

# A Primer to Probability distribution and Parameter Estimation



# Mean and Variance of Random Variable



- Expected value of a random variable  $X$  is the long-run average value of repetitions of the experiment

$$E[X] = x_1p_1 + x_2p_2 + \cdots + x_kp_k .$$

- Variance : Spread of the random variable values

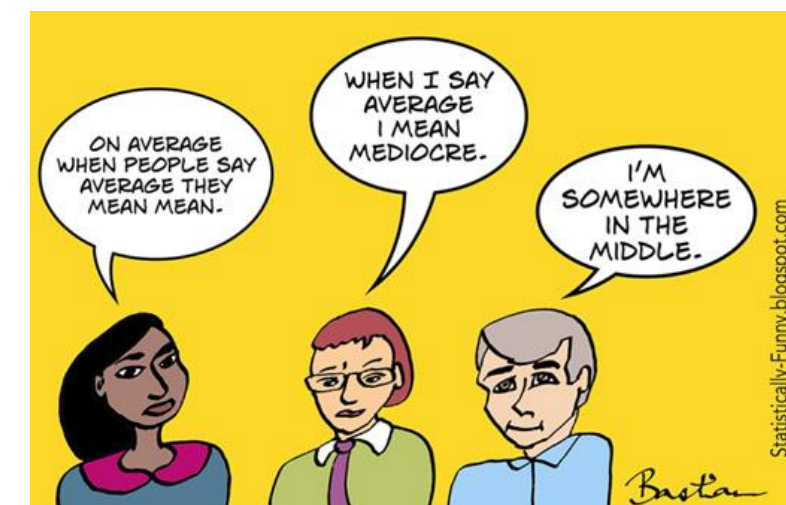
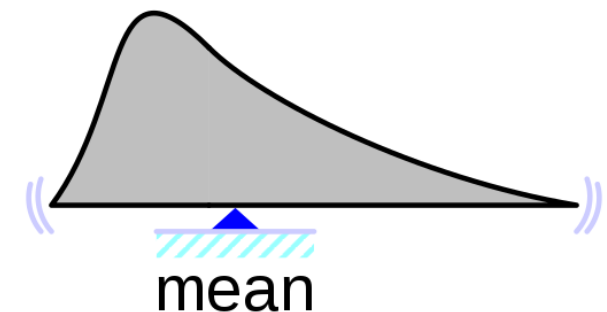
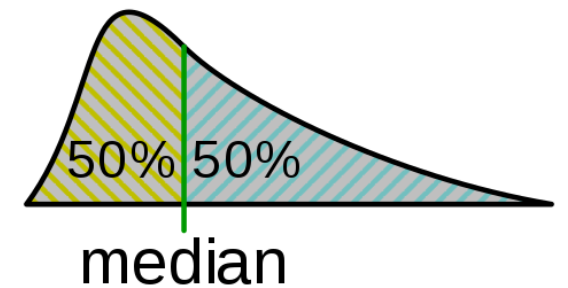
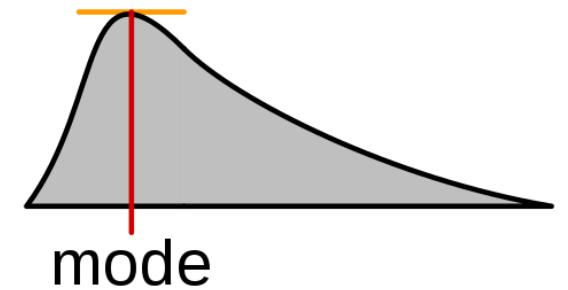
$$\text{Var}(X) = E[(X - \mu)^2] = E[X^2] - (E[X])^2$$

$\sqrt{\text{Var}(X)}$  is called the *standard deviation* of  $X$ .

$W = 0$  with probability 1

$$Y = \begin{cases} -1 & \text{with probability } \frac{1}{2} \\ 1 & \text{with probability } \frac{1}{2} \end{cases}$$

$$Z = \begin{cases} -100 & \text{with probability } \frac{1}{2} \\ 100 & \text{with probability } \frac{1}{2} \end{cases}$$



# Common Discrete Distributions

## Bernoulli and Binomial



- Let  $X \in \{0, 1\}$  be a binary random variable, with probability of “success”  $\theta$ ,  $X$  has a Bernoulli distribution,  $X \sim \text{Ber}(\theta)$

E.g Coin toss, Rain or not

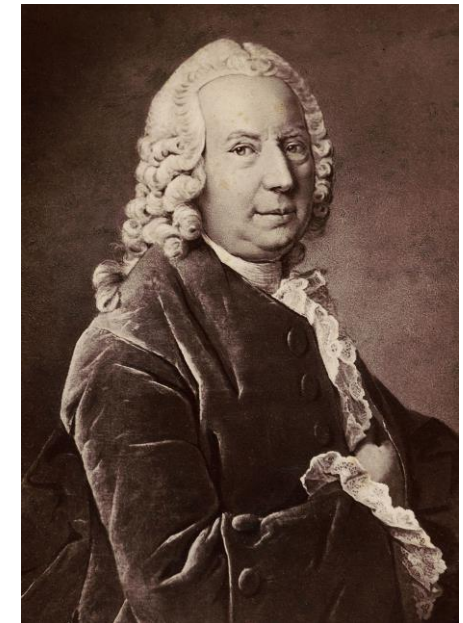
$$\text{Ber}(x|\theta) = \theta^{I(x=1)}(1 - \theta)^{I(x=0)}$$

$$\text{Ber}(x|\theta) = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases}$$



# Common Discrete Distributions

## Bernoulli and Binomial



- Let  $X \in \{0, 1\}$  be a binary random variable, with probability of “success”  $\theta$ ,  $X$  has a Bernoulli distribution,  $X \sim \text{Ber}(\theta)$

E.g Coin toss, Rain or not

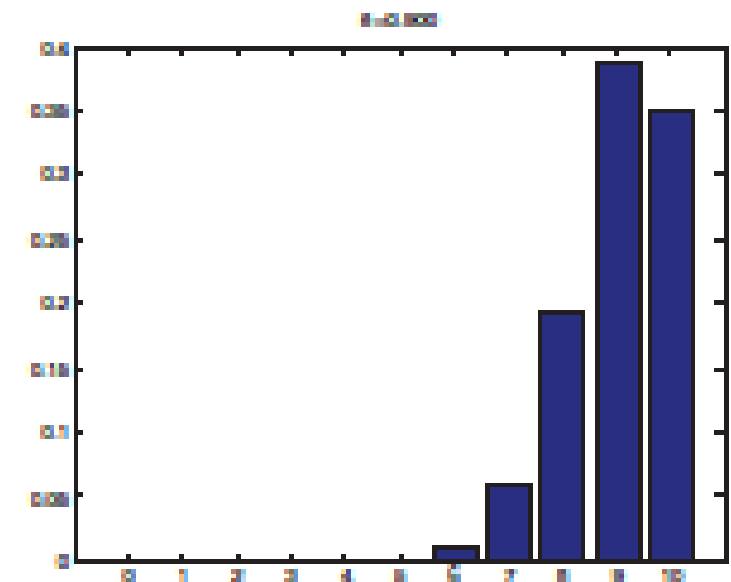
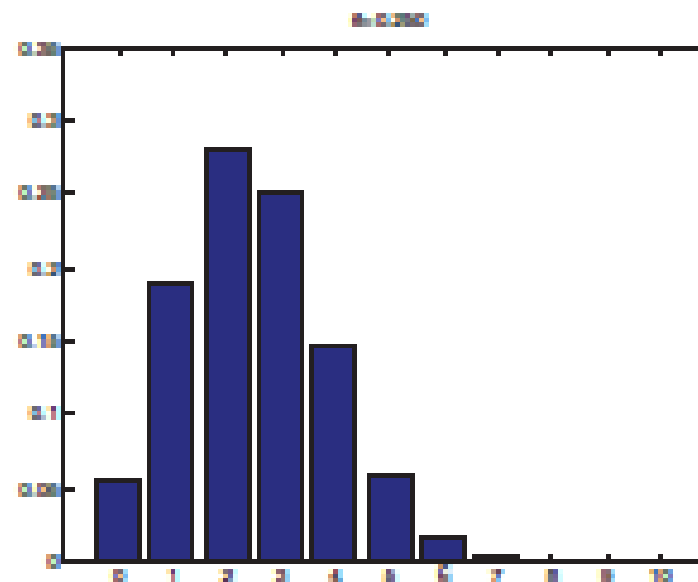
$$\text{Ber}(x|\theta) = \theta^{\mathbb{I}(x=1)}(1 - \theta)^{\mathbb{I}(x=0)} \quad \text{Ber}(x|\theta) = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases}$$

