on  $f$  to get probabilities, which turns into an integral:

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

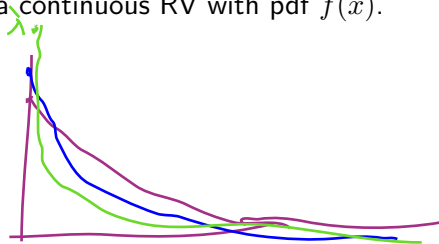
tells us the probability of all outcomes from  $a$  to  $b$  of a continuous RV with pdf  $f(x)$ .

A couple of examples:

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

2. Exponential:  $f(x; \lambda) = \lambda e^{-\lambda x}; \quad x \geq 0$

↑ start     ↑ decay



## On Parameters

Our goal in the GMM problem is to estimate *parameters* - for the Gaussian that's a mean, variance, and possibly covariances - but it may be other terms for other distributions.

The *Poisson* and *Exponential* distributions were based on using *rates* ( $\lambda$ ), and the *binomial* and *geometric* distributions were based on estimating probabilities ( $p$ ).

In the general notation we denote those parameters as the vector  $\Theta$ , and think of  $\Theta$  as the list of “knobs” or “dials” can adjust to change the shape of a distribution in a given model.

The problem with pdfs is that they're backwards of what we want. A pdf gives the “probability of *data* values given *parameter* values.” But the data scientist already has data, we need “probability of *parameter* values given *data* values.”

Bayes:

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

# Likelihoods

Notationally, we have  $f(x|\theta)$ : the pdf is the probabilities of data values  $x$  if we already know  $\theta$ . If we want to “flip” those things, Bayes’ theorem comes to the rescue!

$$P(\theta|x) = \frac{f(x|\theta)f(\theta)}{f(x)}$$

estimates depend on data  $x$ .

$\rightarrow$  Prob data from a distribution pdf

**Definition:** The *likelihood function* of random variable with pdf  $f(x; \Theta)$  and parameters  $\Theta$  with a set of observations  $\vec{x} = [x_1, x_2, \dots, x_n]$  is the probability of the parameters given the data;  $f(\theta|\vec{x})$ . If  $f(\theta)$  is a constant (an “uninformed prior”), it is *proportional* to the joint density  $f(\vec{x}|\theta)$

**Idea:** There’s a lot of theory there and notation here that we won’t cover in great detail, but the summary: to estimate parameters for a data set  $\vec{x}$  and a model/pdf  $f$ , we can ask “which parameter values would have led to the largest value of the pdf at exactly  $\vec{x}$ ?”.

## Likelihood Example



Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

**Example:** Suppose our *model* is that the data is an exponential random variable. Our goal is to estimate  $\lambda$ . We generate a sample of 4 values, and get  $x = [1, 3, 10, 5]$ .

1. What is  $P(\vec{x}|\lambda)$ ?

$$f(x) = \lambda e^{-\lambda x}$$

$$\text{if } \lambda = 4$$

$$f(3) = P(X=3) \text{ when } \lambda = 4$$

$$= 4 \cdot e^{-(3)4}$$

$$f(1), f(10), f(5).$$



2. How would we “choose  $\lambda$  to maximize this?”

## Likelihood Example

$$f(x) = \lambda e^{-\lambda x}$$

Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

**Example:** Suppose our *model* is that the data is an exponential random variable. Our goal is to estimate  $\lambda$ . We generate a sample of 4 values, and get  $x = [1, 3, 10, 5]$ .

1. What is  $P(\vec{x}|\lambda)$ ?

This is secretly an *and* statement: we have 4 *independent*  $x$  values.

$$\begin{aligned}
 P(\vec{x}|\lambda) &= P(\vec{x} = [x_1, x_2, x_3, x_4]|\lambda) = P(x_1 = 1 \cap x_2 = 3 \cap x_3 = 10 \cap x_4 = 5|\lambda) \\
 &= P(x_1=1) \cdot P(x_2=3) \cdot P(x_3=10) \cdot P(x_4=5)
 \end{aligned}$$

*(Handwritten note:  $\lambda = 4$  is circled above the second term)*

2. How would we “choose  $\lambda$  to maximize this?”



## Likelihood Example

Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

**Example:** Suppose our *model* is that the data is an exponential random variable. Our goal is to estimate  $\lambda$ . We generate a sample of 4 values, and get  $x = [1, 3, 10, 5]$ .

