

# Introduction to Econometrics

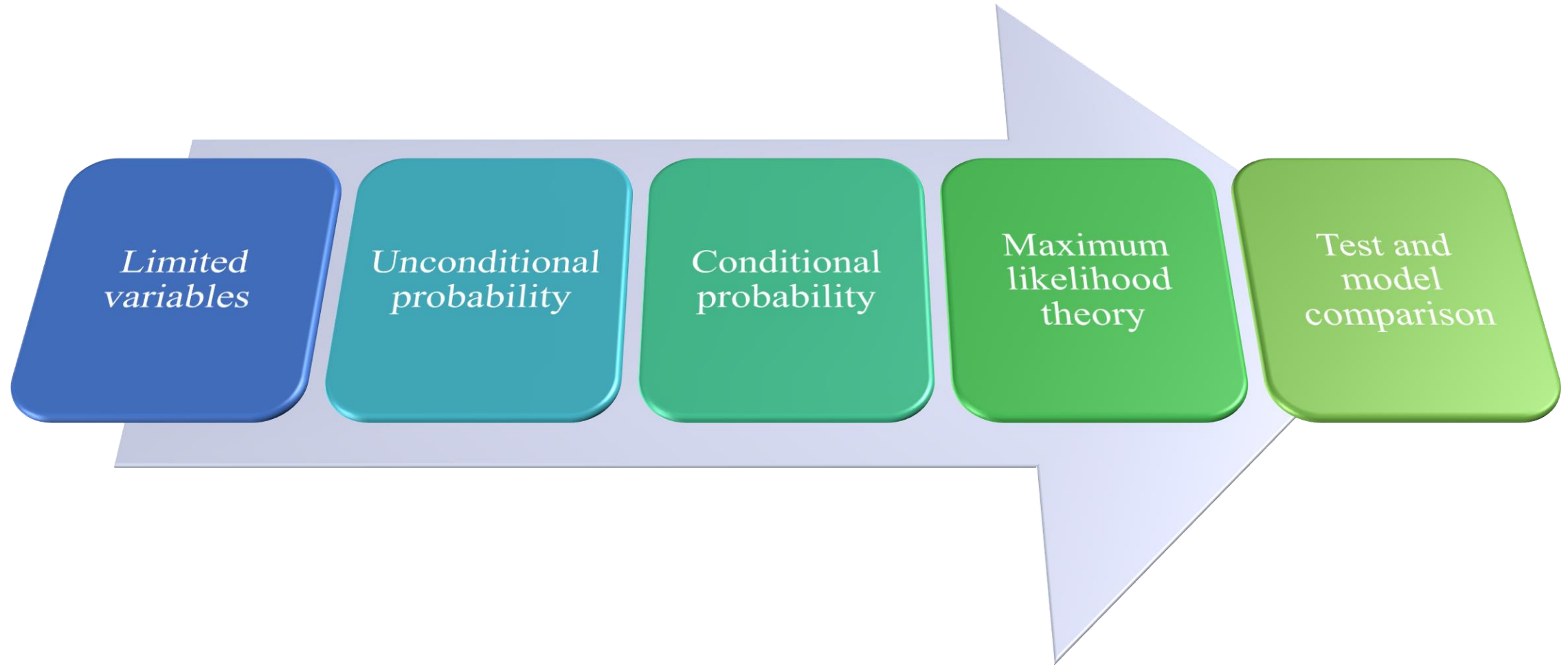
*Limited dependent variables*

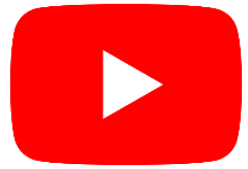


**ESSEC**  
BUSINESS SCHOOL

---

# Outline of this course





*Limited dependent variables*  
**Introduction**

# Limited variable ?

**Limited dependent variable:** Discrete variable takes on finite values.

**Binary variables:** Possible values: 0/1

**Any output of a dichotomous question:**

- Do I feel good today ? (time series)
- Did this bank collapse this year ? (cross-sectional data)
- Will this app become a killer app ? (cross-sectional data or time series)
- Is there a difference between low-cost and high-cost cheese cubes ?