- Suppose we toss a coin  $n$  times. Let  $X \in \{0, \dots, n\}$  be the number of heads. If the probability of heads is  $\theta$ , then we say  $X$  has a binomial distribution, written as  $X \sim \text{Bin}(n, \theta)$ .

$$\text{Bin}(k|n, \theta) \triangleq \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

$$\text{mean} = n\theta, \quad \text{var} = n\theta(1 - \theta)$$

$$\binom{n}{k} \triangleq \frac{n!}{(n-k)!k!}$$





# Discrete Distributions : Multinoulli, Multinomial

- Model the outcomes of tossing a K -sided die :  
categorical/Multinoulli distribution,

$$x \sim \text{Cat}(\theta), \quad p(x = j | \theta) = \theta_j.$$

- Multinomial distribution : Models the outcome of n dice rolls, let  $x = (x_1, \dots, x_K)$  be a random vector, where  $x_j$  number of times side j of the die occurs.

$$\text{Mu}(x | n, \theta) \triangleq \binom{n}{x_1 \dots x_K} \prod_{j=1}^K \theta_j^{x_j}$$

$$\text{Cat}(x | \theta) \triangleq \text{Mu}(x | 1, \theta) \quad \text{Mu}(x | 1, \theta) = \prod_{j=1}^K \theta_j^{x_j - 1}$$

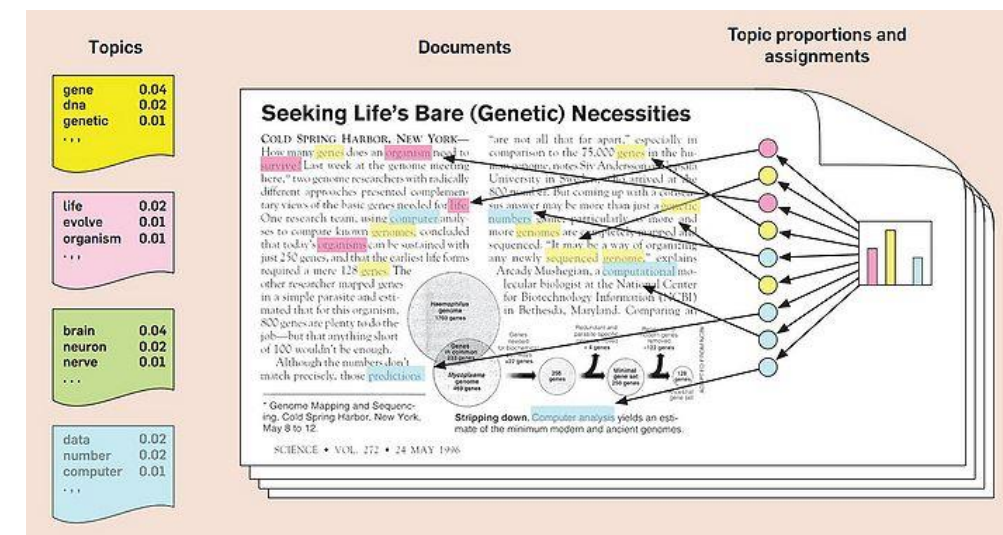
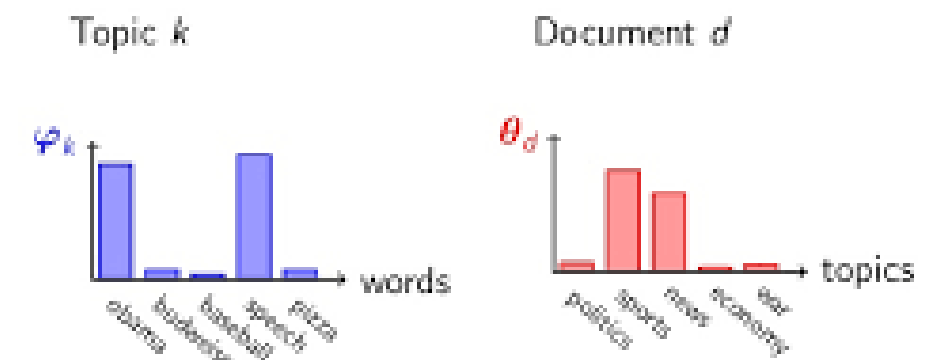
- Probabilistic topic model

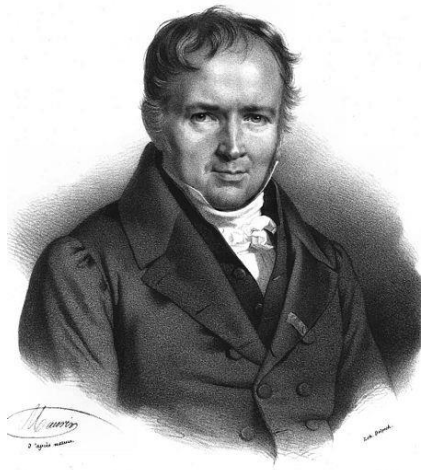
- Text classification

## Latent Dirichlet Allocation

LDA discovers topics into  
a collection of documents.

LDA tags each document  
with topics.





# Poisson distribution



- Model number of events occurring in a fixed interval of time/space

$$P(k \text{ events in interval}) = e^{-\lambda} \frac{\lambda^k}{k!}$$

- $\lambda$  is the average (mean) number of events per interval,  $k = 0, 1, 2, \dots$ , events occur independently, rate is a constant.

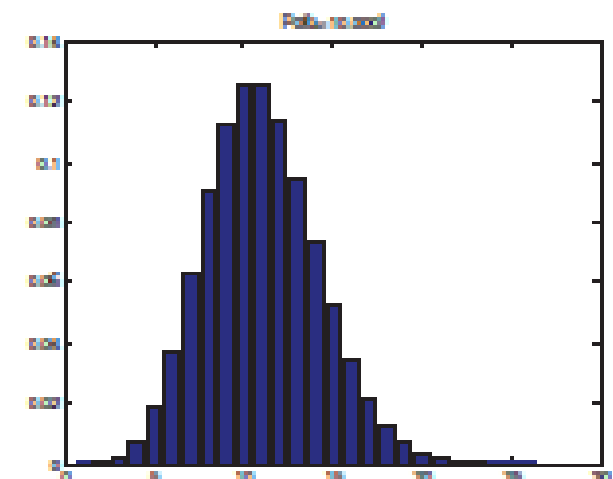
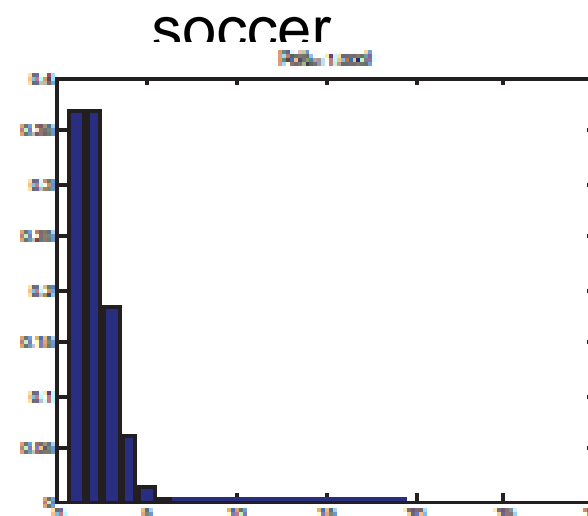
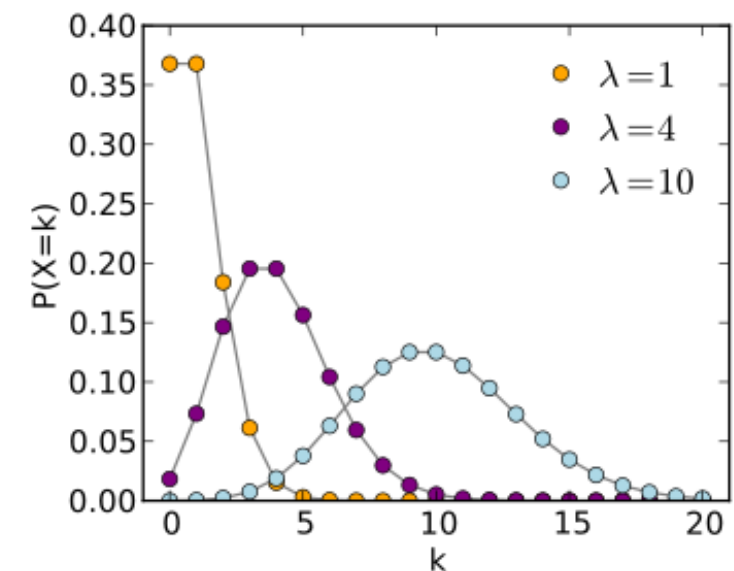
- Models rare events

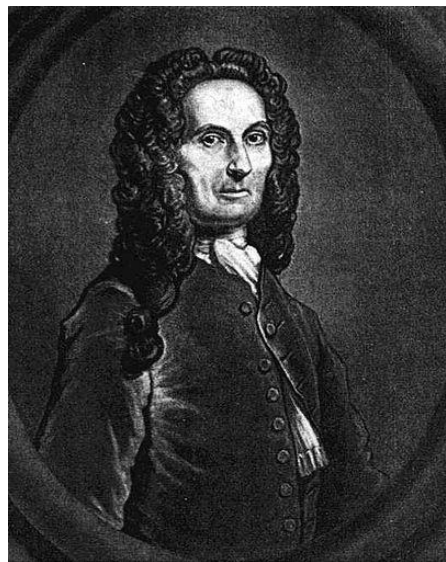
- Number of misprints on a page of a book.

