# Bayes Classification

**Alymzhan Toleu**

*alymzhan.toleu@gmail.com*

# Axioms of probability (Kolmogorov's axioms)

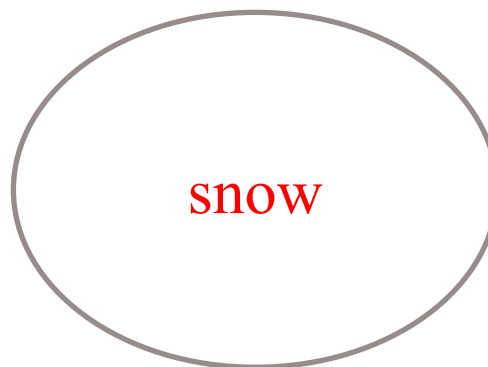A variety of useful facts can be derived from just three axioms:

- $0 \leq P(A) \leq 1$

- $P(\text{true}) = 1$, $P(\text{false}) = 0$

- Union: The union of two events is the probability that either A or B will occur.
  - $P(A \cup B) = P(A) + P(B) - P(A \cup B)$

- Intersection: The intersection of two events is the probability that the two events, A and B, will occur at the same time.
  - Independent events: $P(A \cap B) = P(A) * P(B)$

Two events are mutually exclusive if they cannot occur at the same time.

# Prior

- The **prior probability** of an event refers to the degree of belief assigned to that event before incorporating any additional information.

No snow

snow

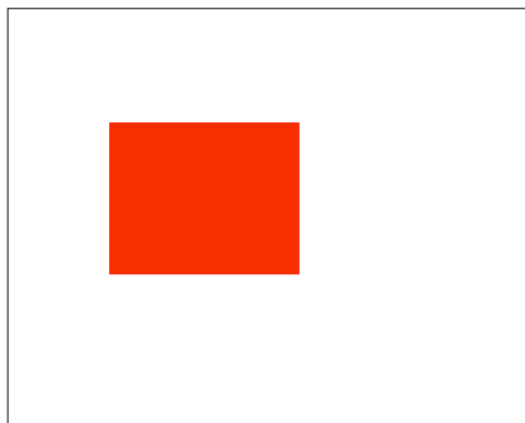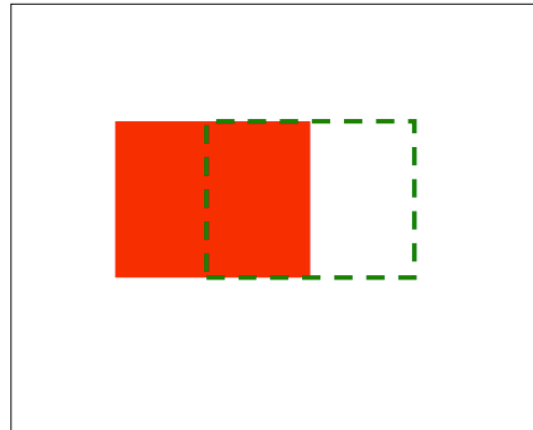P(snow tomorrow) = 0.2
P(no snow tomorrow) = 0.8

# Conditional probability

- P(A = true| B = false): The fraction of cases where A is true if B is false.
- For example:

P(A = 0.2)

P(A|B = 0.5)

# Conditional probability

- The prior belief of a random variable can be improved in some cases by conditioning on one or more other random variables.

- For example:

P(slept in moive) $= \frac{4}{7}$

P(slept in moive | liked movie) $= \frac{1}{4}$

P(did not slept in moive | liked movie) $= \frac{3}{4}$

| Slept | Liked |
|-------|-------|
| 1 | 1 |
| 0 | 1 |
| 1 | 0 |
| 1 | 0 |
| 0 | 1 |
| 1 | 0 |
| 0 | 1 |

# Joint distributions

- The probability that a set of random variables will take a specific value is their joint distribution.

- Notation: P(A ∩ B) or P(A,B)

- Example: P(liked movie, slept) = ?

If we assume independence then

P(A,B)=P(A)P(B)

However, in many cases such an assumption maybe too strong

# Joint distributions

$P(\text{Length} = \text{short}) = \dfrac{3}{7}$

$P(\text{slept in moive}) = \dfrac{4}{7}$

$P(\text{Length} = \text{short} , \text{slept in moive}) = ?$

| Length | Slept | Liked |
|--------|-------|-------|
| Short  | 1     | 1     |
| Long   | 0     | 1     |
| Medium | 1     | 0     |
| Short  | 1     | 0     |
| Medium | 0     | 1     |
| Short  | 1     | 0     |
| Long   | 0     | 1     |

# Joint distributions

$P(\text{Length} = \text{short}) = \dfrac{3}{7} = 0.42$

$P(\text{slept in moive}) = \dfrac{4}{7} = 0.57$

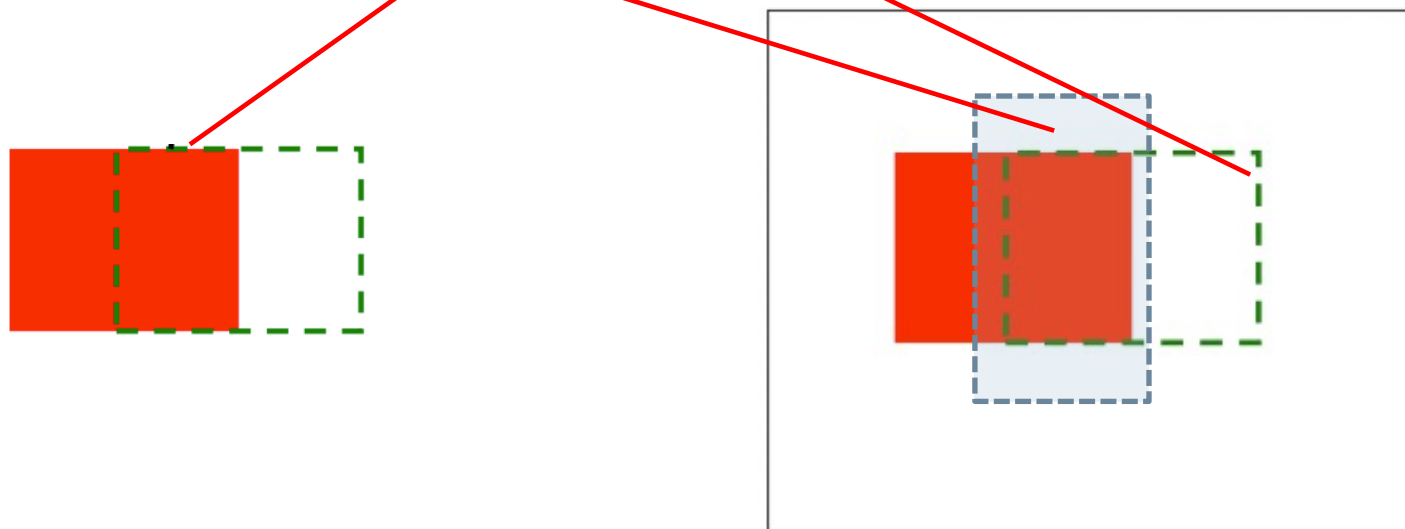$P(\text{Length} = \text{short , slept in moive}) = \dfrac{3}{7} = 0.42$

| Length | Slept | Liked |
|--------|-------|-------|
| Short | 1 | 1 |
| Long | 0 | 1 |
| Medium | 1 | 0 |
| Short | 1 | 0 |
| Medium | 0 | 1 |
| Short | 1 | 0 |
| Long | 0 | 1 |

# Chain rule

- The joint distribution can be specified in terms of conditional probability:

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

# Bayes decision rule

- One of the most important rules in probabilistic theory.
- Derived from the chain rule:

<p style="color:red; text-align:center;">P(A,B) = P(A | B)P(B) = P(B | A)P(A)</p>