9,96 euro/kg



23,27 euro/kg

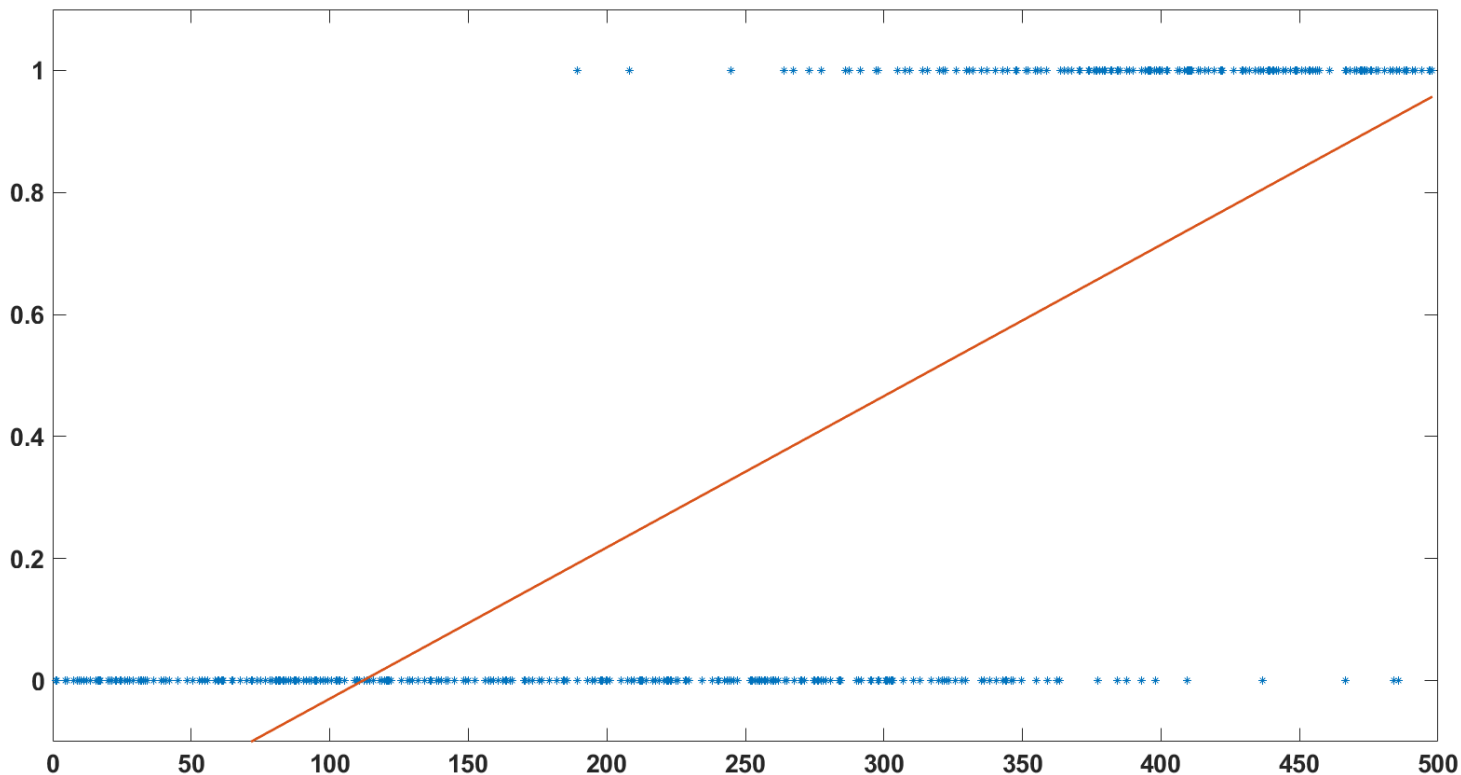


# Linear regression with limited variables

**Dependent variable:** Paid dividend in 2019

**Explanatory variable:** Market capitalisation (in billions \$)

$$\text{Div}_t = -0.27 + 0.0025\text{cap}_t$$



**✗ Interpretation as a probability?**

$$\text{cap}_t = 50 \rightarrow \text{Div}_t = -0.15$$

$$\text{cap}_t = 1000 \rightarrow \text{Div}_t = 2.2$$



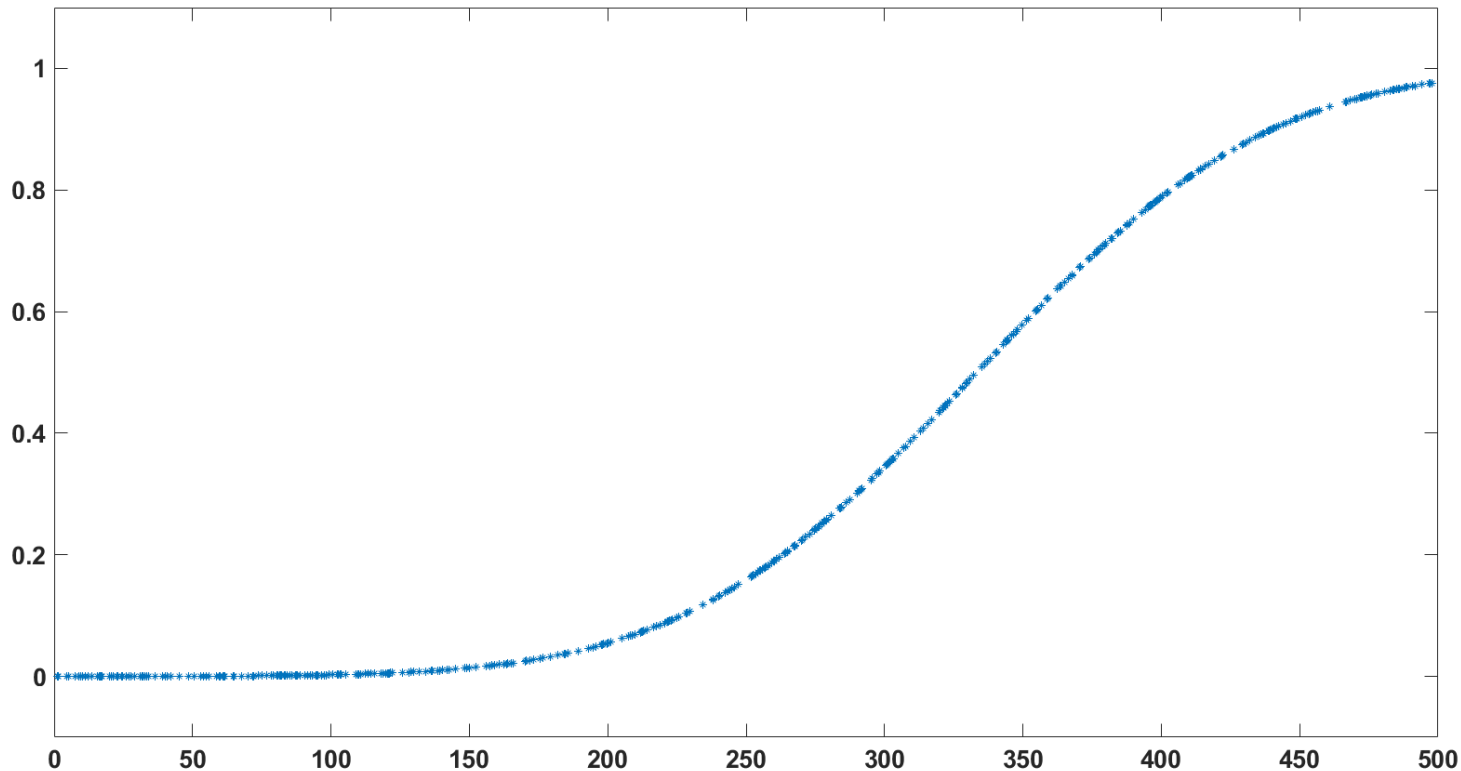
**Useful for quickly checking  
if a relation exists**

# Linear regression with limited variables

Dependent variable: Paid dividend in 2019

Explanatory variable: Market capitalisation (in billions \$)

$$P[\text{Div}_t = 1 | \text{cap}_t] = f(\beta_1 + \beta_2 \text{cap}_t)$$



We do not observe  
the probability



How to infer probabilities  
from binary data ?

*Limited dependent variables*  
**Unconditional probability**

# Model for limited variables

Are you able to distinguish between a low-cost and an expensive cheese cube ?

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

## Unconditional probability

Average probability in the population of detecting a correct cube cheese

$$P[C_t = 1] = p$$

**How to estimate  $p$  ?**

## Conditional probability

Probability given some characteristics of detecting a correct cube cheese

$$P[C_t = 1|x_t] = f(\beta_1 + \beta_2 x_t)$$

**How to estimate  $\beta_1$  and  $\beta_2$  ?**



# Model for limited variables

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

Unconditional probability:  $P[C_t = 1] = p \iff P[C_t = 0] = 1 - p$

Bernoulli



Properties:

$$E(C_t) = p$$

$$V(C_t) = E(C_t^2) - E(C_t)^2 = p - p^2 = p(1 - p)$$

Application: Estimation of a proportion

Survey and polls

When individuals are humans

$$C_t = 1$$



$$C_t = 0$$



Other proportions

When individuals are not humans

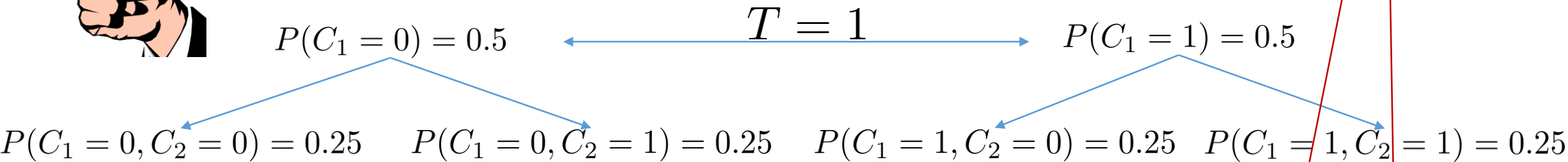


# Unconditional probability

$$P[C_t = 1] = p$$
$$N = \sum_{t=1}^T c_t = \# \text{ of } 1$$

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

Let us assume that  $p = 0.5$



**Likelihood function for two flips:**

$$P(C_1 = 0, C_2 = 1|p) = P(C_1 = 0|p)P(C_2 = 1|p) = (1 - p)p = (1 - p)^{\# \text{ of } 0} p^{\# \text{ of } 1}$$



**General likelihood function:**

$$P(C_1 = c_1, C_2 = c_2, \dots, C_T = c_T|p) = \prod_{t=1}^T P(C_t = c_t|p) = \prod_{t=1}^T p^{c_t} (1 - p)^{1 - c_t} = p^N (1 - p)^{T - N}$$

Red circles highlight  $N$  and  $T - N$  in the final expression, with red lines connecting them to the definitions of  $N$  and  $T$  in the top right corner.