- average number of goals in a World Cup match is approximately 2.5 ;  $\lambda = 2.5$ .

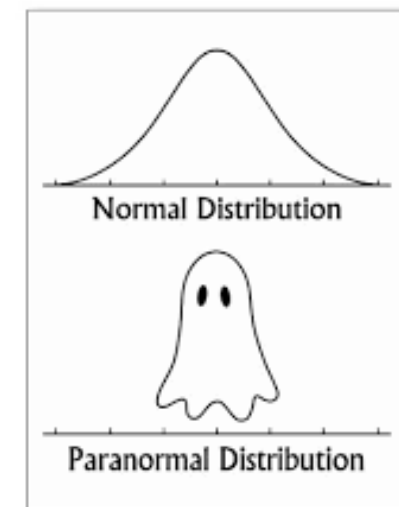
$$P(k \text{ goals in a match}) = \frac{2.5^k e^{-2.5}}{k!}$$

- Number of wrong telephone numbers that are dialed in a day.





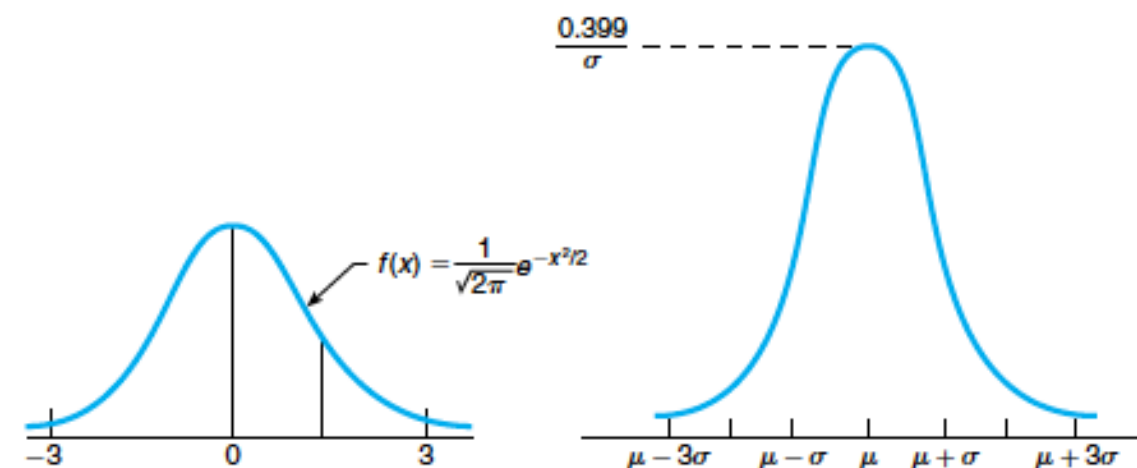
# Normal/Gaussian Random Variables



- 1809 Gauss published his monograph "Theoria motus corporum coelestium in sectionibus conicis solem ambientium"
- All distributions of frequency other than normal are 'abnormal'- Pearson
- A random variable is said to be normally distributed

with parameters  $\mu$  and  $\sigma^2$ ,  $X \sim N(\mu, \sigma^2)$

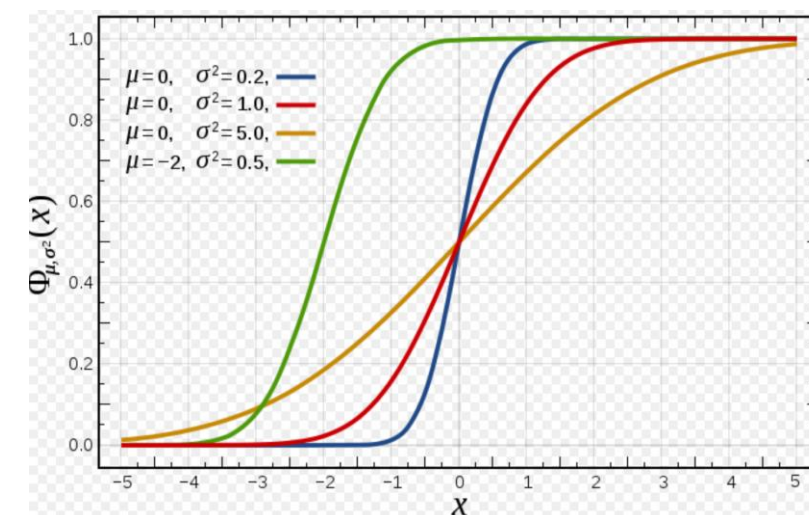
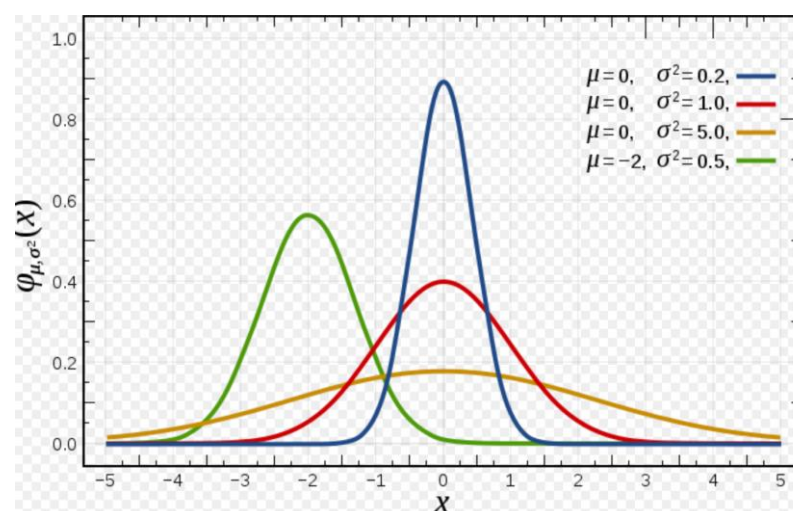
$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}, \quad -\infty < x < \infty$$



- $\mu = E[X]$  is the mean (and mode), and  $\sigma^2 = \text{var}[X]$  is the variance.

$$\Phi(x; \mu, \sigma^2) \triangleq \int_{-\infty}^x \mathcal{N}(z | \mu, \sigma^2) dz$$