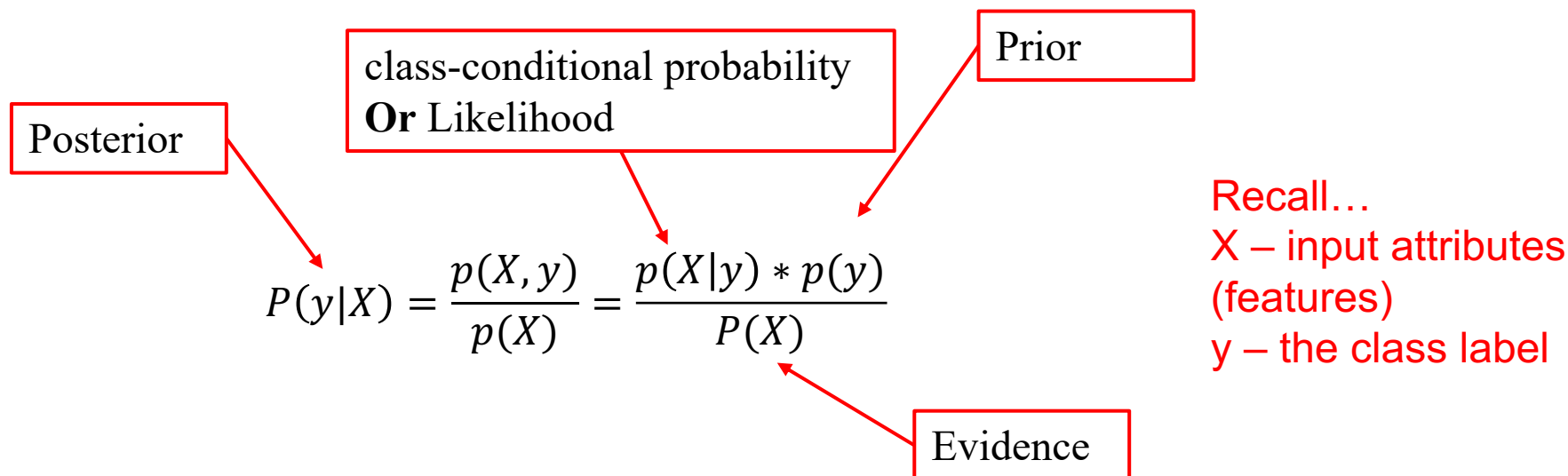- Thus, it becomes generative models

Joint probabilities

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

Thomas Bayes was an English clergyman who set out his theory of probability in 1764

# Bayes decision rule

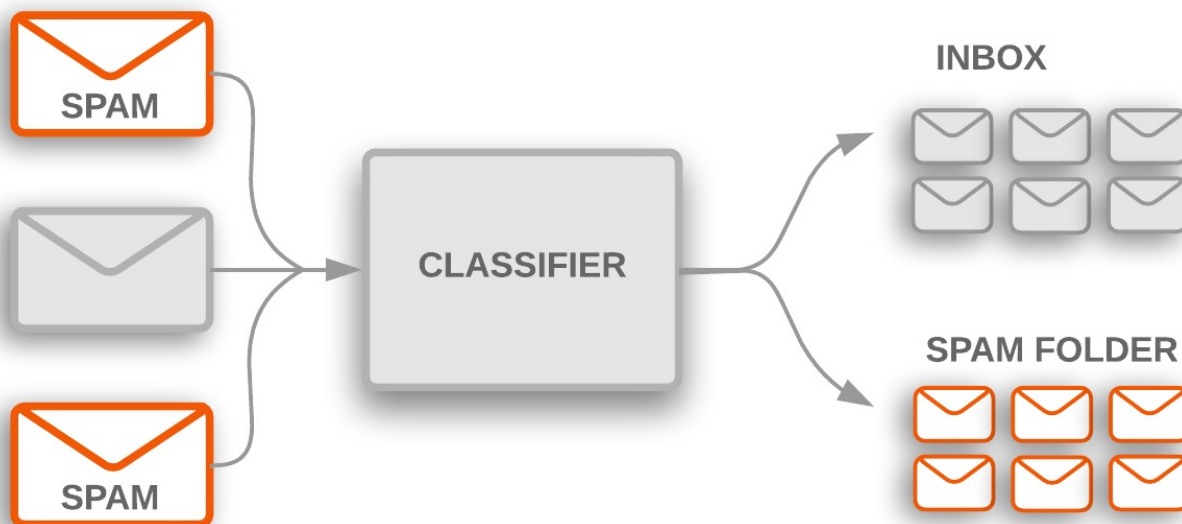- If we know the conditional probability P(X | y) we can determine the appropriate class by using Bayes rule:

Posterior

class-conditional probability
**Or** Likelihood

Prior

$$P(y|X) = \frac{p(X,y)}{p(X)} = \frac{p(X|y) * p(y)}{P(X)}$$

Evidence

Recall…
X – input attributes (features)
y – the class label

But how do we determine p(X|y)?

# Types of classifier

- We can divide the large variety of classification approaches into roughly three main types

    - Discriminative
      - directly estimate a decision rule/boundary - e.g., decision tree

    - Instance based classifiers
      - use observation directly (no models) - e.g. K nearest neighbors

    - <span style="color:red">Generative</span>
      - <span style="color:red">build a generative statistical model - e.g., Bayesian networks</span>

# Spam detection: Example

Normal Emails (8)

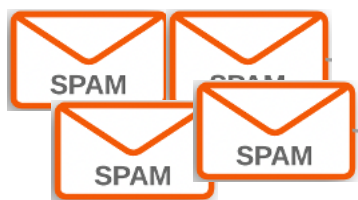| Words | Frequncy |
|-------|----------|
| Dear | 8 |
| Friend | 5 |
| Lunch | 3 |
| Money | 1 |

The probability of each word  we see in the normal messages:

$$P(\textbf{Dear}|Normal) = \frac{8}{17} = 0.47$$

$$P(\textbf{Friend}|Normal) = \frac{5}{17} = 0.29$$

$$P(\textbf{Lunch}|Normal) = \frac{3}{17} = 0.18$$

$$P(\textbf{Money}|Normal) = \frac{1}{17} = 0.06$$

Normal Emails (8)

| Words | Frequncy |
|-------|----------|
| Dear | 8 |
| Friend | 5 |
| Lunch | 3 |
| Money | 1 |

The probability of each word we see in the normal messages:

$$P(\textbf{Dear}|Normal) = \frac{8}{17} = 0.47$$

$$P(\textbf{Friend}|Normal) = \frac{5}{17} = 0.29$$

$$P(\textbf{Lunch}|Normal) = \frac{3}{17} = 0.18$$

$$P(\textbf{Money}|Normal) = \frac{1}{17} = 0.06$$

Likewaise, we calculate the probability of each word we see in the spam:

| Words | Frequncy |
|-------|----------|
| Dear | 2 |
| Friend | 1 |
| Lunch | 0 |
| Money | 4 |

$$P(\textbf{Dear}|Spam) = \frac{2}{7} = 0.29$$

$$P(\textbf{Friend}|Spam) = \frac{1}{7} = 0.06$$

$$P(\textbf{Lunch}|Spam) = \frac{0}{7} = 0$$

$$P(\textbf{Money}|Spam) = \frac{4}{7} = 0.57$$
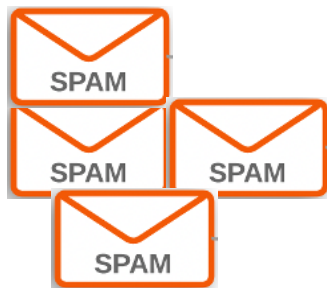
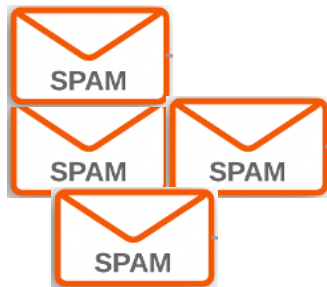Spams(4)

Normal Emails (8)

$$P(\textbf{Dear}|Normal) = \frac{8}{17} = 0.47$$

$$P(\textbf{Friend}|Normal) = \frac{5}{17} = 0.29$$

$$P(\textbf{Lunch}|Normal) = \frac{3}{17} = 0.18$$

$$P(\textbf{Money}|Normal) = \frac{1}{17} = 0.06$$

These probabilities can be called **Likelihood.** Because they are probabilties of discrete, and dot probability of something continuous.

$$P(\textbf{Dear}|Spam) = \frac{2}{7} = 0.29$$

$$P(\textbf{Friend}|Spam) = \frac{1}{7} = 0.06$$

$$P(\textbf{Lunch}|Spam) = \frac{0}{7} = 0$$

$$P(\textbf{Money}|Spam) = \frac{4}{7} = 0.57$$

Spams (4)

## Normal Emails



$$P(\boldsymbol{Dear}|Normal) = \frac{8}{17} = 0.47$$

$$P(\boldsymbol{Friend}|Normal) = \frac{5}{17} = 0.29$$

$$P(\boldsymbol{Lunch}|Normal) = \frac{3}{17} = 0.18$$

$$P(\boldsymbol{Money}|Normal) = \frac{1}{17} = 0.06$$



$$P(\boldsymbol{Dear}|Spam) = \frac{2}{7} = 0.29$$

$$P(\boldsymbol{Friend}|Spam) = \frac{1}{7} = 0.06$$

$$P(\boldsymbol{Lunch}|Spam) = \frac{0}{7} = 0$$

$$P(\boldsymbol{Money}|Spam) = \frac{4}{7} = 0.57$$

Spams

Let's say, we got a new message that said:

Dear Friend

We want to classify if it is a Normal or Spam message.