# Unconditional probability

$$P[C_t = 1] = p$$

## General likelihood function:

$$P(C_1 = c_1, C_2 = c_2, \dots, C_T = c_T | p) = \prod_{t=1}^T P(C_t = c_t | p) = \prod_{t=1}^T p^{c_t} (1 - p)^{1 - c_t} = p^N (1 - p)^{T - N}$$

Likelihood of observing a sample: depends on  $N$  and  $T$

Let us assume that  $p = 0.5$

$$N = 0$$
$$T = 2 \quad P(C_1 = 0, C_2 = 0) = 0.25$$



Prob = 0.25

$$N = 1$$
$$P(C_1 = 1, C_2 = 0) = 0.25$$

+

$$P(C_1 = 0, C_2 = 1) = 0.25$$



Prob = 0.5

$$N = 2$$
$$P(C_1 = 1, C_2 = 1) = 0.25$$



Prob = 0.25



*The likelihood of the number of successes can be expressed as*

$$f(N|T, p) = \sum_{i=1}^{\# \text{success} = N} p^N (1 - p)^{T - N}$$

# Unconditional probability

$$P[C_t = 1] = p$$
$$N = \sum_{t=1}^T C_t = \# \text{ of } 1$$

## Additional examples

	$p = 0.5$			
	$N = 0$	$N = 1$	$N = 2$	$N = 3$
$T = 1$	Prob = 0.5	Prob = 0.5		
$T = 2$	Prob = 0.25	Prob = 0.5	Prob = 0.25	
$T = 3$	Prob = 0.125	Prob = 0.375	Prob = 0.375	Prob = 0.125

## Binomial distribution

1' to 6'



$$f(N|T, p) = \sum_{i=1}^{\# \text{success} = N} p^N (1 - p)^{T-N} = \binom{T}{N} p^N (1 - p)^{T-N}$$

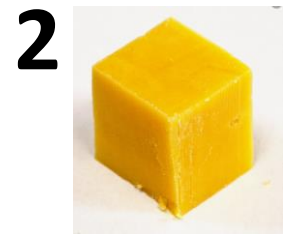
*The sum of Bernoulli random variables follows a Binomial distribution:*

$$C_t \sim \text{Be}(p) \text{ then } N = \sum_{t=1}^T C_t \sim \text{Bino}(p, T)$$

# Cheese tasting



One trial:



$$P[C_t = 1] = p$$

Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

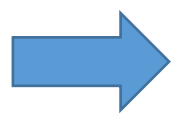
Unable to tell them apart  $\rightarrow$  Random guess:

$p = 0.5$	$N = 0$	$N = 1$	$\bullet \bullet \bullet$	$N = 8$	$N = 9$	$N = 10$
$T = 10$	Prob = 0	Prob = 0.01		Prob = 0.04	Prob = 0.01	Prob = 0

Link to statistical tests

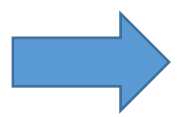
$$P[N \geq 8] = 0.05$$

Significant level of 99%



One mistake is allowed!

Significant level of 95%



Two mistake is allowed

Are we able to statistically detect the expensive cheese ?

$P[C_t = 1] = p$

# Unconditional probability

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

How to estimate  $p$  ?

$$p = 0.5, p = 0.25$$

	$N = 0$	$N = 1$	$N = 2$	$N = 3$
$T = 1$	Prob = 0.5 Prob = 0.75	Prob = 0.5 Prob = 0.25		
$T = 2$	Prob = 0.25 Prob = 0.562	Prob = 0.5 Prob = 0.375	Prob = 0.25 Prob = 0.0625	
$T = 3$	Prob = 0.125 Prob = 0.421	Prob = 0.375 Prob = 0.421	Prob = 0.375 Prob = 0.14	Prob = 0.125 Prob = 0.015

**Binomial distribution:**  $\max_p f(N|T, p) = \max_p \binom{T}{N} p^N (1 - p)^{T - N}$

Maximum likelihood estimator (MLE)

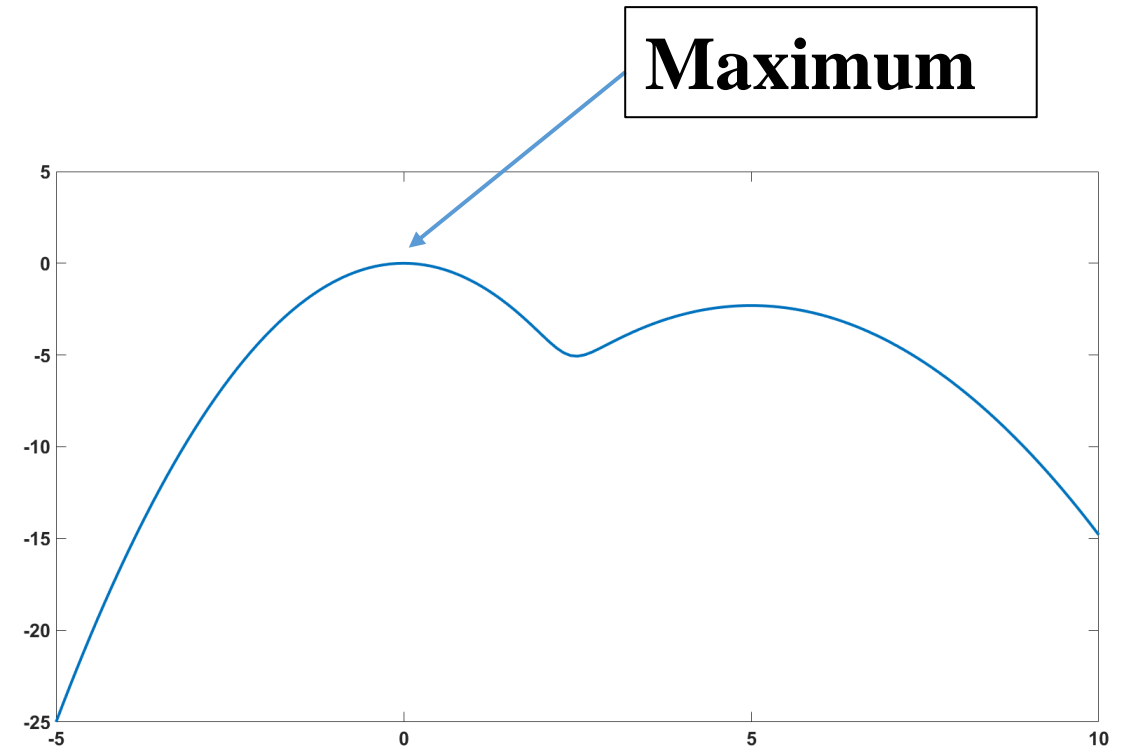
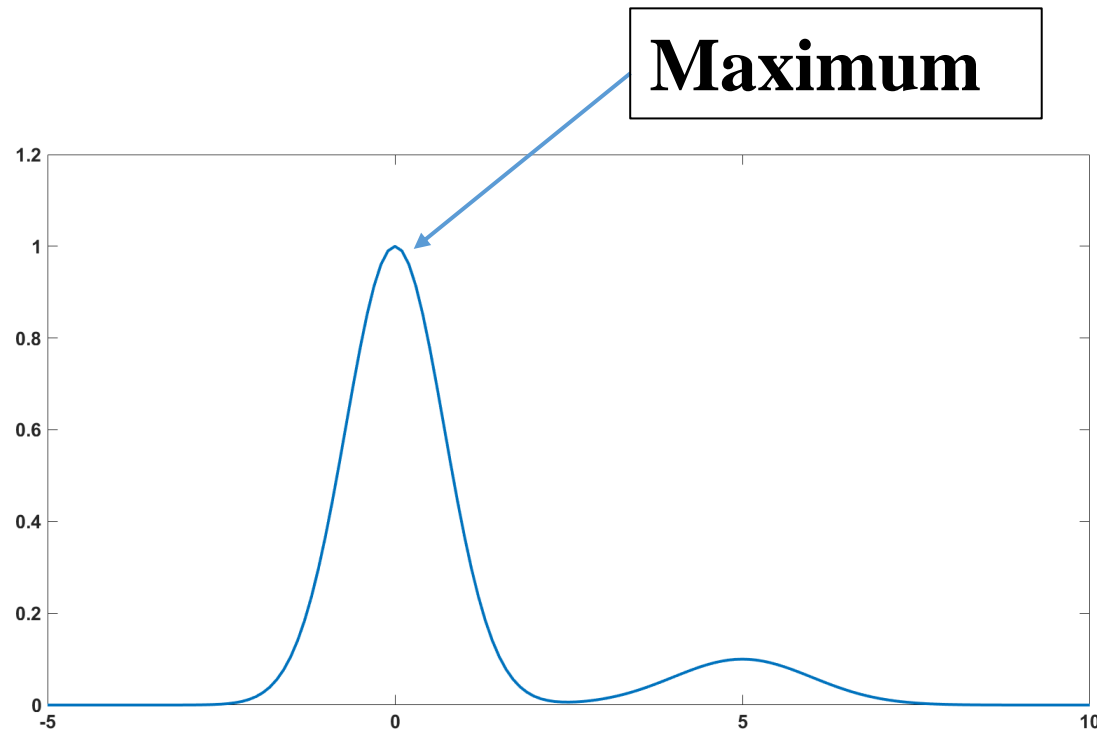
# Unconditional probability