- CDF of the Gaussian



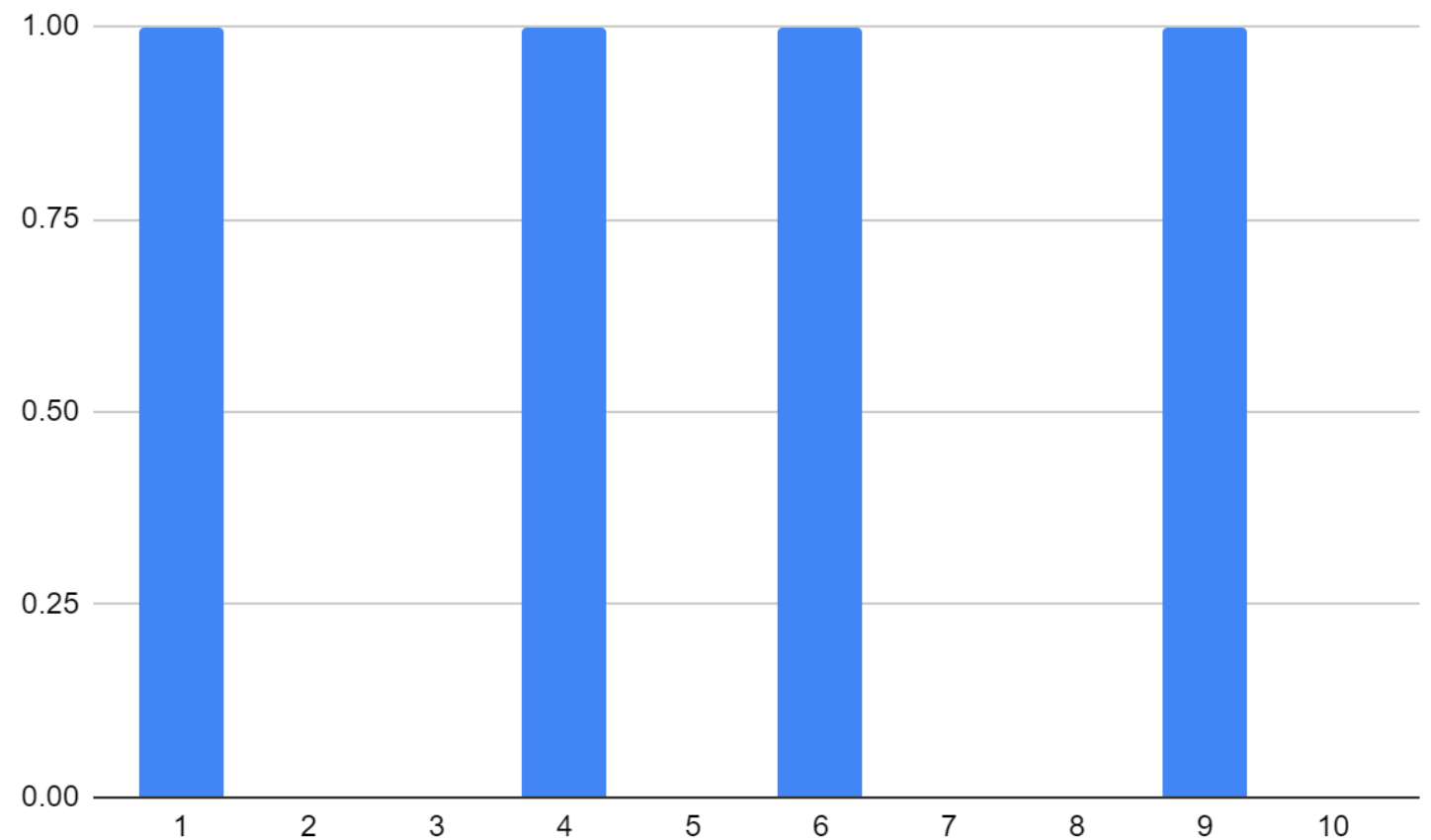
# Probability Distribution Summary

- X : Discrete
  - Binary valued scalar (0/1) : Bernoulli
  - Binary valued vector (one of K): Multinoulli/categorical
  - Multivalued scalar (M of N ) : Binomial
  - Multivalued vector (M1, M2, ... MK) : Multinomial
  - Integer valued scalar (1 to infinity) : Poisson
- X : continuous, real valued
  - Interval [a,b] : Uniform, Interval [0,1] : Beta
  - non-negative (0,infinity) : Exponential, Gamma
  - real line (-infinity, infinity) : Normal, students, Laplace
  - Vector : Real valued : Gaussian ; Simplex : Dirichlet



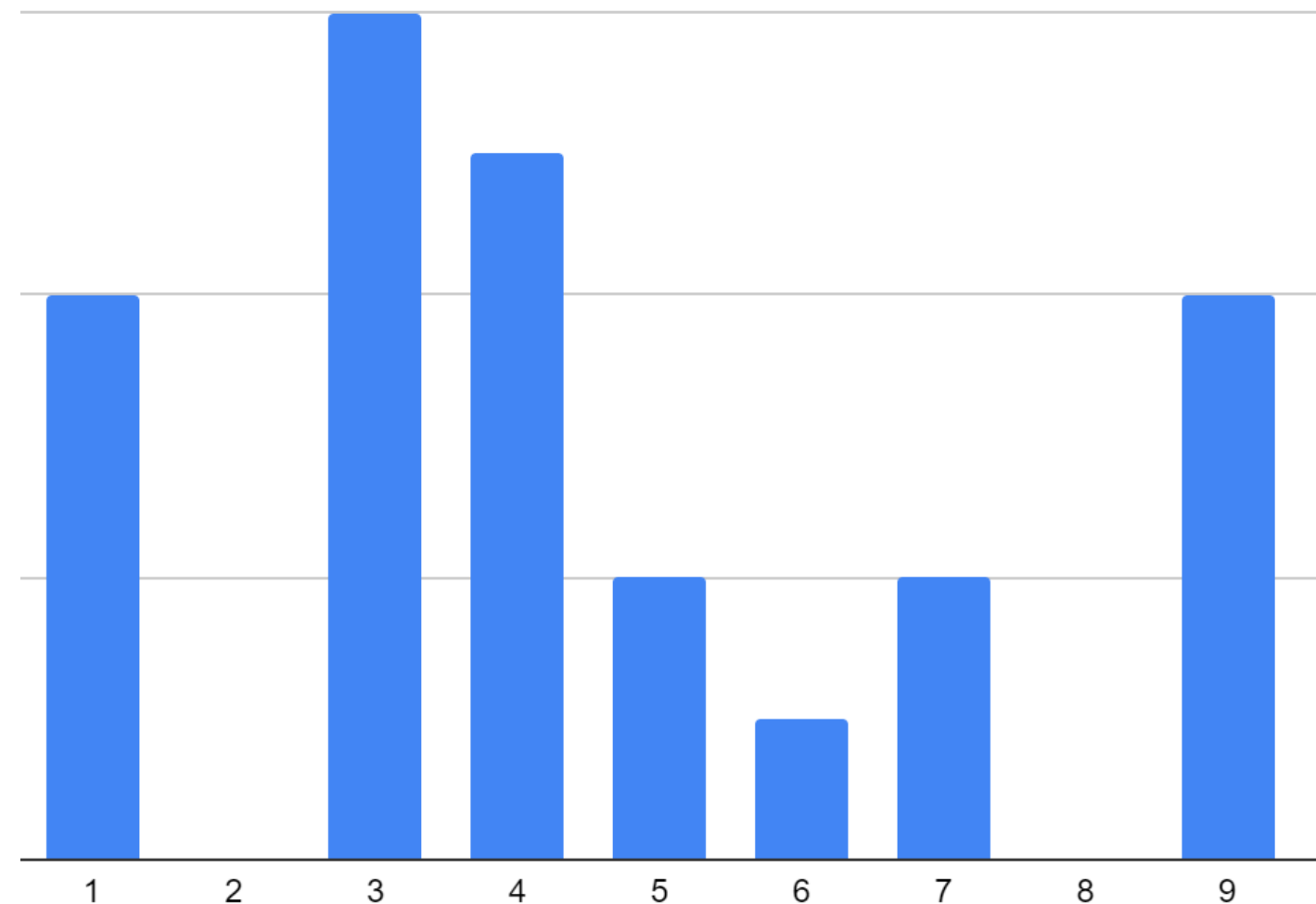
# Parameter Estimation

- Data points representing the if it has rained or not in last 10 days.
- Whats the probability that it will rain tomorrow ?
- How many days will it rain in next 5 days ?