What we do?

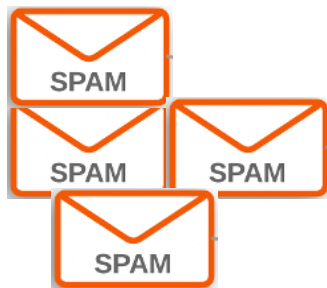$$P(N) * P(Dear|N) * P(Friend|N) = 0.09$$

## Normal Emails (8)

$$P(Dear|N) = \frac{8}{17} = 0.47$$

$$P(Friend|N) = \frac{5}{17} = 0.29$$

$$P(Lunch|N) = \frac{3}{17} = 0.18$$

$$P(Money|N) = \frac{1}{17} = 0.06$$

$$P(Dear|S) = \frac{2}{7} = 0.29$$

$$P(Friend|S) = \frac{1}{7} = 0.06$$

$$P(Lunch|S) = \frac{0}{7} = 0$$

$$P(Money|S) = \frac{4}{7} = 0.57$$

Spams (4)

Let's say, we got a new message that said:

Dear Friend
………..
……..

First, we need to calculate the initial guess (**prior probability**) of a message, regardless what it says, is a normal message.

$$P(N) = \frac{8}{8 + 4} = 0.66$$

$$P(N) * P(Dear|N) * P(Friend|N)$$
$$= 0.66. * 0.47 * 0.29 = 0.09$$

It is the score is given to the message "Dear Friend" if it is a normal.

$$P(N) * P(Dear|N) * P(Friend|N) = 0.09$$

## Normal Emails (8)



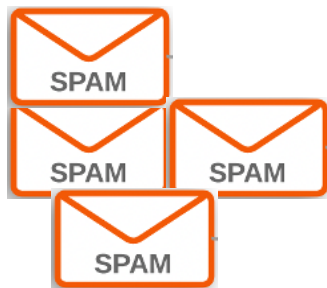$$P(Dear|N) = \frac{8}{17} = 0.47$$

$$P(Friend|N) = \frac{5}{17} = 0.29$$

$$P(Lunch|N) = \frac{3}{17} = 0.18$$

$$P(Money|N) = \frac{1}{17} = 0.06$$

$$P(Dear|S) = \frac{2}{7} = 0.29$$

$$P(Friend|S) = \frac{1}{7} = 0.06$$

$$P(Lunch|S) = \frac{0}{7} = 0$$

$$P(Money|S) = \frac{4}{7} = 0.57$$

Spams (4)

Let's say, we got a new message that said:

> Dear Friend
> ………..
> ……..

First, we need to calculate the initial guess **(prior probability)** of a message, regardless what it says, is a Spam message.

$$P(S) = \frac{4}{8 + 4} = 0.33$$

$$P(S) * P(Dear|S) * P(Friend|S)$$
$$= 0.33 * 0.29 * 0.06 = \boxed{0.01}$$

It is the score is given to the message "Dear Friend" if it is a Spam.

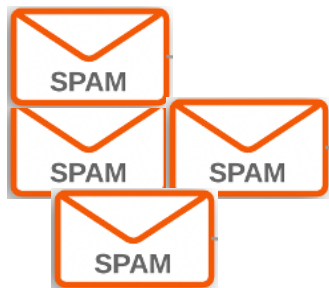$$P(N) * P(\textbf{Dear}|N) * P(\textbf{Friend}|N) = 0.09$$

## Normal Emails (8)

$$P(\textbf{Dear}|N) = \frac{8}{17} = 0.47$$

$$P(\textbf{Friend}|N) = \frac{5}{17} = 0.29$$

$$P(\textbf{Lunch}|N) = \frac{3}{17} = 0.18$$

$$P(\textbf{Money}|N) = \frac{1}{17} = 0.06$$

$$P(\textbf{Dear}|S) = \frac{2}{7} = 0.29$$

$$P(\textbf{Friend}|S) = \frac{1}{7} = 0.06$$

$$P(\textbf{Lunch}|S) = \frac{0}{7} = 0$$

$$P(\textbf{Money}|S) = \frac{4}{7} = 0.57$$

Spams (4)

Let's say, we got a new message that said:

Dear Friend
………..
……..

$$P(N) * P(\textbf{Dear}|N) * P(\textbf{Friend}|N) = 0.09$$

$$P(S) * P(\textbf{Dear}|S) * P(\textbf{Friend}|S) = 0.01$$

$$0.09 \ (N) > 0.01 \ (S)$$

This is a normal message.

This is the basics how **Naïve Bayes Classification** works.
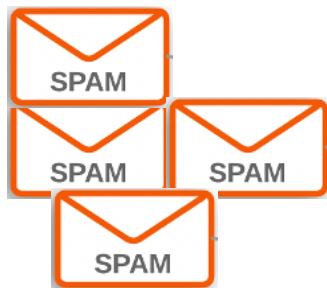
# The problem of zero probability

Normal Emails (8)

$$P(\boldsymbol{Dear}|N) = \frac{8}{17} = 0.47$$

$$P(\boldsymbol{Friend}|N) = \frac{5}{17} = 0.29$$

$$P(\boldsymbol{Lunch}|N) = \frac{3}{17} = 0.18$$

$$P(\boldsymbol{Money}|N) = \frac{1}{17} = 0.06$$

$$P(\boldsymbol{Dear}|S) = \frac{2}{7} = 0.29$$

$$P(\boldsymbol{Friend}|S) = \frac{1}{7} = 0.06$$

$$P(\boldsymbol{Lunch}|S) = \frac{0}{7} = 0$$

$$P(\boldsymbol{Money}|S) = \frac{4}{7} = 0.57$$

Spams (4)

Let's say, we got a new message that said:

Lunch Money Money Money Money
………..

$$P(N) * P(\boldsymbol{Lunch}|N) * P(\boldsymbol{Money}|N)^4 = 0.000002$$

$$P(S) * P(\boldsymbol{Lunch}|S) * P(\boldsymbol{Money}|S)^4 = 0$$

$$0.00002(N) > 0 \ (S)$$

This is a normal message.

# The problem of zero probability

- To avoid this issue, there is an approach called **smoothing technique.**

- A small number of counts ($\alpha$, alpha) will add to each sample (word).

- Make sure there are no zero probability.

# Smoothing technique

$set\ \alpha = 1$

Normal Emails (8)

| Words | Frequncy |
|-------|----------|
| Dear | $8+\alpha$ |
| Friend | $5+\alpha$ |
| Lunch | $3+\alpha$ |
| Money | $1+\alpha$ |

$$P(\textbf{Dear}|N) = \frac{9}{21} = 0.43$$

$$P(\textbf{Friend}|N) = \frac{6}{21} = 0.29$$

$$P(\textbf{Lunch}|N) = \frac{4}{21} = 0.19$$

$$P(\textbf{Money}|N) = \frac{2}{21} = 0.1$$

| Words | Frequncy |
|-------|----------|
| Dear | $2+\alpha$ |
| Friend | $1+\alpha$ |
| Lunch | $0+\alpha$ |
| Money | $4+\alpha$ |

Spams(4)

$$P(\textbf{Dear}|S) = \frac{3}{11} = 0.27$$

$$P(\textbf{Friend}|S) = \frac{2}{11} = 0.18$$

$$P(\textbf{Lunch}|S) = \frac{1}{11} = 0.09$$

$$P(\textbf{Money}|S) = \frac{5}{11} = 0.45$$
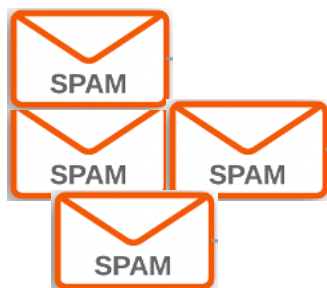
# Smoothing technique

Normal Emails (8)



$$P(\textbf{Dear}|N) = \frac{9}{21} = 0.43$$

$$P(\textbf{Friend}|N) = \frac{6}{21} = 0.29$$

$$P(\textbf{Lunch}|N) = \frac{4}{21} = 0.19$$

$$P(\textbf{Money}|N) = \frac{2}{21} = 0.1$$



$$P(\textbf{Dear}|S) = \frac{3}{11} = 0.27$$

$$P(\textbf{Friend}|S) = \frac{2}{11} = 0.18$$

$$P(\textbf{Lunch}|S) = \frac{1}{11} = 0.09$$

$$P(\textbf{Money}|S) = \frac{5}{11} = 0.45$$

Spams (4)

Let's say, we got a new message that said:

Lunch Money Money Money Money
………..

$$P(N) * P(\textbf{Lunch}|N) * P(\textbf{Money}|N)^4 = 0.00001$$

$$P(S) * P(\textbf{Lunch}|S) * P(\textbf{Money}|S)^4 = 0.00122$$

$$0.00001(N) < 0.00122 \text{ (S)}$$

This is a spam message.