$$P[C_t = 1] = p$$
$$N = \sum_{t=1}^T C_t$$

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

$\max_p f(N|T, p) = \max_p \binom{T}{N} p^N (1-p)^{T-N}$   **Complicated function to optimize**

If  $\hat{p}$  maximizes  $f(N|T, p)$  then  $\hat{p}$  maximizes  $\ln f(N|T, p)$ .



# Unconditional probability

$$P[C_t = 1] = p$$
$$N = \sum_{t=1}^T C_t$$

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

$$\max_p \ln f(N|T, p) = \max_p \ln \binom{T}{N} + N \ln(p) + (T - N) \ln(1 - p)$$

**Estimator of p:**  $\hat{p} = \frac{N}{T}$

**Number of successes over the number of trials**

**Proof:**

$$\begin{aligned} \frac{d \ln f(N|T, p)}{dp} &= \frac{N}{p} - \frac{(T - N)}{1 - p} (= 0), \\ \frac{N}{\hat{p}} &= \frac{(T - N)}{1 - \hat{p}}, \\ N - N\hat{p} &= (T - N)\hat{p}, \\ \hat{p} &= \frac{N}{T} \end{aligned}$$



# Unconditional probability

$$P[C_t = 1] = p$$
$$N = \sum_{t=1}^T C_t$$

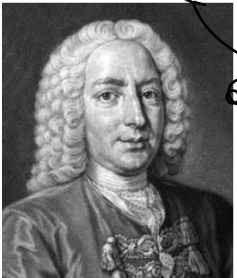
## Properties of the estimator ?

Note that:  $\hat{p} = \frac{N}{T} = \frac{1}{T} \sum_{t=1}^T C_t = \bar{C}$

**Average of random variables!**

What kind of random variable is  $C_t$  ?

Bernoulli



Properties:

$$E(C_t) = p$$

$$V(C_t) = E(C_t^2) - E(C_t)^2 = p - p^2 = p(1 - p)$$

Assumption: No dependence between the trials



e.g. No learning curve

# Unconditional probability

**Properties of the estimator:**  $\hat{p} = \bar{C}$

## Unbiasedness

$$E(\hat{p}) = \frac{1}{T} \sum_{t=1}^T E(C_t)$$



$$E(\hat{p}) = p$$

## Consistency

$$V(\hat{p}) = \frac{1}{T^2} \sum_{t=1}^T V(C_t)$$

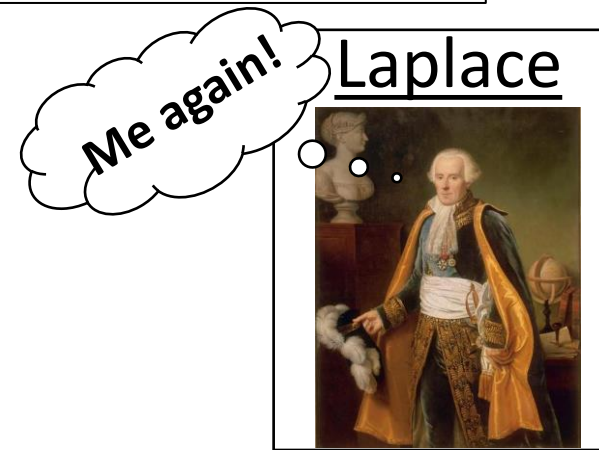


$$V(\hat{p}) = \frac{p(1-p)}{T}$$

**Also: Implication of the LLN**

## **Distribution of the estimator ?**

**Central limit theorem:**  $\hat{p} = \frac{1}{T} \sum_{t=1}^T C_t \rightarrow N\left(p, \frac{p(1-p)}{T}\right)$




Laplace

# Statistical test when T large

$$E(C_t) = p$$
$$V(C_t) = p(1 - p)$$

Central limit theorem:  $\hat{p} = \frac{1}{T} \sum_{t=1}^T C_t \rightarrow N(p, \frac{p(1-p)}{T})$

Hypothesis:  $H_0 : p = 0.4$  vs  $H_1 : p \neq 0.4$

Under the null:  $\hat{p} \sim N(0.4, \frac{0.24}{T})$    $\frac{\hat{p}-0.4}{\sqrt{\frac{0.24}{T}}} \sim N(0, 1)$

$$\text{For large } T: Z_p \equiv \frac{\hat{p}-p}{\sqrt{\frac{p(1-p)}{T}}} \sim N(0, 1)$$

**Any estimate of the statistic is a realization of a N(0,1)**

**Take  
Away**

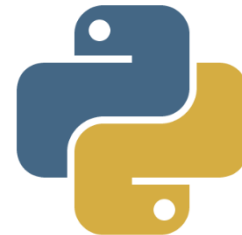
**For large T, no need to use the binomial cumulative density function for performing statistical tests.**

# Empirical exercise

**Iphone application:** What is the proportion of Game applications ?

**What do we learn in this exercise ?**

- How to estimate the unconditional probability from a limited variable.
- How to maximize a likelihood function.
- How to test a probability.