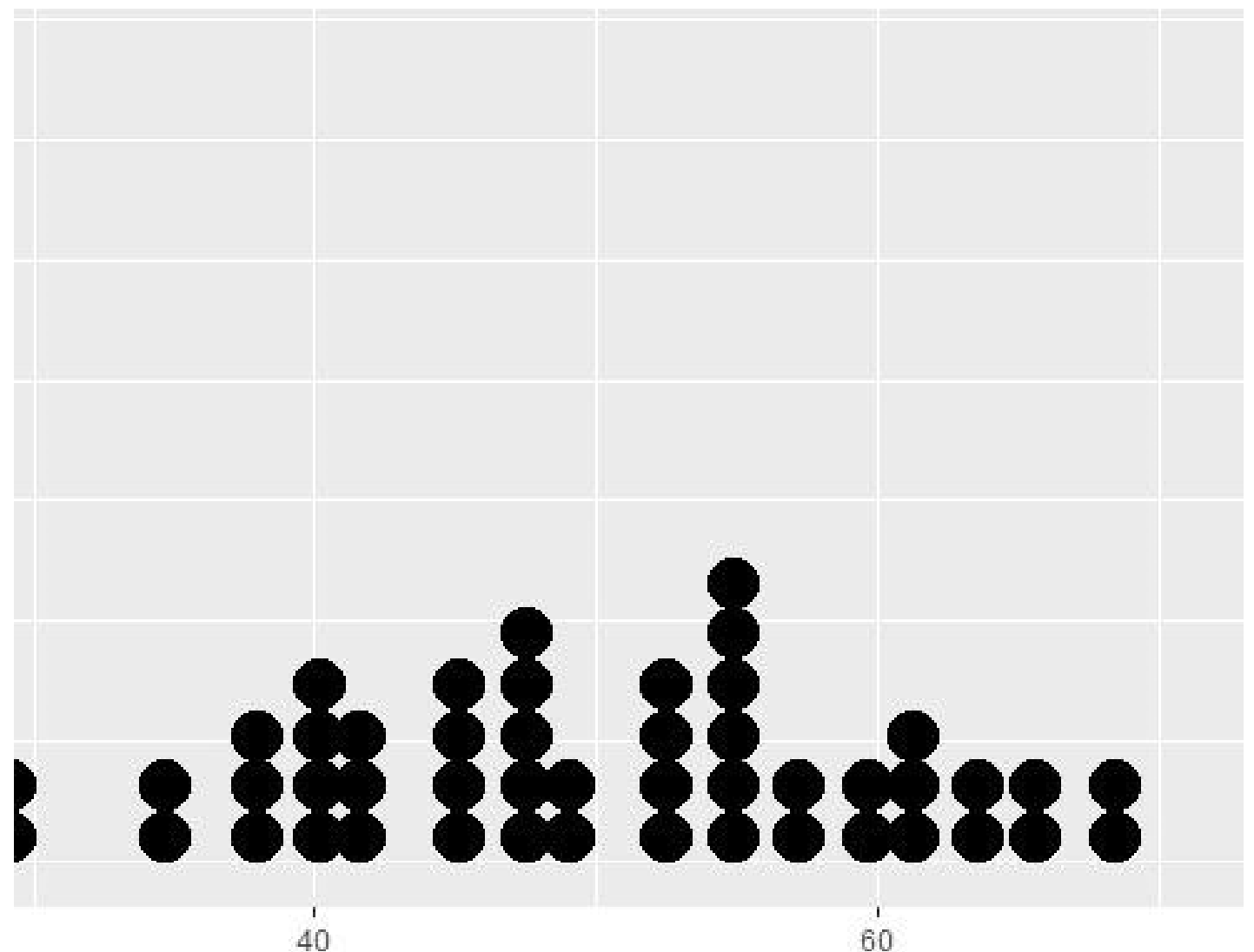
# Parameter Estimation

- The number of traffic accidents in Berkeley, California, in 10 randomly chosen nonrainy days in 1998 is as follows:
- 4, 0, 6, 5, 2, 1, 2, 0, 4, 3
- Use these data to estimate the proportion of nonrainy days that had 2 or fewer accidents that year.



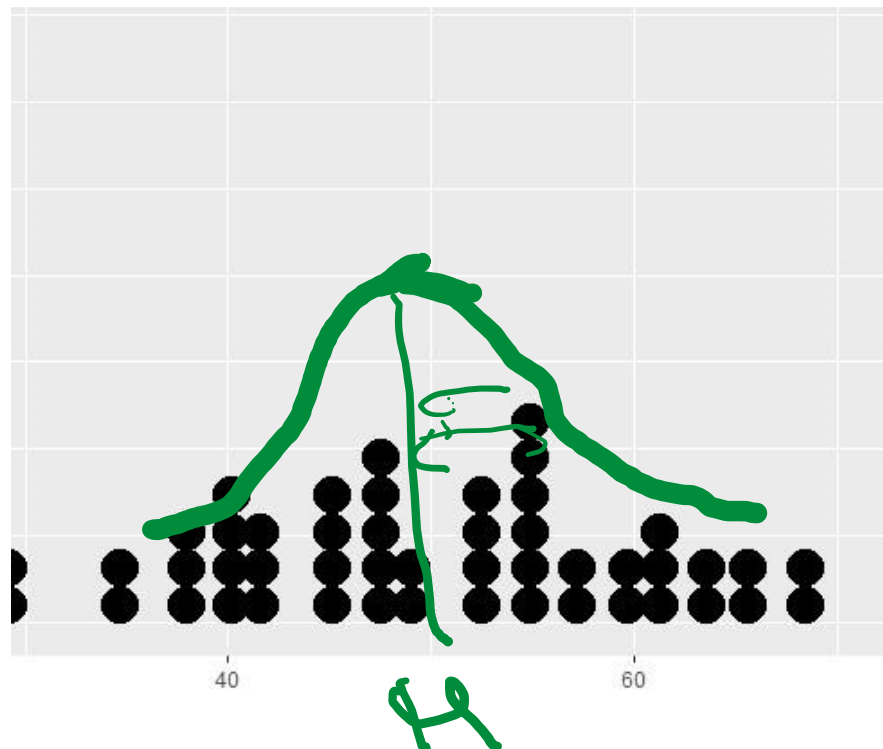
# Parameter Estimation

- Data points representing the weight (in kgs) of students in a class.
- Whats mean and std deviation of the data ?
- Whats the probability that weight  $> 60$



# Parameter Estimation

- Any statistic used to estimate the value of an unknown parameter  $\theta$  is called an estimator of  $\theta$ .
  - mean and variance for Normal, rate ( $\lambda$ ) for Poisson, etc.
- Maximum likelihood** estimator
- MLE can be defined as a method for estimating parameters of a distribution from sample data such that the likelihood of obtaining the observed data is maximized.
- Provides optimal way to fit a distribution to the data

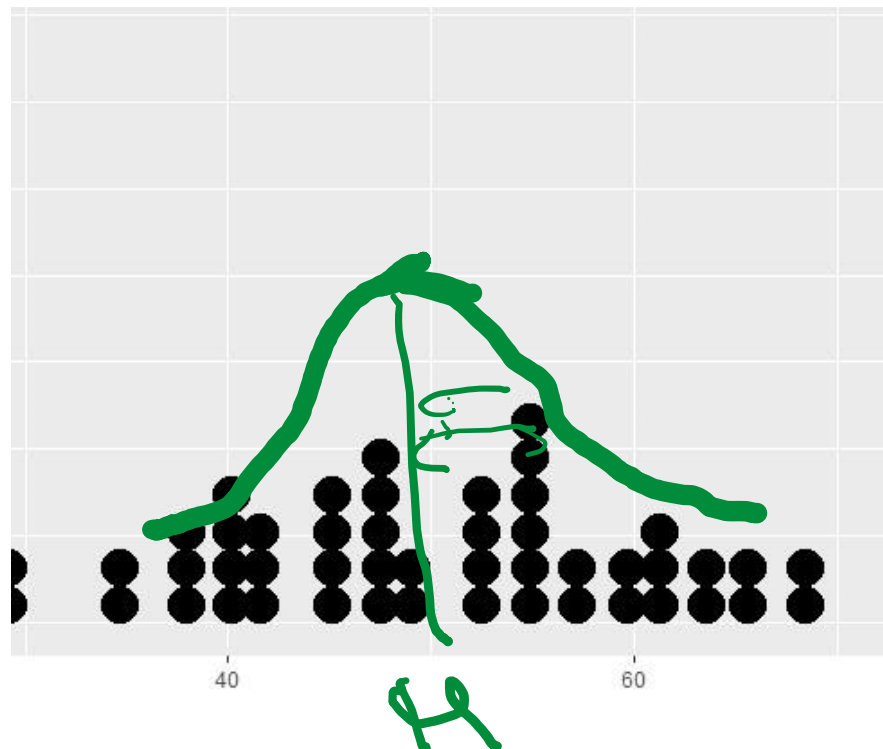


$$X \sim N(\mu, \sigma^2)$$
$$\theta = (\mu, \sigma)$$



# Parameter Estimation

- Any statistic used to estimate the value of an unknown parameter  $\theta$  is called an estimator of  $\theta$ .
  - mean and variance for Normal, rate ( $\lambda$ ) for Poisson, etc.
- Maximum likelihood estimator**
- MLE can be defined as a method for estimating parameters of a distribution from sample data such that the likelihood of obtaining the observed data is maximized.
- Provides optimal way to fit a distribution to the data



$$X \sim N(\mu, \sigma^2)$$
$$\theta = (\mu, \sigma)$$

# Parameter Estimation

## Maximum likelihood estimator

- $f(x_1, \dots, x_n|\theta)$  represents the probability that the values  $x_1, x_2, \dots, x_n$  will be observed when  $\theta$  is the true value of the parameter
- **Maximum Likelihood estimation** : maximum likelihood estimate  $\hat{\theta}$  is defined to be that value of  $\theta$  maximizing  $L(\theta) = f(x_1, \dots, x_n|\theta)$

$$\operatorname{argmax}_{\theta} L(\theta) = \operatorname{argmax}_{\theta} f(x_1, \dots, x_n|\theta) = \operatorname{argmax}_{\theta} \log[f(x_1, \dots, x_n|\theta)].$$