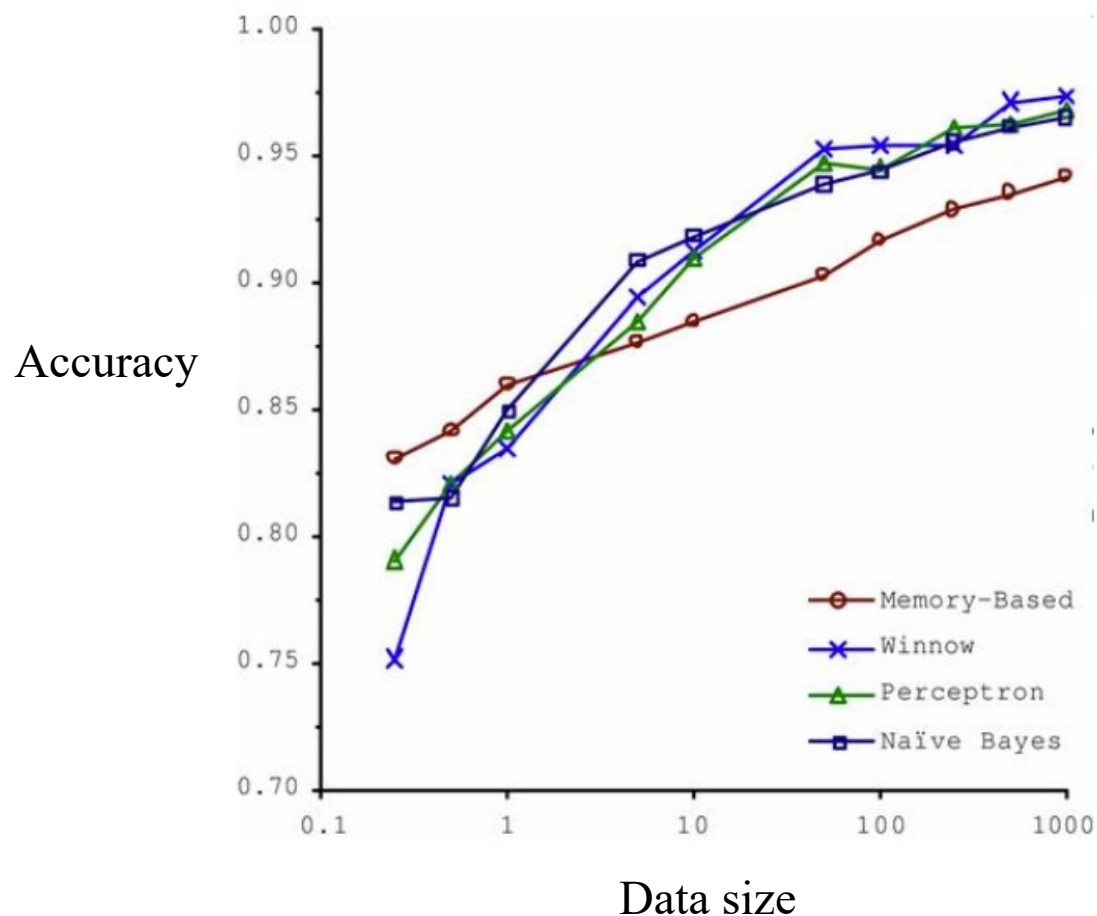- ## Why Naïve Bayes Classification is **Naïve?**

# Word order issue

- It treats each email as a bag of words.

- In other words, it treats all word orders the same.

- It ignores word orders.

- For example:

  - Score fir **Dear Friend** and **Friend Dear** are same.

$$P(N) * P(\boldsymbol{Dear}|N) * P(\boldsymbol{Friend}|N) = 0.09$$

$$P(N) * P(\boldsymbol{Friend}|N) * P(\boldsymbol{Dear}|N) = 0.09$$

# Data size is matter

# Maximum Likelihood Estimation

# Likelihood vs. Probability

- **Probability** corresponds to finding the chance of something given a sample distribution of the data.

- **Likelihood** refers to finding the best distribution of the data given a particular value of some feature or some situation in the data.
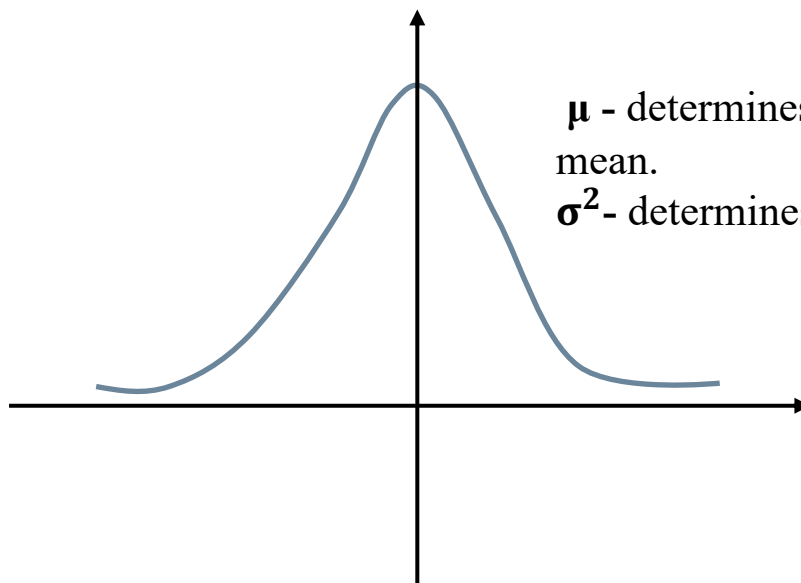
**Probability** $\quad P(data \,|\, distribution)$

**Likelihood** $\quad$ Likelihood$(distribution \,|\, data )$

# The Normal Distribution

- The Gaussian, also known as the normal distribution, is a widely used model for the distribution of continuous variables.

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

where $\mu$ is the mean and $\sigma^2$ is the variance.

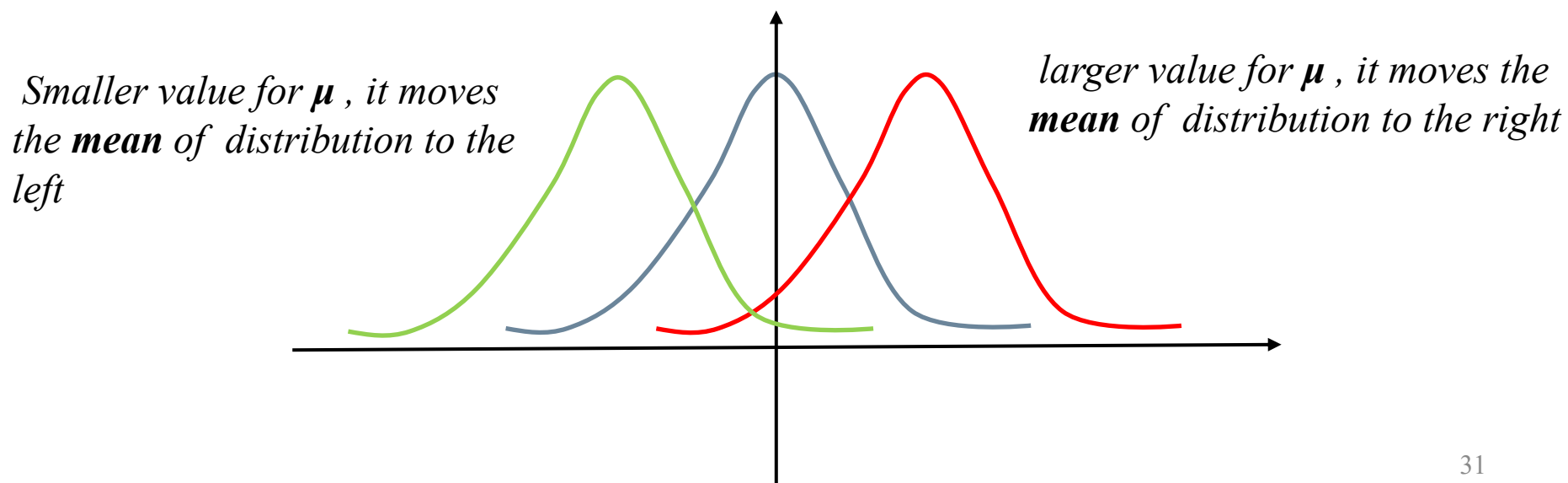**μ** - determines the location of the normal distribution's mean.