*Limited dependent variables*  
**Conditional probability**

# Conditional probability

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct

Conditional probability:  $P[C_t = 1|x_t] = f(\beta_1 + \beta_2 x_t)$

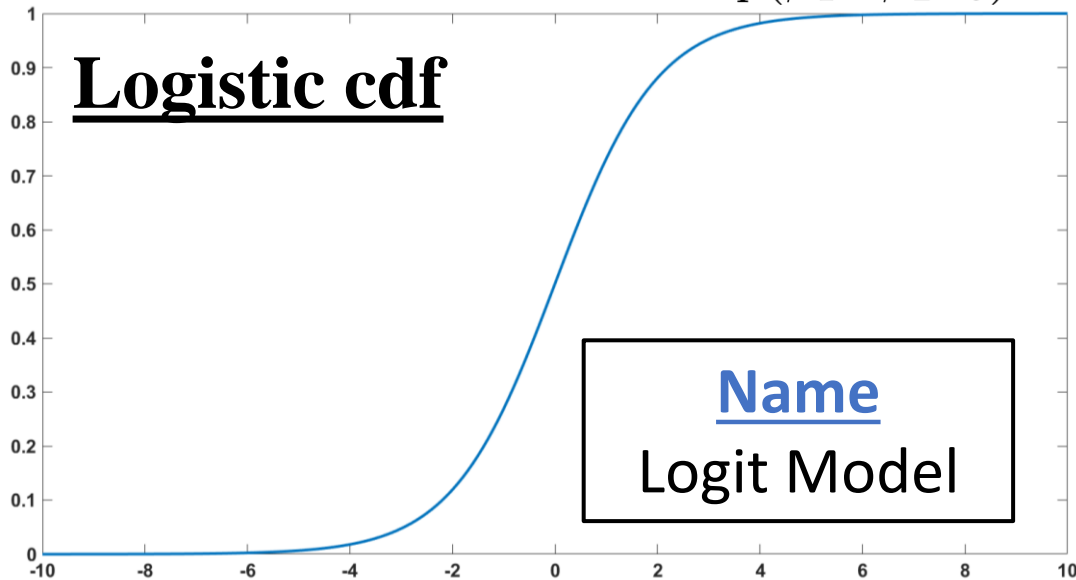
**Which function to choose ?**



**Any cumulative density function on the real support**

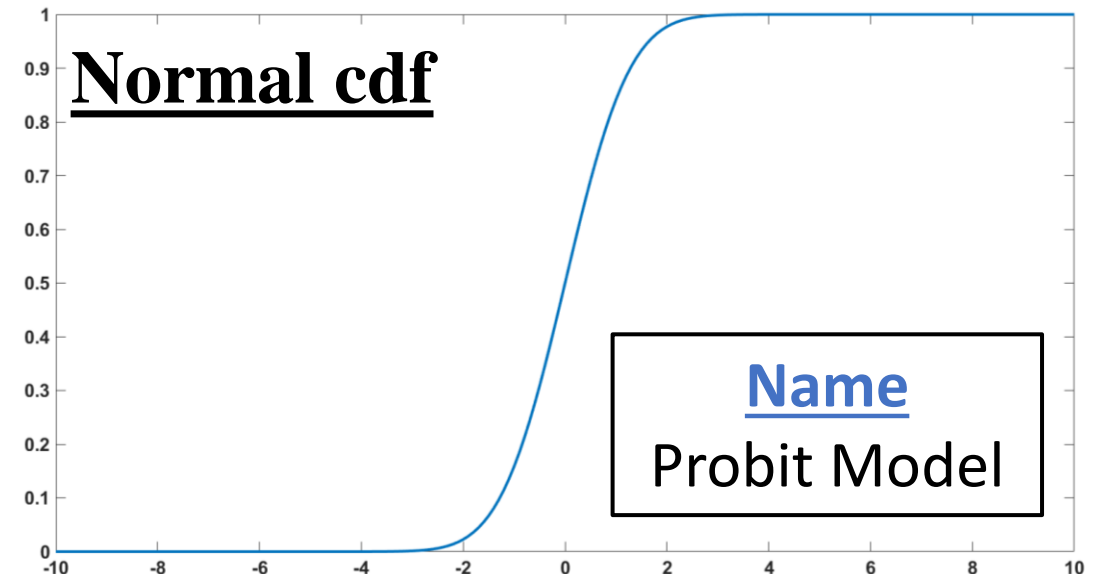
$$f(\beta_1 + \beta_2 x_t) = \frac{\exp(\beta_1 + \beta_2 x_t)}{1 + \exp(\beta_1 + \beta_2 x_t)}$$

**Logistic cdf**



$$f(\beta_1 + \beta_2 x_t) = \Phi(\beta_1 + \beta_2 x_t) = \int_{-\infty}^{\beta_1 + \beta_2 x_t} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) dz$$

**Normal cdf**



# Conditional probability

**Dependent variable:** Distinguishing the cheese cube  $\rightarrow C_t = 1$  if correct



**How to estimate a conditional probability ?**

$$P[C_t = 1|x_t] = f(\beta_1 + \beta_2 x_t) = f_\beta(x_t)$$

$$P(C_1 = 0|x_1) = 1 - f_\beta(x_1) \xleftrightarrow{T=1} P(C_1 = 1|x_1) = f_\beta(x_1)$$

$T = 2$

$$\begin{aligned} P(C_1 = 0, C_2 = 0|x_1, x_2) &= (1 - f_\beta(x_1))(1 - f_\beta(x_2)) & P(C_1 = 1, C_2 = 0|x_1, x_2) &= f_\beta(x_1)(1 - f_\beta(x_2)) \\ P(C_1 = 0, C_2 = 1|x_1, x_2) &= (1 - f_\beta(x_1))f_\beta(x_2) & P(C_1 = 1, C_2 = 1|x_1, x_2) &= f_\beta(x_1)f_\beta(x_2) \end{aligned}$$

**General likelihood function:**

$$P(C_1 = c_1, C_2 = c_2, \dots, C_T = c_T|x_1, \dots, x_T) = \prod_{t=1}^T P(C_t = c_t|x_t) = \prod_{t=1}^T f_\beta(x_t)^{c_t} (1 - f_\beta(x_t))^{1-c_t}$$

**Cannot be simplified with the number of successes  $N$  and of trials  $T$**

$$\prod_{t=1}^T f_\beta(x_t)^{c_t} (1 - f_\beta(x_t))^{1-c_t} \neq p^N (1 - p)^{T-N}$$

# Conditional probability

$$P[C_t = 1|x_t] = f_\beta(x_t)$$

## Maximum likelihood estimator:

$$\max_{\beta_1, \beta_2} P(C_1 = c_1, C_2 = c_2, \dots, C_T = c_T | x_1, \dots, x_T) = \max_{\beta_1, \beta_2} \prod_{t=1}^T f_\beta(x_t)^{c_t} (1 - f_\beta(x_t))^{1-c_t}$$



**Complicated function to optimize**



$$\max_{\beta_1, \beta_2} \sum_{t=1}^T c_t \ln[f_\beta(x_t)] + (1 - c_t) \ln(1 - f_\beta(x_t))$$

**Logistic function:**  $f_\beta(x_t) = \frac{\exp(\beta_1 + \beta_2 x_t)}{1 + \exp(\beta_1 + \beta_2 x_t)}$

**Normal cdf:**  $f(\beta_1 + \beta_2 x_t) = \Phi(\beta_1 + \beta_2 x_t) = \int_{-\infty}^{\beta_1 + \beta_2 x_t} \frac{1}{\sqrt{2\pi}} \exp(-\frac{z^2}{2}) dz$

**Likelihood function cannot be maximized analytically**



# Conditional probability

$$P[C_t = 1|x_t] = \Phi(\beta_1 + \beta_2 x_t)$$

**Normal cdf:**  $f(\beta_1 + \beta_2 x_t) = \Phi(\beta_1 + \beta_2 x_t) = \int_{-\infty}^{\beta_1 + \beta_2 x_t} \frac{1}{\sqrt{2\pi}} \exp(-\frac{z^2}{2}) dz$

**Maximum likelihood estimator with Normal cdf:**

$$\max_{\beta_1, \beta_2} \sum_{t=1}^T c_t \ln[\Phi(\beta_1 + \beta_2 x_t)] + (1 - c_t) \ln(1 - \Phi(\beta_1 + \beta_2 x_t))$$

**We maximize the function using a statistical software:**

**Numerical algorithms:**

1. Steepest gradient (Based on the first derivative)
2. Newton-Raphson method (Based on the first two derivatives)
3. Nelder-Mead algorithm (derivative-free)
4. Heuristic algorithms (Particle swarm, Differential Evolution, ...)