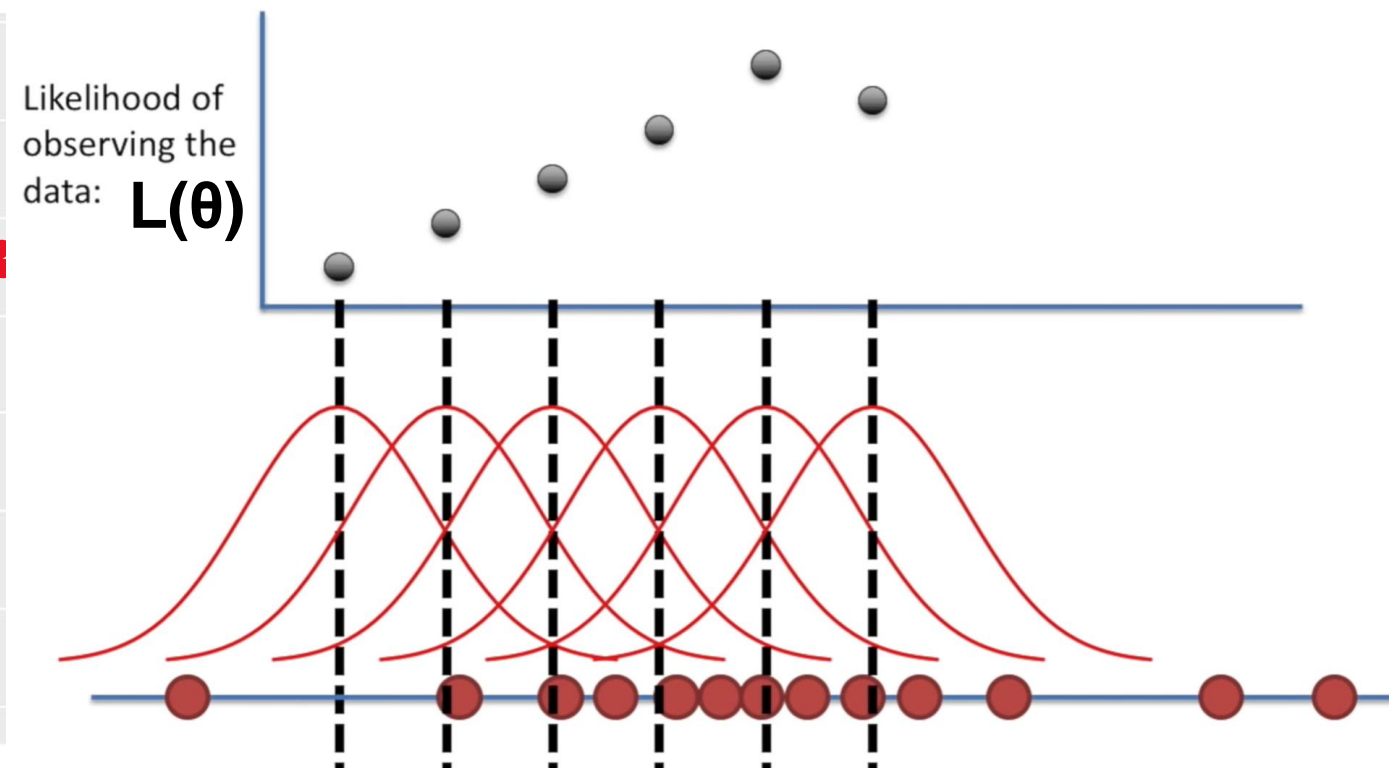
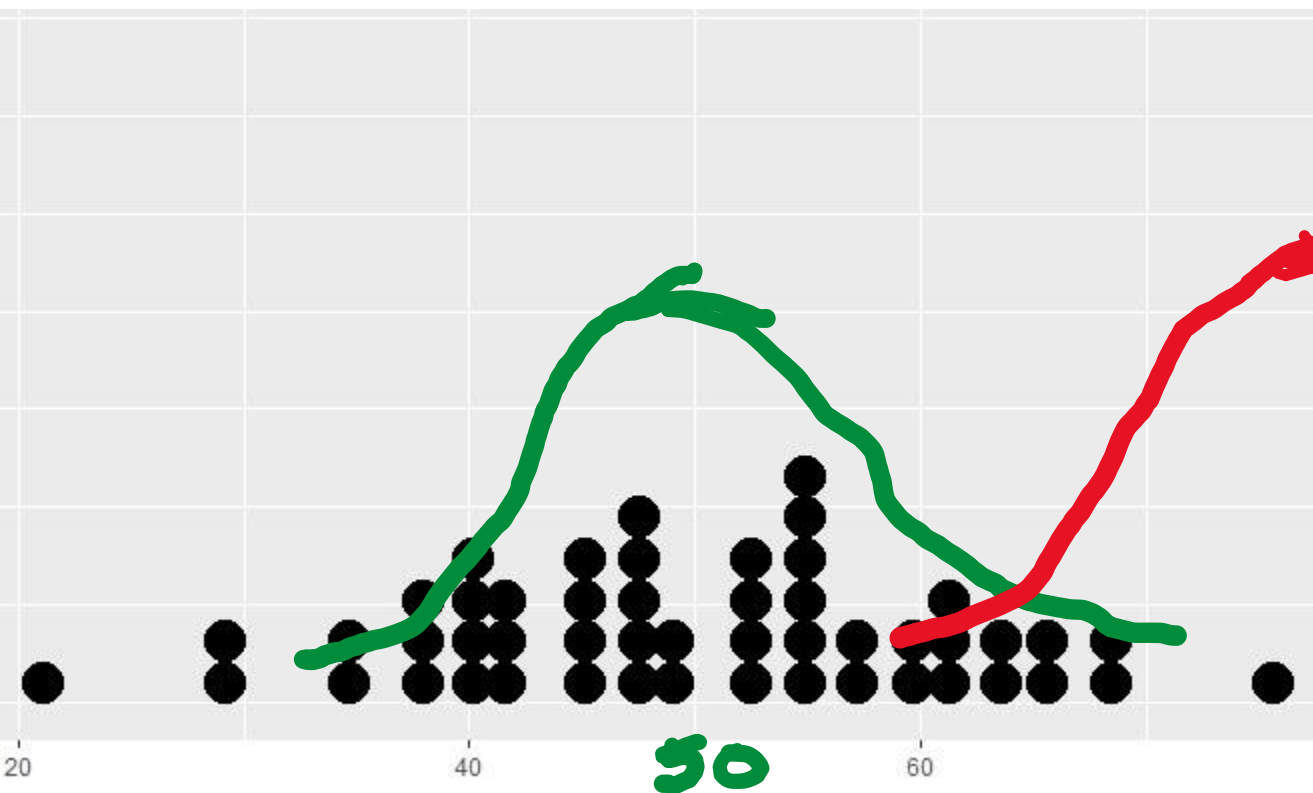
Note that  $L(\theta)$  is not a distribution over  $\theta$  but just a function of  $\theta$ .

Independent and identically distributed (i.i.d.) assumption

$$f(x_1, x_2, \dots, x_n) = f_{X_1}(x_1)f_{X_2}(x_2) \cdots f_{X_n}(x_n)$$

# Parameter estimation

- which of the following would maximize the probability of observing the data
  - Mean = 100, SD = 10
  - Mean = 50, SD = 10



# Maximum likelihood estimation

- (Maximum Likelihood Estimator of a Bernoulli Parameter) Suppose you have data from  $n$  independent Bernoulli trials,  $X_1, \dots, X_n$ . Assuming the success probability is  $p$  what is the maximum likelihood estimator of  $p$ ?

$$X_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \text{otherwise} \end{cases}$$

$$P\{X_i = 1\} = p = 1 - P\{X_i = 0\}$$

$$P\{X_i = x\} = p^x (1 - p)^{1-x}, \quad x = 0, 1$$





# Maximum likelihood estimation

- (Maximum Likelihood Estimator of a Bernoulli Parameter) Suppose you have data from  $n$  independent Bernoulli trials,  $X_1, \dots, X_n$ . Assuming the success probability is  $p$  what is the maximum likelihood estimator of  $p$ ?