**σ²** - determines the normal distribution's width.

30

# The Normal Distribution

- The Gaussian, also known as the normal distribution, is a widely used model for the distribution of continuous variables.

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

where $\mu$ is the <span style="color:red">mean</span> and $\sigma^2$ is the <span style="color:blue">variance</span>.

*Smaller value for $\boldsymbol{\mu}$ , it moves the **mean** of distribution to the left*

*larger value for $\boldsymbol{\mu}$ , it moves the **mean** of distribution to the right*
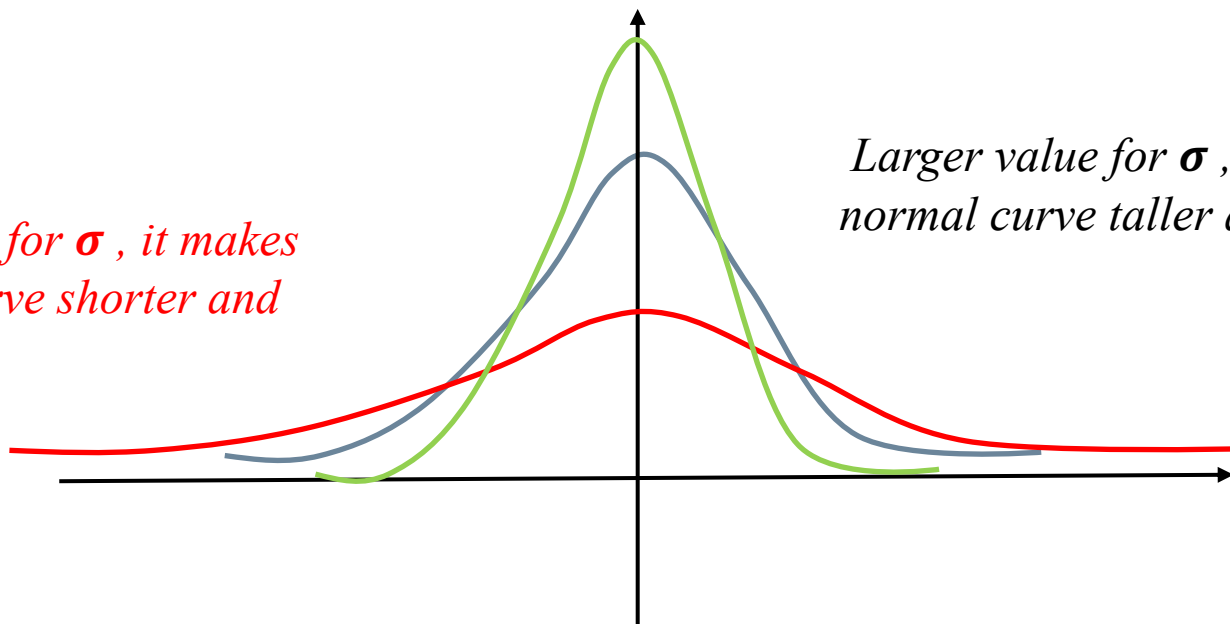
# The Normal Distribution

- The Gaussian, also known as the normal distribution, is a widely used model for the distribution of continuous variables.

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

where $\mu$ is the mean and $\sigma^2$ is the variance.

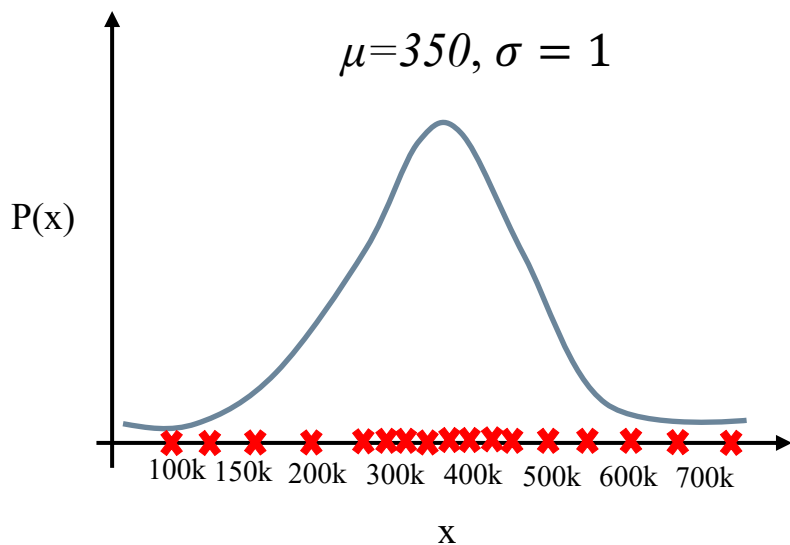*Larger value for $\boldsymbol{\sigma}$, it makes the normal curve taller and narrower.*

*Smaller value for $\boldsymbol{\sigma}$, it makes the normal curve shorter and wider.*

# Likelihood vs. Probability: Example

Distribution of employee's income in a company.



$\mu=350, \sigma = 1$

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

# Likelihood vs. Probability: Example

Distribution of employee's income in a company.



$\mu=0, \sigma = 1$

P(x)

100k 150k  200k   300k  400k   500k  600k  700k

x

Total area under the curve.

$$\int p(x) = 1$$

# Likelihood vs. Probability: Example

Distribution of employee's income in a company.

$\mu=0, \sigma = 1$

P(x)

100k 150k  200k  300k  400k  500k  600k  700k

x

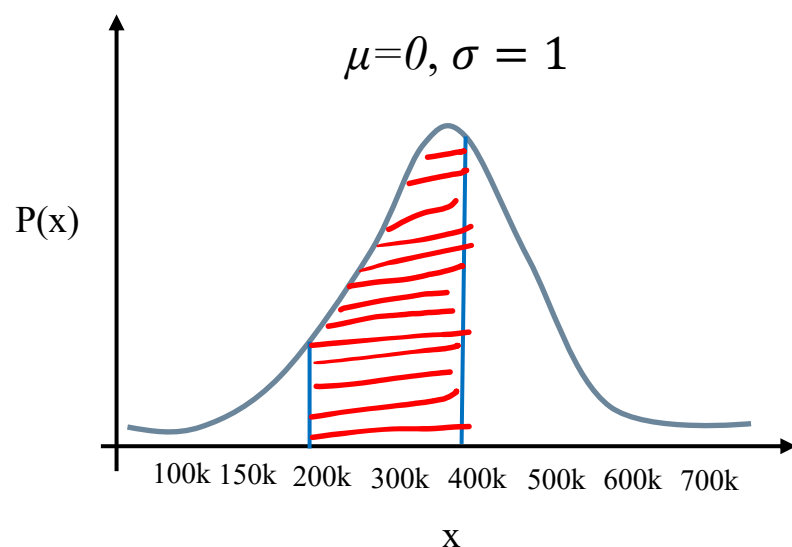Probability of employee's income between 200k and 400k:

$$\int_{200k}^{400k} p(x)dx$$

What is the relation between likelihood and probability here?