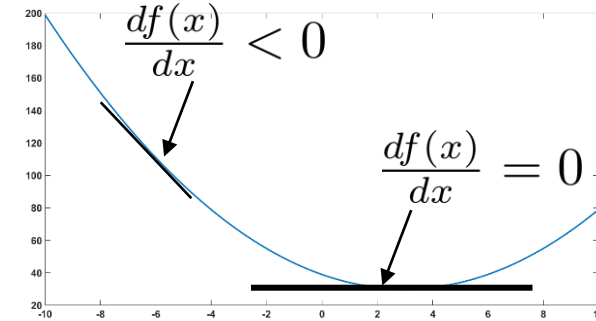
Focus on the Newton-Raphson method

# Numerical optimization

e.g. function:  $f(x) = x^2 + 4x + 2$

**Newton-Raphson:**  $\hat{x} = \operatorname{argmin}_x f(x)$

**Idea:** At a minimum, the derivative must be zero:  $\frac{df(x)}{dx} \big|_{x=\hat{x}} = 0$



**Algorithm:**

1. Start from an initial point:  $x_0$
2. Find a promising new point:  $x_1 = x_0 + t$

**If the promising point is a minimum, its derivative is zero:**  $\frac{df(x)}{dx} \big|_{x=x_1} = 0$

**The derivative at the promising point is approximated by a Taylor expansion:**

$$\underbrace{\frac{df(x)}{dx} \big|_{x=x_1}}_{=0} \approx \frac{df(x)}{dx} \big|_{x=x_0} + \underbrace{(x_1 - x_0)}_{=t} \frac{d^2 f(x)}{dx^2} \big|_{x=x_0} \Rightarrow t = - \left( \frac{d^2 f(x)}{dx^2} \big|_{x=x_0} \right)^{-1} \frac{df(x)}{dx} \big|_{x=x_0}$$

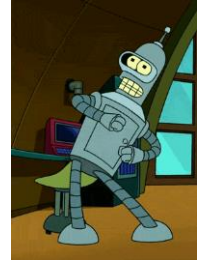
3. Check if the promising point is a minimum, if not come back to 1 with  $x_0 = x_1$

# Numerical optimization

**Illustration with a function:**  $f(x) = x^2 + 4x + 2 \longrightarrow \frac{df(x)}{dx} = 2x + 4$  and  $\frac{d^2 f(x)}{dx^2} = 2$

## Newton-Raphson algorithm:

1. Start from an initial point:  $x_0 = 5$
2. Find a promising new point:  $x_1 = x_0 + t$  with  $t = -\left(\frac{d^2 f(x)}{dx^2} \Big|_{x=x_0}\right)^{-1} \frac{df(x)}{dx} \Big|_{x=x_0}$   
 $\longrightarrow t = -\frac{1}{2}(2x_0 + 4) = -7 \longrightarrow x_1 = -2$
3. Check if the promising point is a minimum:  $\frac{df(x)}{dx} \Big|_{x=x_1} = 2(-2) + 4 = 0$



For more complex function, it requires several iterations

## Why does it work in one iteration for quadratic functions ?

The derivative is linear so the Taylor expansion is exact.

# Conditional probability

$$P[C_t = 1|x_t] = \Phi(\beta_1 + \beta_2 x_t)$$

## Maximum likelihood estimator with Normal cdf:

$$\max_{\beta_1, \beta_2} \sum_{t=1}^T c_t \ln[\Phi(\beta_1 + \beta_2 x_t)] + (1 - c_t) \ln(1 - \Phi(\beta_1 + \beta_2 x_t))$$

Using a Numerical algorithm, we get our estimates:  $\hat{\beta}_1, \hat{\beta}_2$  ✓

### Two issues:

1. How to interpret these estimates ?
2. What are the statistical properties of our estimator ?

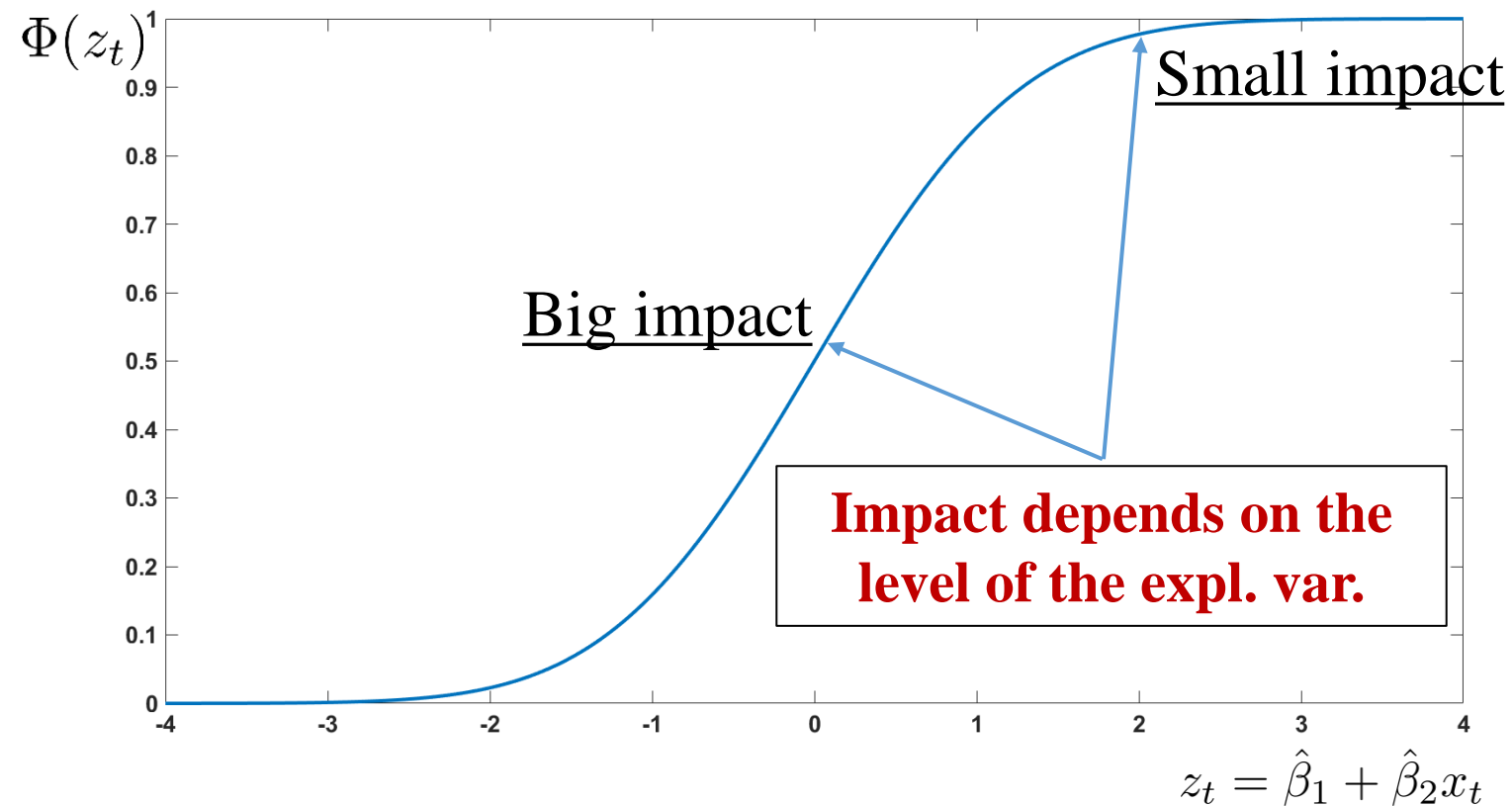


# Interpretation

We have our estimates:  $\hat{\beta}_1, \hat{\beta}_2$

**How do I interpret them ?**

Estimates are related to the probability:  $P[C_t = 1|x_t] = \Phi(\hat{\beta}_1 + \hat{\beta}_2 x_t)$