$$X_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \text{otherwise} \end{cases} \quad \begin{aligned} P\{X_i = 1\} &= p = 1 - P\{X_i = 0\} \\ P\{X_i = x\} &= p^x(1 - p)^{1-x}, \quad x = 0, 1 \end{aligned}$$

$$\begin{aligned} f(x_1, \dots, x_n | p) &= P\{X_1 = x_1, \dots, X_n = x_n | p\} \\ &= p^{x_1} (1 - p)^{1-x_1} \dots p^{x_n} (1 - p)^{1-x_n} \\ &= p^{\sum_{i=1}^n x_i} (1 - p)^{n - \sum_{i=1}^n x_i}, \quad x_i = 0, 1, \quad i = 1, \dots, n \end{aligned}$$

# Maximum likelihood estimation

- (Maximum Likelihood Estimator of a Bernoulli Parameter) Suppose you have data from  $n$  independent Bernoulli trials,  $X_1, \dots, X_n$ . Assuming the success probability is  $p$  what is the maximum likelihood estimator of  $p$ ?

$$X_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \text{otherwise} \end{cases} \quad P\{X_i = 1\} = p = 1 - P\{X_i = 0\}$$

To determine the value of  $p$  that maximizes the likelihood,

$$\log f(x_1, \dots, x_n | p) = \sum_{i=1}^n x_i \log p + \left( n - \sum_{i=1}^n x_i \right) \log(1 - p)$$

# Maximum likelihood estimation

- (Maximum Likelihood Estimator of a Bernoulli Parameter) Suppose you have data from  $n$  independent Bernoulli trials,  $X_1, \dots, X_n$ . Assuming the success probability is  $p$  what is the maximum likelihood estimator of  $p$ ?

$$X_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \text{otherwise} \end{cases} \quad P\{X_i = 1\} = p = 1 - P\{X_i = 0\}$$

To determine the value of  $p$  that maximizes the likelihood,

$$\log f(x_1, \dots, x_n | p) = \sum_{i=1}^n x_i \log p + \left( n - \sum_{i=1}^n x_i \right) \log(1 - p)$$

$$\frac{d}{dp} \log f(x_1, \dots, x_n | p) = \frac{\sum_{i=1}^n x_i}{p} - \frac{\left( n - \sum_{i=1}^n x_i \right)}{1 - p} \quad \hat{p} = \frac{\sum_{i=1}^n x_i}{n}$$

# Maximum likelihood estimation

- (Maximum Likelihood Estimator of a Bernoulli Parameter) Suppose you have data from  $n$  independent Bernoulli trials,  $X_1, \dots, X_n$ . Assuming the success probability is  $p$  what is the maximum likelihood estimator of  $p$ ?

$$X_i = \begin{cases} 1 & \text{if trial } i \text{ is a success} \\ 0 & \text{otherwise} \end{cases} \quad P\{X_i = 1\} = p = 1 - P\{X_i = 0\}$$

To determine the value of  $p$  that maximizes the likelihood,

proportion of the observed trials that result in successes.

$$\frac{d}{dp} \log f(x_1, \dots, x_n | p) = \frac{\sum_{i=1}^n x_i}{p} - \frac{\left(n - \sum_{i=1}^n x_i\right)}{1 - p} \quad \hat{p} = \frac{\sum_{i=1}^n x_i}{n}$$

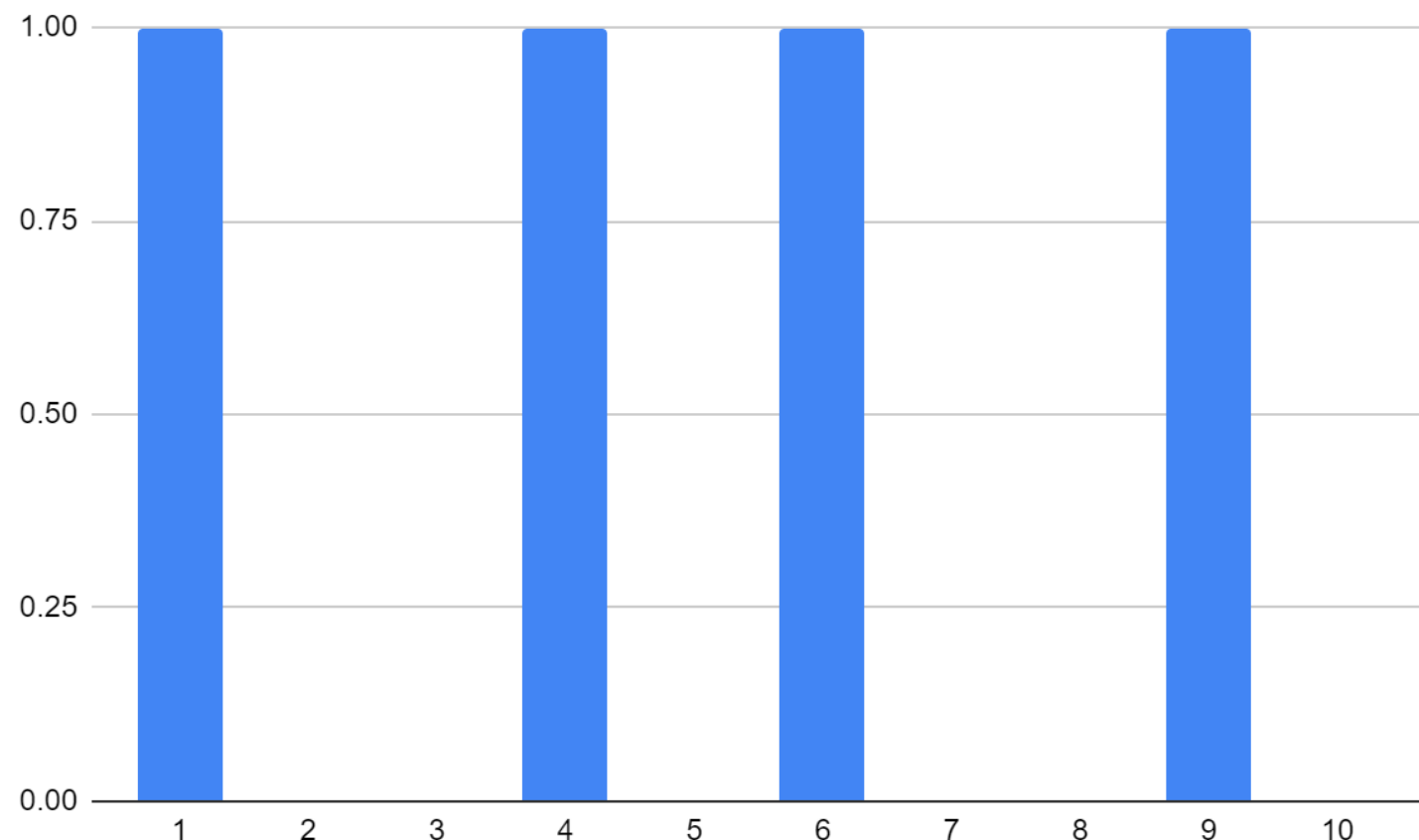




Suppose that each RAM (random access memory) chip produced by a certain manufacturer is, independently, of acceptable quality with probability  $p$ . Then if out of a sample of 1,000 tested 921 are acceptable, what is the maximum likelihood estimate of  $p$  ?



- Data points representing the if it has rained or not in last 10 days.
- Whats the probability that it will rain tomorrow ?
- How many days will it rain in next 5 days ?



# Parameter Estimation

- Multinomial
- 3,1,2,4,3,5,6,1,3,4



$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad \mathbf{x} = (0, 0, 1, 0, 0, 0)^T.$$

data set  $\mathcal{D}$  of  $N$  independent observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}. \quad m_k = \sum_n x_{nk}$$

# Parameter Estimation

- Multinomial

data set  $\mathcal{D}$  of  $N$  independent observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}.$$

# Parameter Estimation : Multinomial

data set  $\mathcal{D}$  of  $N$  independent observations  $\mathbf{x}_1, \dots, \mathbf{x}_N$ .

$$p(\mathcal{D}|\boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{(\sum_n x_{nk})} = \prod_{k=1}^K \mu_k^{m_k}.$$

$$p(\mathbf{x}|\boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k}$$

Constraint

$$\mu_k^{\text{ML}} = \frac{m_k}{N}$$

ML solution !