# Likelihood vs. Probability: Example

## Probability
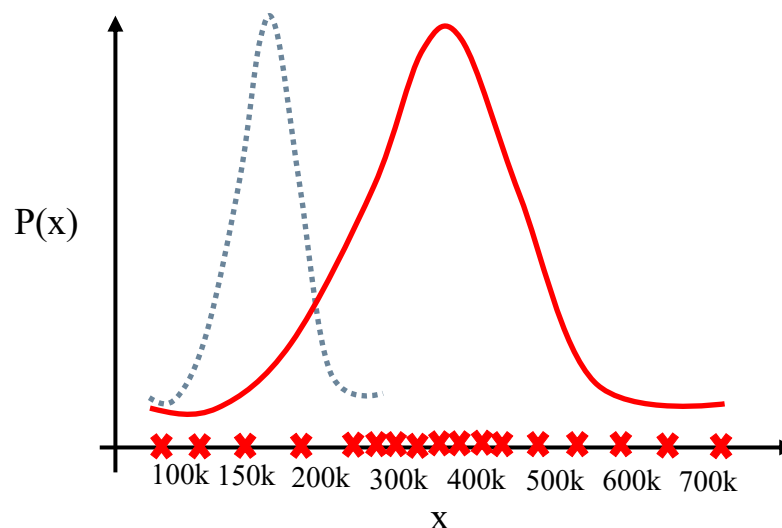
Distribution of income in an organization.

$\mu=0, \sigma = 1$

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

Probability of employee's income between 200k and 400k:

$$\int_{200k}^{400k} p(x)dx$$

## Likelihood

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

We want to find some distribution fits to this data.
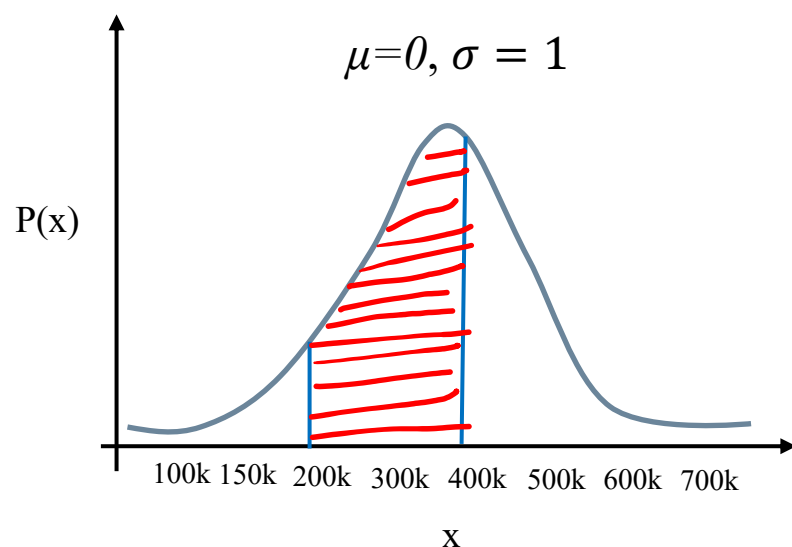
$\mu=?, \sigma =?$

For example:

$\mu=10, \sigma = 0.5$

$\mu=20, \sigma = 2$

36

# Likelihood vs. Probability: Example
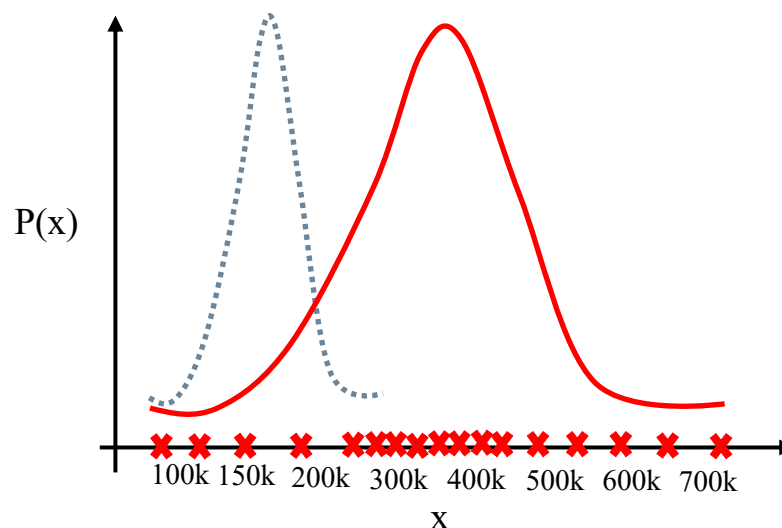
## Probability

Distribution of income in an organization.

$\mu=0,\ \sigma = 1$

P(x)



100k 150k 200k 300k 400k 500k 600k 700k

x

The probability of employee's income between 200k and 400k:

$$\int_{200k}^{400k} p(x)dx$$

## Likelihood

P(x)



100k 150k 200k 300k 400k 500k 600k 700k

x

We want to find some distribution fit to this data.

$\mu=?,\ \sigma =?$
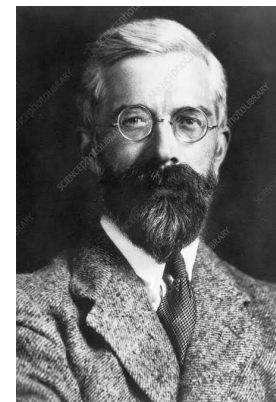
Add likelihood function:

$L\ (\mu=10,\ \sigma = 0.5) = ?$
$L\ (\mu=20,\ \sigma = 2) = ?$

# Definition: Likelihood

"The likelihood that any parameter (or set of parameters) should have any assigned value (or set of values) is proportional to the probability that if this were so, the totality of observations should be that observed."

— Fisher, 1922

Ronald Fisher

# Definition: Likelihood

"The likelihood that <u>any parameter (or set of parameters) should have any assigned value (or set of values)</u> is proportional to the probability that if this were so, the totality of observations should be that observed."
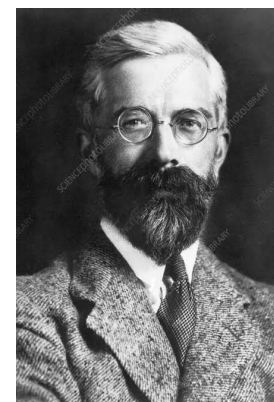
— Fisher, 1922

Ronald Fisher

$$L(\mu{=}10, \sigma = 0.5) \quad \propto \quad P\big(x^{(1)} = 100k, x^{(2)} = 200k, x^{(3)} = 300k \ldots x^{(n)}\big| \mu{=}10, \sigma = 0.5\big)$$

We pass in the parameters of the distribution

The likelihood value going to be proportional to the probability of observing all of these examples given the parameters of the assumed distribution

# Maximize the Likelihood

- Find the values of theta that is going to maximize the likelihood function.

$$\hat{\theta}^{MLE} = argmax\ \boldsymbol{L}(\theta)$$

# Maximize the Likelihood

- Likelihood function

$$L(\theta) \propto P(D \mid \theta) \quad \Longrightarrow \quad \hat{\theta}^{MLE} = argmax\ L(\theta)$$

$$L(\theta) \propto P\big(x^{(1)} = 100k, x^{(2)} = 200k, x^{(3)} = 300k\ ...\ x^{(n)}\big|\ \theta\big)$$

Assume these incomes are independent and identically distributed (i.i.d.).

$$L(\theta) \propto P\big(x^{(1)}\big|\ \theta\big) * P\big(x^{(2)}\big|\ \theta\big)\ * P\big(x^{(3)}\big|\ \theta\big)\ * \cdots * P\big(x^{(n)}\big|\ \theta\big))$$

$$L(\theta) \propto \prod_{1}^{n} P\big(x^{(i)}\big|\ \theta\big) \quad \longleftarrow \quad \text{Multiplications lead to the arithmetic underflow}$$

The term arithmetic underflow is a condition in a computer program where the result of a calculation is a number of more precise absolute value than the computer can actually represent in memory on its central processing unit (CPU).