## Signs

$\hat{\beta}_2 > 0$  : prob increases with  $x_t$

$\hat{\beta}_2 < 0$  : prob decreases with  $x_t$

## Exact interpretation:

$$\frac{dP[C_t=1|x_t]}{dx_t} = \underbrace{\frac{d\Phi(\hat{\beta}_1 + \hat{\beta}_2 x_t)}{d(\hat{\beta}_1 + \hat{\beta}_2 x_t)}}_{\text{Density function of } N(0,1)} \hat{\beta}_2$$

Density function of  $N(0,1)$

# Statistical test

We have our estimates:  $\hat{\beta}_1, \hat{\beta}_2$

**Without analytical formulas, how can we study the statistical properties ?**



**Maximum likelihood theory provides a solution**

*Limited dependent variables*

**Introduction to the maximum likelihood theory  
(Relevant for empirical works)**

# Maximum likelihood estimator

$$E(C_t) = p$$
$$V(C_t) = p(1 - p)$$

## Maximum likelihood estimator (MLE):

$$\hat{p} = \operatorname{argmax}_p P(C_1, C_2, \dots, C_T) = \operatorname{argmax}_p \prod_{t=1}^T p^{C_t} (1 - p)^{1 - C_t}$$

### Intuition:

**Curvature gives a sense of the variance**

$$\hat{p} \sim N(p, \sigma_p^2)$$



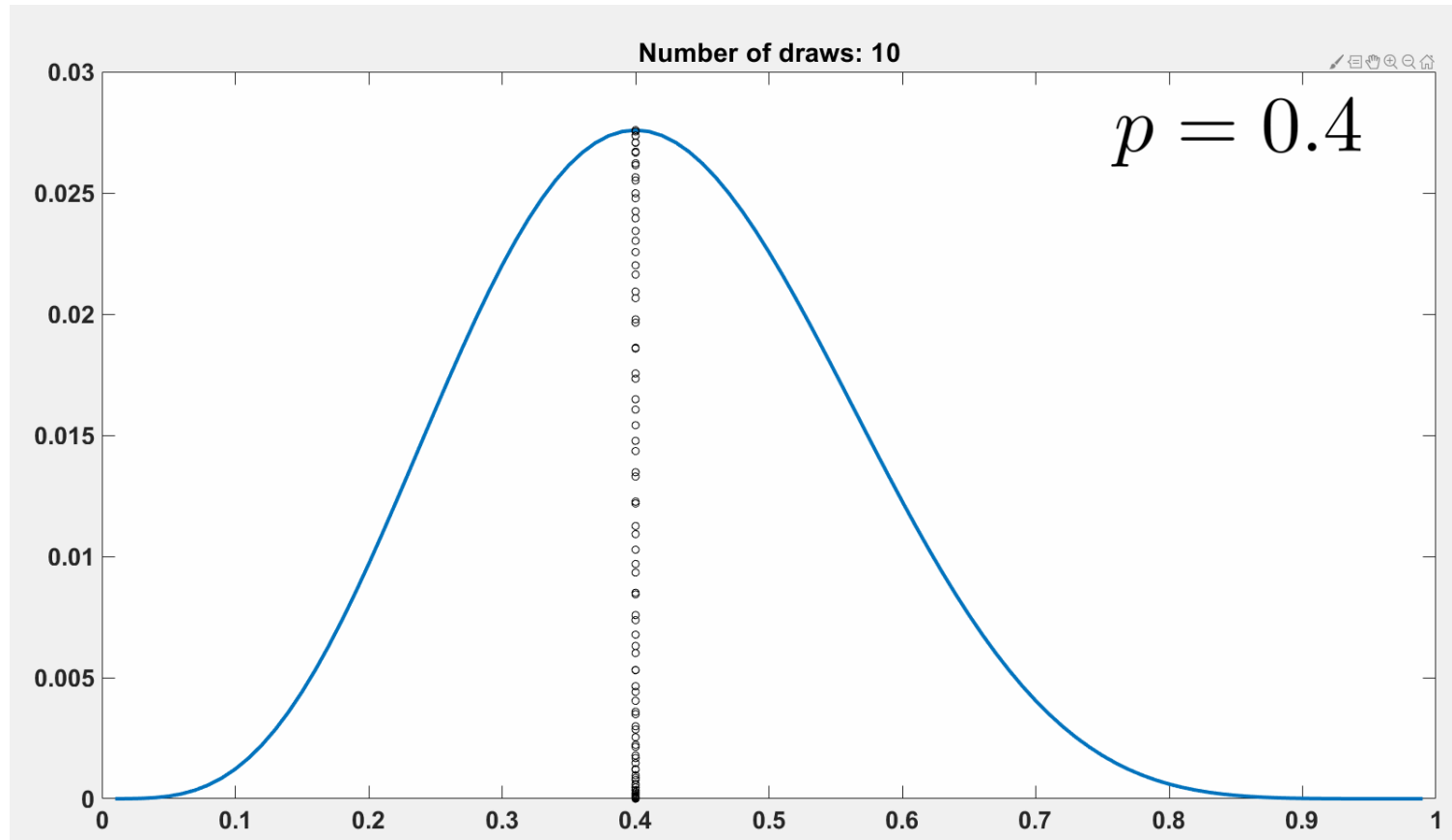
$$\ln f(\hat{p}|p, \sigma_p^2) = -\frac{1}{2} \ln(2\pi\sigma_p^2) - \frac{(\hat{p}-p)^2}{2\sigma_p^2}$$



$$\frac{d^2 \ln f(\hat{p}|p, \sigma_p^2)}{dp^2} = -\frac{1}{\sigma_p^2}$$



$$\sigma_p^2 = \left( -\frac{d^2 \ln f(\hat{p}|p, \sigma_p^2)}{dp^2} \right)^{-1}$$





# In a Nutshell: Maximum likelihood theory

**Given the log-likelihood function:**

$$\ln f(c_{1:T}|\beta = \{\beta_1, \beta_2\}) = \sum_{t=1}^T c_t \ln[\Phi(\beta_1 + \beta_2 x_t)] + (1 - c_t) \ln(1 - \Phi(\beta_1 + \beta_2 x_t))$$

**1. Find the maximum likelihood estimate:**  $\frac{d \ln f(c_{1:T}|\beta)}{d\beta} = 0 \quad \longrightarrow \quad \hat{\beta} = \{\hat{\beta}_1, \hat{\beta}_2\}$

**2. Compute the second derivative:**  $H(\beta) = \frac{d^2 \ln f(c_{1:T}|\beta)}{d\beta^2}$

**3. Maximum likelihood estimator:**

$$\hat{\beta}_1 \sim N(\beta_1, ((-H(\hat{p}))^{-1})_{11}) \text{ for large } T.$$



# In a Nutshell: Maximum likelihood theory

## Example: Unconditional probability

### Slide 11: Log-likelihood function is given by

$$\ln P(C|p) = \sum_{t=1}^T c_t \ln[p] + (1 - c_t) \ln(1 - p)$$

1. Find the maximum likelihood estimate:  $\frac{d \ln P(C|p)}{dp} = 0$

$$\frac{d \ln P(C|p)}{dp} = \frac{\sum_{t=1}^T c_t}{p} - \frac{\sum_{t=1}^T (1 - c_t)}{1 - p} \quad \Rightarrow \quad \frac{\hat{p}}{1 - \hat{p}} = \frac{N}{T - N} \quad \Rightarrow \quad \hat{p} = \frac{\sum_{t=1}^T c_t}{T} = \frac{N}{T}$$

2. Compute the second derivative:  $H(p) = \frac{d^2 \ln P(C|p)}{dp^2}$

$$\frac{d^2 \ln P(C|p)}{dp^2} = -\frac{N}{p^2} - \frac{T - N}{(1 - p)^2} \quad \Rightarrow \quad H(\hat{p}) = -\frac{T}{\hat{p}^2} \left[ \frac{N}{T} + \frac{T - N}{T} \frac{\hat{p}^2}{(1 - \hat{p})^2} \right] \quad \Rightarrow \quad H(\hat{p}) = -\frac{T}{\hat{p}(1 - \hat{p})}$$

3. Maximum likelihood estimator:

$$\hat{p} \sim N(p, (-H(\hat{p}))^{-1}) \text{ for large } T \quad \Rightarrow \quad \hat{p} \sim N(p, \frac{\hat{p}(1 - \hat{p})}{T})$$

# Empirical exercise

**Iphone application:** Does the app score help for becoming a killer app ?

**What do we learn in this exercise ?**

- How to estimate the conditional probability from a limited variable.
- How to maximize a likelihood function.
- How to test the value of a parameter.