1. What is  $P(\vec{x}|\lambda)$ ?

This is secretly an *and* statement: we have 4 *independent*  $x$  values.

$$\begin{aligned} P(\vec{x}|\lambda) &= P(\vec{x} = [x_1, x_2, x_3, x_4]|\lambda) = P(x_1 = 1 \cap x_2 = 3 \cap x_3 = 10 \cap x_4 = 5|\lambda) \\ &= P(x_1 = 1|\lambda) \cdot P(x_2 = 3|\lambda) \cdot P(x_3 = 10|\lambda) \cdot P(x_4 = 5|\lambda) \\ &= \lambda e^{-\lambda 1} \cdot \lambda e^{-\lambda 3} \cdot \lambda e^{-\lambda 10} \cdot \lambda e^{-\lambda 5} \end{aligned}$$

2. How would we “choose  $\lambda$  to maximize this?”

## Likelihood Example

Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

**Example:** Suppose our *model* is that the data is an exponential random variable. Our goal is to estimate  $\lambda$ . We generate a sample of 4 values, and get  $x = [1, 3, 10, 5]$ .

1. What is  $P(\vec{x}|\lambda)$ ?

This is secretly an *and* statement: we have 4 *independent*  $x$  values.

$$\begin{aligned} P(\vec{x}|\lambda) &= P(\vec{x} = [x_1, x_2, x_3, x_4]|\lambda) = P(x_1 = 1 \cap x_2 = 3 \cap x_3 = 10 \cap x_4 = 5|\lambda) \\ &= P(x_1 = 1|\lambda) \cdot P(x_2 = 3|\lambda) \cdot P(x_3 = 10|\lambda) \cdot P(x_4 = 5|\lambda) \\ &= \lambda e^{-\lambda 1} \cdot \lambda e^{-\lambda 3} \cdot \lambda e^{-\lambda 10} \cdot \lambda e^{-\lambda 5} \end{aligned}$$

2. How would we “choose  $\lambda$  to maximize this?” **Hit it with the Calculus stick!**

## Likelihood Example

$$P(\vec{x}|\lambda) = \lambda e^{-\lambda 1} \cdot \lambda e^{-\lambda 3} \cdot \lambda e^{-\lambda 10} \cdot \lambda e^{-\lambda 5}$$

How would we “choose  $\lambda$  to maximize this?”

1. We'd probably collect the  $\lambda$ s at the start into  $\lambda^4$ .
2. We'd probably combine exponentials into  $e^{-\lambda(1+3+10+5)}$
3. We'd possibly realize that our goal is to pick the  $\lambda$  value that returns **maximum**. But products and exponentials are harder to deal with than sums, so we could take the *logarithm* of our function and find its maximum instead.

$$\ln[\lambda^4 \cdot e^{-\lambda(1+3+5+10)}]$$





## Likelihood Example

$$P(\vec{x}|\lambda) = \lambda e^{-\lambda 1} \cdot \lambda e^{-\lambda 3} \cdot \lambda e^{-\lambda 10} \cdot \lambda e^{-\lambda 5}$$

How would we “choose  $\lambda$  to maximize this?”

1. We'd probably collect the  $\lambda$ s at the start into  $\lambda^4$ .
2. We'd probably combine exponentials into  $e^{-\lambda(1+3+10+5)}$
3. We'd possibly realize that our goal is to pick the  $\lambda$  value that returns **maximum**. But products and exponentials are harder to deal with than sums, so we could take the *logarithm* of our function and find its maximum instead.

The **log-likelihood** of this problem is

$$\log \left( \lambda^4 e^{-\lambda(1+3+10+5)} \right) = \log \prod_{i=1}^4 f(x_i, \lambda)$$

## Likelihood Example

The **log-likelihood** of this problem is

$$\log \left( \lambda^4 e^{-\lambda(1+3+10+5)} \right) = \log \prod_{i=1}^4 f(x_i, \lambda)$$

To maximize with respect to  $\lambda$ , we differentiate and set equal to zero!

## Likelihood Example

The **log-likelihood** of this problem is

fire/cum

$$\log \left( \lambda^4 e^{-\lambda(1+3+10+5)} \right) = \log \prod_{i=1}^4 f(x_i, \lambda)$$

To maximize with respect to  $\lambda$ , we differentiate and set equal to zero!

$$LL(x, \theta) = \log \left( \lambda^4 e^{-\lambda(1+3+10+5)} \right);$$

$$\frac{d}{d\lambda} LL(x, \theta) = \frac{d}{d\lambda} \log \lambda^4 + \frac{d}{d\lambda} (-\lambda(1+3+10+5))$$

$$= \frac{4}{\lambda} - 19; \quad \text{now set to zero;}$$

$$0 = \frac{4}{\lambda} - 19;$$

$$\lambda = \frac{4}{19}$$

rate: events  
time.

So our estimate of the rate of these things is 4 events per 19 units of time.