# Maximize the Likelihood

- Likelihood function

$$L(\theta) \propto P(D \mid \theta) \implies \hat{\theta}^{MLE} = argmax \, L(\theta)$$

$$L(\theta) \propto P\left(x^{(1)} = 100k, x^{(2)} = 200k, x^{(3)} = 300k \ldots x^{(n)} \middle| \theta\right)$$

Assume these incomes are independent and identically distributed (i.i.d.).

$$L(\theta) \propto P\left(x^{(1)} \middle| \theta\right) * P\left(x^{(2)} \middle| \theta\right) * P\left(x^{(3)} \middle| \theta\right) * \cdots * P\left(x^{(n)} \middle| \theta\right))$$

$$L(\theta) \propto \prod_{1}^{n} P\left(x^{(i)} \middle| \theta\right) \quad \longleftarrow$$

Multiplications lead to the arithmetic underflow

Taking the logarithms on both sides: $\quad log(L(\theta)) \propto log(\prod_{i=1}^{n} P\left(x^{(i)} \middle| \theta\right))$

$$log(L(\theta)) \propto \sum_{i=1}^{n} log(P\left(x^{(i)} \middle| \theta\right))$$

# Maximize the Likelihood

- Likelihood function

$$L(\theta) \propto P(D \mid \theta) \quad \Longrightarrow \quad \hat{\theta}^{MLE} = argmax\ L(\theta)$$

$$log(L(\theta)) \propto \sum_{i=1}^{n} log(P(x^{(i)}\mid \theta))$$

$$\hat{\theta}^{MLE} = \ argmax \sum_{i=1}^{n} log(P(x^{(i)}\mid \theta))$$

Maximize Likelihood Estimation

# Maximum Likelihood
# For the Normal Distribution

# Maximum Likelihood For the Normal Distribution

- the Normal Distribution

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

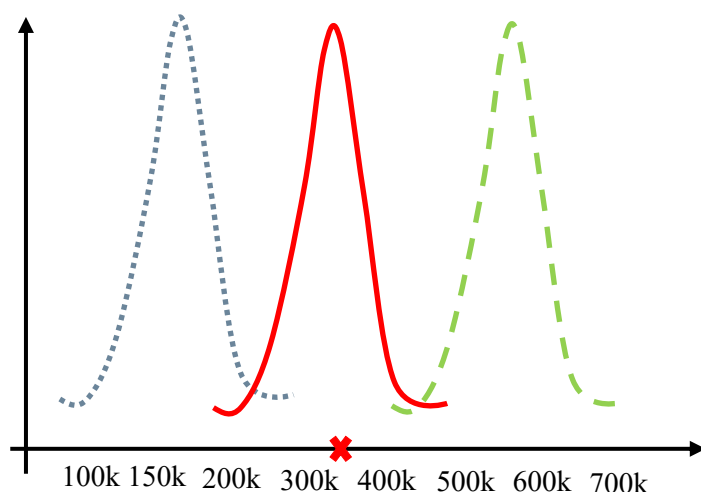where $\mu$ is the mean and $\sigma^2$ is the variance.

- the Likelihood of the normal Distribution

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

# Example (cont.)

Distribution of income in an organization.

For simplicity, we assume there is only one employee: $x^{(1)}=350$



```
1  mus = [200,300,350,400,500]
2  sigma = 0.5
3
4  pvs = []
5  for mu in mus:
6      pv = scipy.stats.norm(mu, sigma).pdf(350
7      pvs.append(pv)
```

```
1  plt.xlabel("u")
2  plt.ylabel("likelihood")
3  plt.plot(mus,pvs)
```

`[<matplotlib.lines.Line2D at 0x166b06490>]`

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

$L\ (\mu=100, \sigma = 0.5) = ?$

$L\ (\mu=300, \sigma = 0.5) = ?$
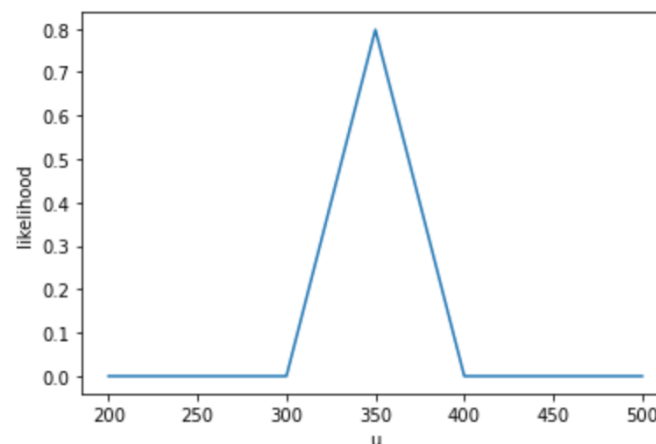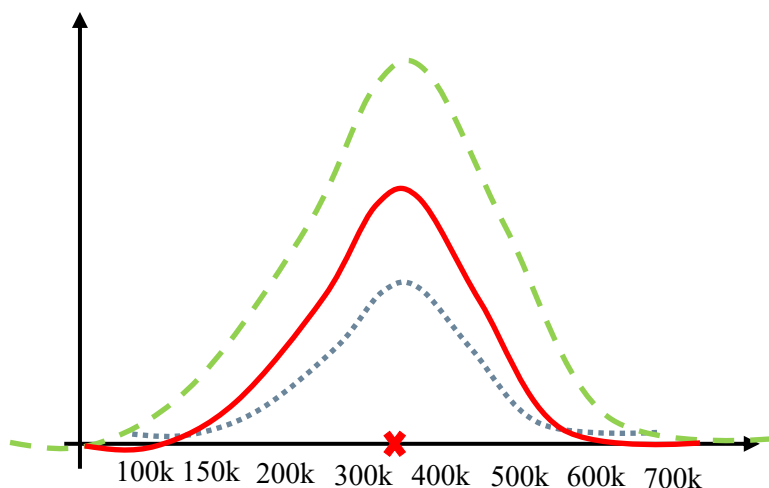
$L\ (\mu=500, \sigma = 0.5) = ?$
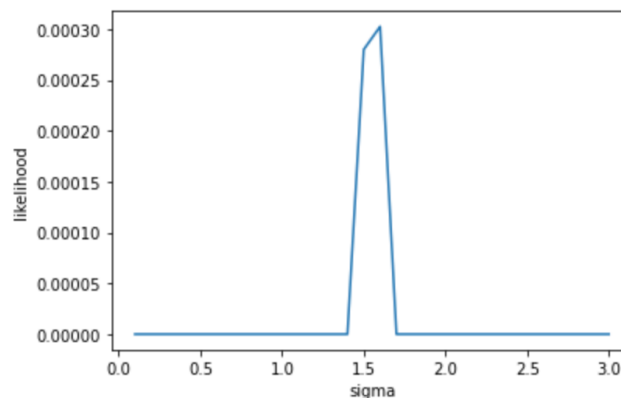
46

# Example (cont.)

Distribution of income in an organization.

For simplicity, we assume there is only one employee: $x^{(1)}=350$



```
1
2  sigmas = np.linspace(start=0.1, stop=3, num=30)
3
4  xvs = np.linspace(start=100, stop=700, num=30,dtype=int)
5  mu = np.sum(xvs)/len(xvs)
6
7  pvs = []
8  for x in xvs:
9      _pvs = []
10     for sigma in sigmas:
11         pv = scipy.stats.norm(mu, sigma).pdf(x)
12         _pvs.append(pv)
13     pvs.append(np.max(_pvs))
```

```
1  plt.xlabel("sigma")
2  plt.ylabel("likelihood")
3  plt.plot(sigmas,pvs)
```

]: [<matplotlib.lines.Line2D at 0x166b38bb0>]

$$L(\sigma \mid \mu, x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$
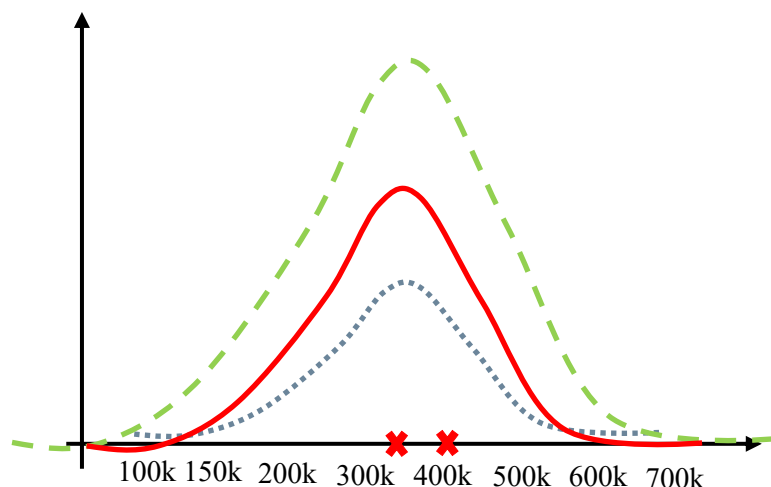
$L(\mu{=}200, \sigma = 0.5) = ?$

$L(\mu{=}200, \sigma = 1) = ?$

$L(\mu{=}200, \sigma = 2) = ?$



47

# Example (cont.)

Distribution of income in an organization.