# MLE : Poisson !

- (Maximum Likelihood Estimator of a Poisson Parameter) Suppose  $X_1, \dots, X_n$  are independent Poisson random variables each having mean  $\lambda$ . Determine the maximum likelihood estimator of  $\lambda$ .

$$\begin{aligned} f(x_1, \dots, x_n | \lambda) &= \frac{e^{-\lambda} \lambda^{x_1}}{x_1!} \dots \frac{e^{-\lambda} \lambda^{x_n}}{x_n!} \\ &= \frac{e^{-n\lambda} \lambda^{\sum_{i=1}^n x_i}}{x_1! \dots x_n!} \end{aligned}$$

$$\log f(x_1, \dots, x_n | \lambda) = -n\lambda + \sum_{i=1}^n x_i \log \lambda - \log c$$

# MLE : Poisson !

- (Maximum Likelihood Estimator of a Poisson Parameter) Suppose  $X_1, \dots, X_n$  are independent Poisson random variables each having mean  $\lambda$ . Determine the maximum likelihood estimator of  $\lambda$ .

$$\frac{d}{d\lambda} \log f(x_1, \dots, x_n | \lambda) = -n + \frac{\sum_{i=1}^n x_i}{\lambda}$$

$$\hat{\lambda} = \frac{\sum_{i=1}^n x_i}{n}$$

ML solution !



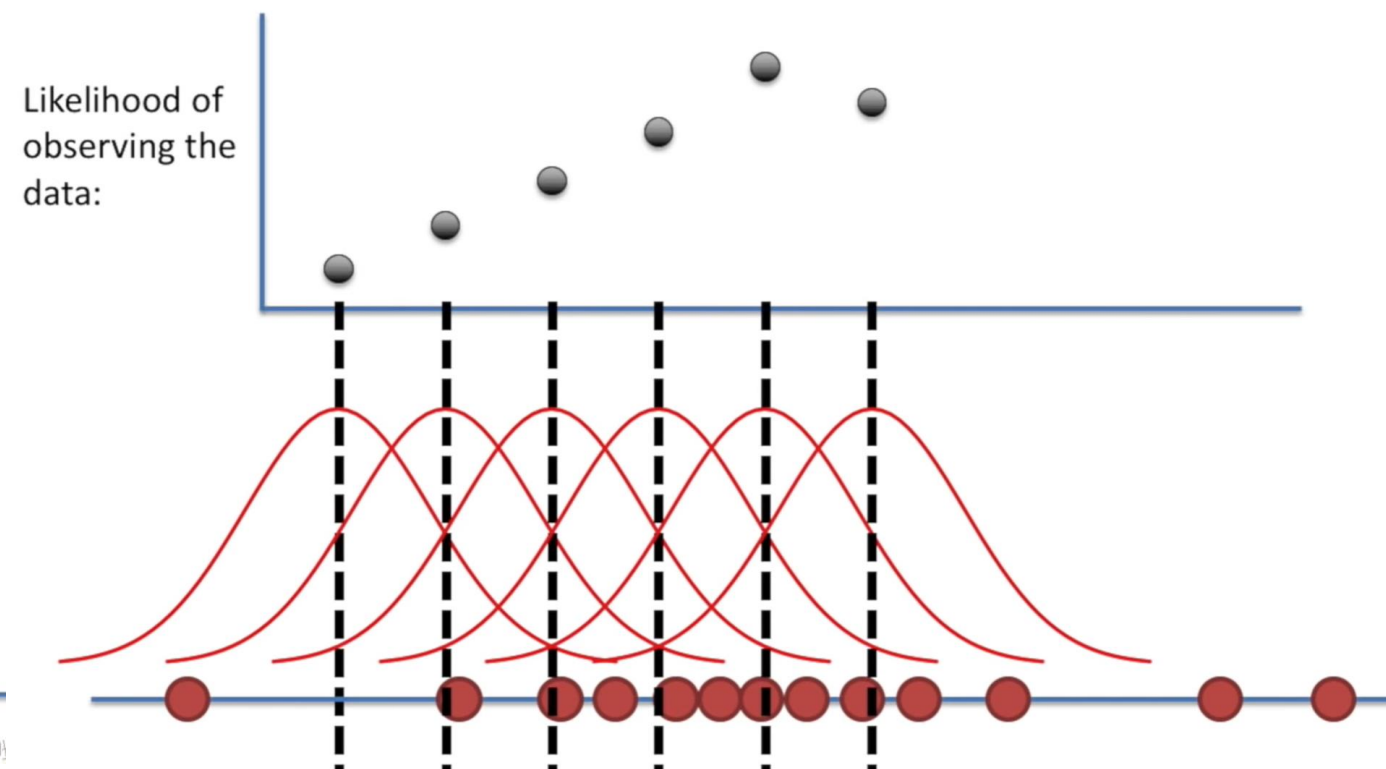
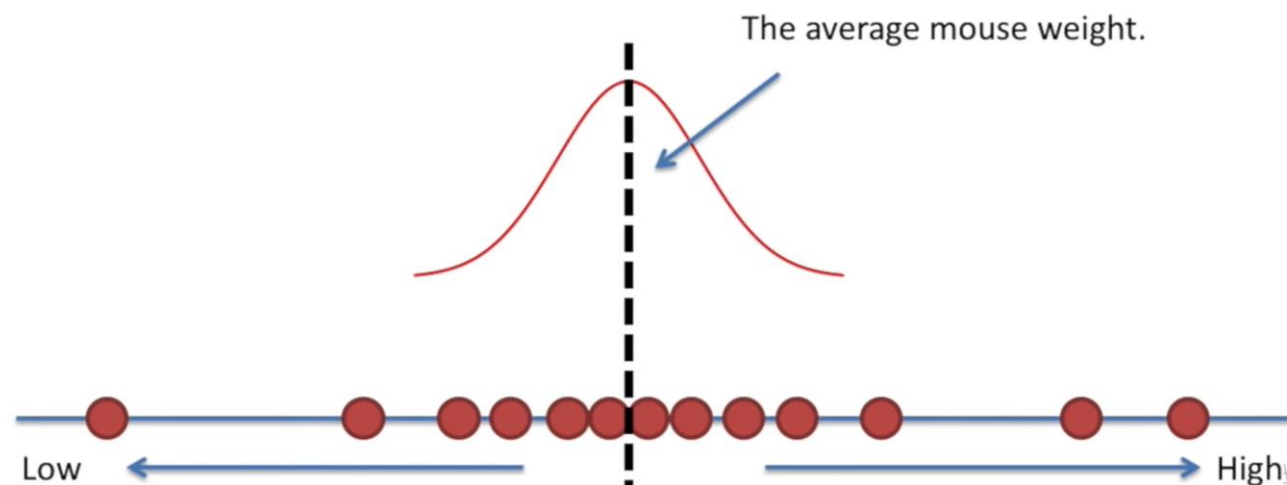
# Maximum likelihood estimation

- The number of traffic accidents in Berkeley, California, in 10 randomly chosen nonrainy days in 1998 is as follows:
- 4, 0, 6, 5, 2, 1, 2, 0, 4, 3
- Use these data to estimate the proportion of nonrainy days that had 2 or fewer accidents that year.

# MLE : Normal

- (Maximum Likelihood Estimator in a Normal Population) Suppose  $X_1, \dots, X_n$  are independent, normal random variables each with unknown mean  $\mu$  and unknown standard deviation  $\sigma$ .

- 1) We expect most of the measurements (mouse weights) to be close to the mean (average).



# MLE : Normal

- (Maximum Likelihood Estimator in a Normal Population) Suppose  $X_1, \dots, X_n$  are independent, normal random variables each with unknown mean  $\mu$  and unknown standard deviation  $\sigma$ .

$$\begin{aligned} f(x_1, \dots, x_n | \mu, \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left[\frac{-(x_i - \mu)^2}{2\sigma^2}\right] \\ &= \left(\frac{1}{2\pi}\right)^{n/2} \frac{1}{\sigma^n} \exp\left[\frac{-\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right] \end{aligned}$$

$$\log f(x_1, \dots, x_n | \mu, \sigma) = -\frac{n}{2} \log(2\pi) - n \log \sigma - \frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}$$



# MLE : Normal

- (Maximum Likelihood Estimator in a Normal Population) Suppose  $X_1, \dots, X_n$  are independent, normal random variables each with unknown mean  $\mu$  and unknown standard deviation  $\sigma$ .

$$\frac{\partial}{\partial \mu} \log f(x_1, \dots, x_n | \mu, \sigma) = \frac{\sum_{i=1}^n (x_i - \mu)}{\sigma^2}$$

$$\frac{\partial}{\partial \sigma} \log f(x_1, \dots, x_n | \mu, \sigma) = -\frac{n}{\sigma} + \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^3}$$

$$\hat{\mu} = \sum_{i=1}^n x_i / n \quad \hat{\sigma} = \left[ \sum_{i=1}^n (x_i - \hat{\mu})^2 / n \right]^{1/2}$$

# Model Selection

- Given some observations  $X_1, X_2, \dots, X_N$ , how do you decide which probability distribution to model it ?

[illegible]