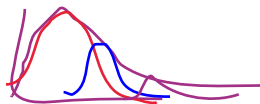
## The Gaussian Likelihood

Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

The setup to maximum likelihood is always the same: we look at all  $n$  of our data points  $x_1, x_2, \dots, x_n$  and ask about the probability of  $f(x = x_i | \mu, \sigma^2)$  for each one, then multiply them all together as an “and” or *joint* probability. Then we hit it with a logarithm to make maximization easier.

$$f(\vec{x} | \mu, \sigma^2) = \prod_{i=1}^n \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x_i - \mu)^2}{2\sigma^2}} \right)$$

Handwritten annotations: A purple circle around  $f(\vec{x} | \mu, \sigma^2)$  with an arrow pointing to the word "estimate". A red box around the product term with an arrow pointing to the text "each data point".



## The Gaussian Likelihood

Time to get to the real payoff: how do we estimate parameters if our data is **Gaussian**?

The setup to maximum likelihood is always the same: we look at all  $n$  of our data points  $x_1, x_2, \dots, x_n$  and ask about the probability of  $f(x = x_i | \mu, \sigma^2)$  for each one, then multiply them all together as an “and” or *joint* probability. Then we hit it with a logarithm to make maximization easier.

$$f(\vec{x} | \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x_i - \mu)^2}{2\sigma^2}}$$

$$LL(\mu, \sigma^2 | x) = \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{\frac{-(x_i - \mu)^2}{2\sigma^2}}$$

$$= \frac{-n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

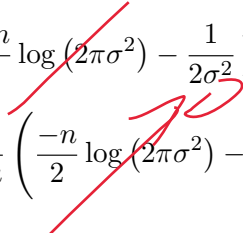
$\frac{\partial}{\partial \mu}$  AND  $\frac{\partial}{\partial \sigma^2}$



## The Gaussian Likelihood

Finally, we have to take two derivatives and set both  $\frac{d}{d\mu}LL$  and  $\frac{d}{d\sigma^2}LL$  equal to zero.

$$LL(\mu, \sigma^2 | x) = \frac{-n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{d}{d\mu}LL = \frac{d}{d\mu} \left( \frac{-n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$


## The Gaussian Likelihood

Finally, we have to take two derivatives and set both  $\frac{d}{d\mu}LL$  and  $\frac{d}{d\sigma^2}LL$  equal to zero.

$$LL(\mu, \sigma^2 | x) = \frac{-n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{d}{d\mu}LL = \frac{d}{d\mu} \left( \frac{-n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$

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

$$0 = \sum_{i=1}^n (x_i - \mu) = \sum_{i=1}^n (x_i) - \sum_{i=1}^n \mu$$

$$\Rightarrow n\mu = \sum_{i=1}^n (x_i) \Rightarrow \mu = \bar{X}$$

## The Gaussian Likelihood

What does that even mean? It means that the “best guess” of the mean  $\mu$  of the *probability density function* giving rise to our data was the sample mean  $\bar{x}$ . This is one of the measures as to why the sample mean is a **great measure** of population mean.

There are others, including the central limit theorem and the result that the sample mean is closest point in terms of sums of squared deviations (Euclidean distance) to each other point.

What about  $\sigma^2$ ?



## The Gaussian Likelihood

What does that even mean? It means that the “best guess” of the mean  $\mu$  of the *probability density function* giving rise to our data was the sample mean  $\bar{x}$ . This is one of the measures as to why the sample mean is a **great measure** of population mean.

There are others, including the central limit theorem and the result that the sample mean is closest point in terms of sums of squared deviations (Euclidean distance) to each other point.

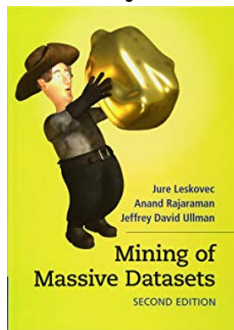
What about  $\sigma^2$ ?

The **maximum likelihood estimate** (MLEs) for  $\mu, \sigma^2$  of a Gaussian is  $\bar{x}, \frac{\sum (x_i - \bar{x})^2}{n}$ ... but it turns out MLEs are often slightly worse - especially for small samples - than the estimates you learned in your intro class. They're just a way to arrive at similar results: if you want to estimate the mean and variance of a normal, you use *sample* mean and *sample* variance.

And you can use *sample* covariance for covariances: 
$$Cov(\vec{x}, \vec{y}) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y}$$

## Acknowledgments

Some material is adapted/adopted from Mining of Massive Data Sets, by Jure Leskovec, Anand Rajaraman, Jeff Ullman (Stanford University) <http://www.mmds.org>



Special thanks to Tony Wong for sharing his original adaptation and adoption of slide material.