For simplicity, we assume there are two employees: : $x^{(1)}$=350 and $x^{(2)}$=400



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

$$L(\mu, \sigma \mid x^{(1)} = 350, x^{(2)} = 400)$$

$$= L(\mu, \sigma \mid x^{(1)} = 350) * L(\mu, \sigma \mid x^{(2)} = 400)$$

If there are $n$ employees

$$L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \ldots, x^{(n)}) = L(\mu, \sigma \mid x^{(1)}) * L(\mu, \sigma \mid x^{(2)}) * \cdots * L(\mu, \sigma \mid x^{(n)})$$

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

# Example (cont.)

$$L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \ldots, x^{(n)})$$

$$= L(\mu, \sigma \mid x^{(1)}) \, * \, L(\mu, \sigma \mid x^{(2)}) * \cdots * L(\mu, \sigma \mid x^{(n)})$$

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

Taking the logarithms on both sides

$$ln(L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \ldots, x^{(n)})) \; = ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^{(1)}-\mu)^2}{2\sigma^2}} * \cdots * \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^{(n)}-\mu)^2}{2\sigma^2}} \right)$$

… apply log transormations…

$$= \frac{1}{2}\ln(2\pi) - \ln(\sigma) - \frac{\left(x^{(1)}-\mu\right)^2}{2\sigma^2} - \cdots - \frac{1}{2}\ln(2\pi) - \ln(\sigma) - \frac{\left(x^{(1)}-\mu\right)^2}{2\sigma^2}$$

$$= \frac{n}{2}\ln(2\pi) - n\ln(\sigma) - \frac{\left(x^{(1)}-\mu\right)^2}{2\sigma^2} - \cdots - \frac{\left(x^{(1)}-\mu\right)^2}{2\sigma^2}$$

Likelihood function of the Normal Distribution: Example (cont.)

$$L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \frac{n}{2}\ln(2\pi) - n\ln(\sigma) - \frac{(x^{(1)} - \mu)^2}{2\sigma^2} - \dots - \frac{(x^{(1)} - \mu)^2}{2\sigma^2}$$

# Partial derivatives:

$$\frac{\partial L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \dots, x^{(n)})}{\partial \mu} = 0 - 0 + \frac{(x^{(1)} - \mu)}{\sigma^2} + \dots + \frac{(x^{(n)} - \mu)}{\sigma^2}$$

$$= \frac{1}{\sigma^2}[(x^{(1)} + \dots + x^{(n)}) - n\,\mu]$$

$$\frac{\partial L(\mu, \sigma \mid x^{(1)}, x^{(2)}, \dots, x^{(n)})}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3}[(x^{(1)} - \mu)^2 + \dots + (x^{(n)} - \mu)^2]$$

# Likelihood function of the Normal Distribution: Example (cont.)

$$\frac{\partial L(\mu,\sigma \mid x^{(1)},x^{(2)},\ldots,x^{(n)})}{\partial \mu} = \frac{1}{\sigma^2}[(x^{(1)}+\ldots+x^{(n)})-n\,\mu] = 0$$

Multiply both sides by $\sigma^2$

$$\mu = \frac{(x^{(1)}+\ldots+x^{(n)})}{n}$$

$$\frac{\partial L(\mu,\sigma \mid x^{(1)},x^{(2)},\ldots,x^{(n)})}{\partial \sigma} = -\frac{n}{\sigma} + \frac{1}{\sigma^3}[\left(x^{(1)}-\mu\right)^2 + \cdots + \left(x^{(n)}-\mu\right)^2] = 0$$
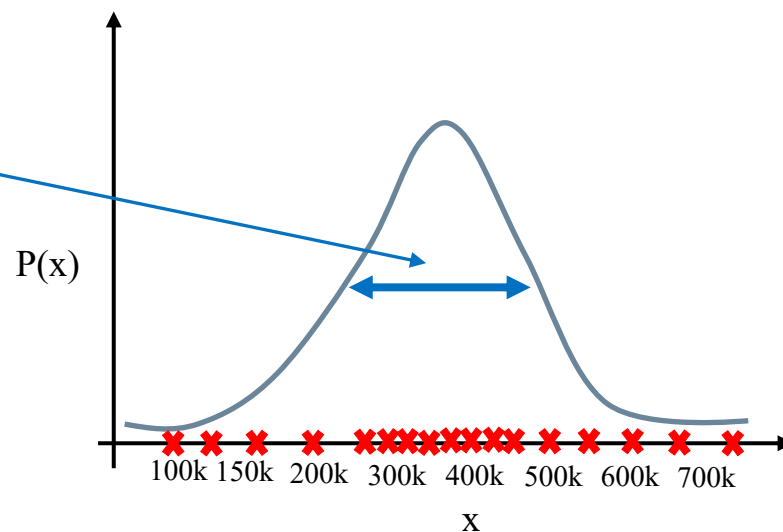
$$\sigma = \sqrt{\frac{(x^{(1)}-\mu)^2 + \cdots + (x^{(n)}-\mu)^2}{n}}$$

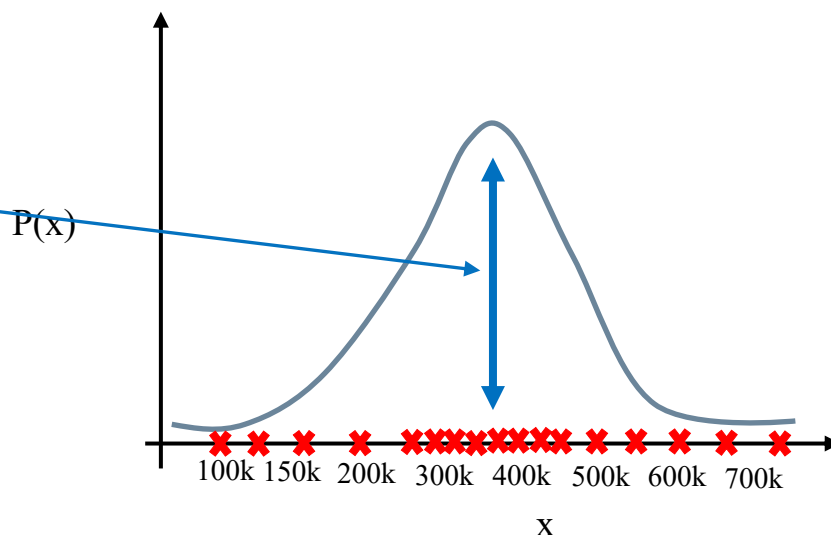# Likelihood function of the Normal Distribution: Example (cont.)

$$\mu = \frac{(x^{(1)} + \ldots + x^{(n)})}{n}$$

The mean of the data is the maximum likelihood estimate for where the center of the normal distribution

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

$$\sigma = \sqrt{\frac{(x^{(1)} - \mu)^2 + \cdots + (x^{(n)} - \mu)^2}{n}}$$

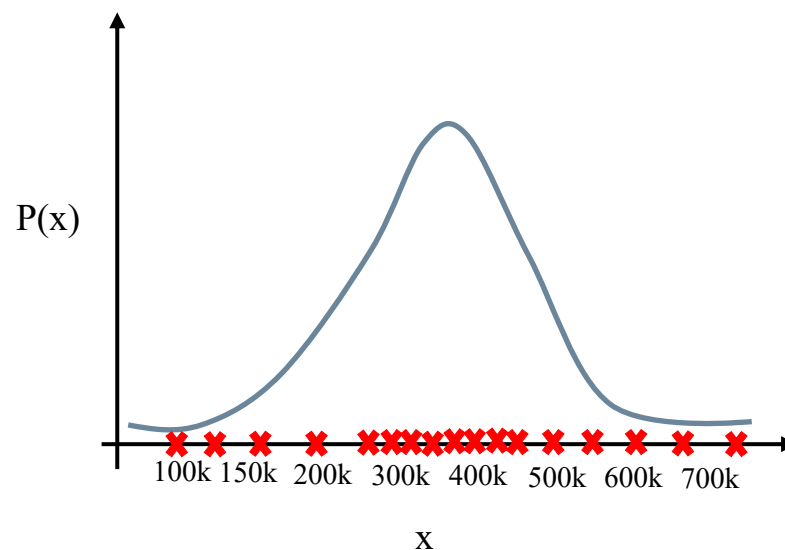The standard deviation of the data is the maximum likelihood estimate how wide the normal distribution shoul be

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

# Likelihood function of the Normal Distribution: Example (cont.)

$$\mu = \frac{(x^{(1)} + \dots + x^{(n)})}{n}$$

$$\sigma = \sqrt{\frac{(x^{(1)} - \mu)^2 + \dots + (x^{(n)} - \mu)^2}{n}}$$

P(x)

100k 150k 200k 300k 400k 500k 600k 700k

x

These solutions may be obvious, but
from maximum likelihood estimation,
we prove that our intuition are correct.

- Thank you!