*Limited dependent variables*

**Introduction to the maximum likelihood theory  
(mathematical part)**

# Maximum likelihood estimator

$$E(C_t) = p$$
$$V(C_t) = p(1 - p)$$

Central limit theorem:  $\hat{p} = \frac{1}{T} \sum_{t=1}^T C_t \rightarrow N(p, \frac{p(1-p)}{T})$

Maximum likelihood theory:

Likelihood function:  $P(C|p) = \prod_{t=1}^T p^{C_t} (1-p)^{1-C_t}$

Log-likelihood function:  $\ln P(C|p) = \sum_{t=1}^T C_t \ln p + (1 - C_t) \ln(1 - p)$

First derivative:  $\frac{d \ln P(C|p)}{dp} = \sum_{t=1}^T \frac{C_t}{p} - \frac{(1-C_t)}{1-p} \Rightarrow E\left[\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right] = \frac{E(C_t)}{p} - \frac{E(1-C_t)}{1-p} = 0$

Asymptotically, the true probability should maximize the likelihood function!

$$\frac{d \ln P(C|p)}{dp} \Big|_{p=p} = \sum_{t=1}^T \frac{C_t}{p} - \frac{(1-C_t)}{1-p} \approx 0 \text{ for large } T.$$

$$\frac{d \ln P(C|p)}{dp} \Big|_{p=p} \sim N\left(0, T V\left(\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right)\right) \text{ for large } T.$$

# Maximum likelihood estimator

$$\frac{d \ln P(C|p)}{dp} \Big|_{p=p} \sim N\left(0, T V\left(\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right)\right) \text{ for large } T.$$

$$\begin{aligned} V\left(\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right) &= E\left(\left(\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right)^2\right), \\ &= \frac{1}{p^2}p + \frac{1}{(1-p)^2}(1-p), \\ &= \frac{1}{p(1-p)}. \end{aligned}$$

Maximum likelihood estimator:  $\frac{d \ln P(C|p)}{dp} \Big|_{p=\hat{p}} = 0 \Rightarrow \hat{p} = \frac{1}{T} \sum_{t=1}^T C_t$

## Maximum likelihood theory

Mean Value Theorem:

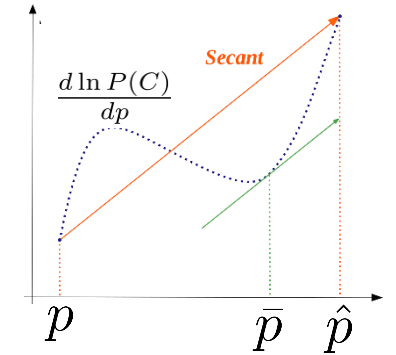
$$\underbrace{\frac{d \ln P(C|p)}{dp} \Big|_{p=\hat{p}}}_{=0} = \frac{d \ln P(C|p)}{dp} \Big|_{p=p} + (\hat{p} - p) \frac{d^2 \ln P(C|p)}{dp^2} \Big|_{p=\bar{p}}$$



$$(\hat{p} - p) = \left( - \frac{d^2 \ln P(C|p)}{dp^2} \Big|_{p=\bar{p}} \right)^{-1} \frac{d \ln P(C|p)}{dp} \Big|_{p=p}$$



$$(\hat{p} - p) \sim N\left(0, \left( \frac{d^2 \ln P(C|p)}{dp^2} \Big|_{p=p} \right)^{-1} T V\left(\frac{C_t}{p} - \frac{(1-C_t)}{1-p}\right) \left( \frac{d^2 \ln P(C|p)}{dp^2} \Big|_{p=p} \right)^{-1} \right) \text{ for large } T.$$



# Maximum likelihood estimator

$$V\left(\frac{C_t}{p} - \frac{(1 - C_t)}{1 - p}\right) = \frac{1}{p(1 - p)}.$$

$$(p - \hat{p}) \sim N\left(0, \left(\frac{-d^2 \ln P(C)}{dp^2}\right)\bigg|_{p=p}\right)^{-1} \frac{T}{p(1-p)} \left(\frac{-d^2 \ln P(C)}{dp^2}\right)\bigg|_{p=p}^{-1} \text{ for large } T.$$

Variance of the estimator:

$$\begin{aligned} \frac{d^2 \ln P(C)}{dp^2} &= \frac{d}{dp} \left[ \sum_{t=1}^T \frac{C_t}{p} - \frac{(1 - C_t)}{1 - p} \right], \\ &= T \frac{1}{T} \sum_{t=1}^T \frac{C_t - 1}{(1 - p)^2} - \frac{C_t}{p^2}, \\ \frac{d^2 \ln P(C)}{dp^2} \bigg|_{p=p} &= TE\left(\frac{C_t - 1}{(1 - p)^2} - \frac{C_t}{p^2}\right), \text{ for large } T, \\ &= -\frac{T}{p(1 - p)} \text{ for large } T. \end{aligned}$$

$$(\hat{p} - p) \sim N\left(0, \frac{p(1-p)}{T}\right) \text{ for large } T.$$

# Maximum likelihood theory

## Maximum likelihood theory:

Log-likelihood function:  $\ln f(y_1, \dots, y_T | \theta) = \sum_{t=1}^T \ln f(y_t | \theta)$

Maximum likelihood estimator:  $\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{t=1}^T \ln f(y_t | \theta) \quad \Rightarrow \quad \sum_{t=1}^T \frac{d \ln f(y_t | \theta)}{d\theta} = 0$

Mean value Theorem: 
$$\underbrace{\frac{d \ln f(y_1, \dots, y_T)}{d\theta} \Big|_{\theta=\hat{\theta}}}_{=0} = \underbrace{\frac{d \ln f(y_1, \dots, y_T)}{d\theta} \Big|_{\theta=\theta_0}}_{\rightarrow_d N\left(0, TV\left(\frac{d \ln f(y_t | \theta)}{d\theta}\right)\right)} + (\hat{\theta} - \theta_0) \underbrace{\frac{d^2 \ln f(y_1, \dots, y_T)}{d\theta^2} \Big|_{\theta=\bar{\theta}}}_{\rightarrow_p TV\left(\frac{d \ln f(y_t | \theta)}{d\theta}\right)}$$

Properties of the MLE:  $\hat{\theta} \sim N\left(\theta_0, \left(\frac{d^2 \ln f(y_1, \dots, y_T)}{d\theta^2} \Big|_{\theta=\theta_0}\right)^{-1}\right)$  for large T.



**CAN: Consistent and Asymptotically Normally distributed**



# Maximum likelihood theory: Example

Linear regression with a constant:  $y_t = \beta_1 + \epsilon_t$

ML: Additional assumption:  $\epsilon_t \sim N(0, \sigma^2)$  (independent)

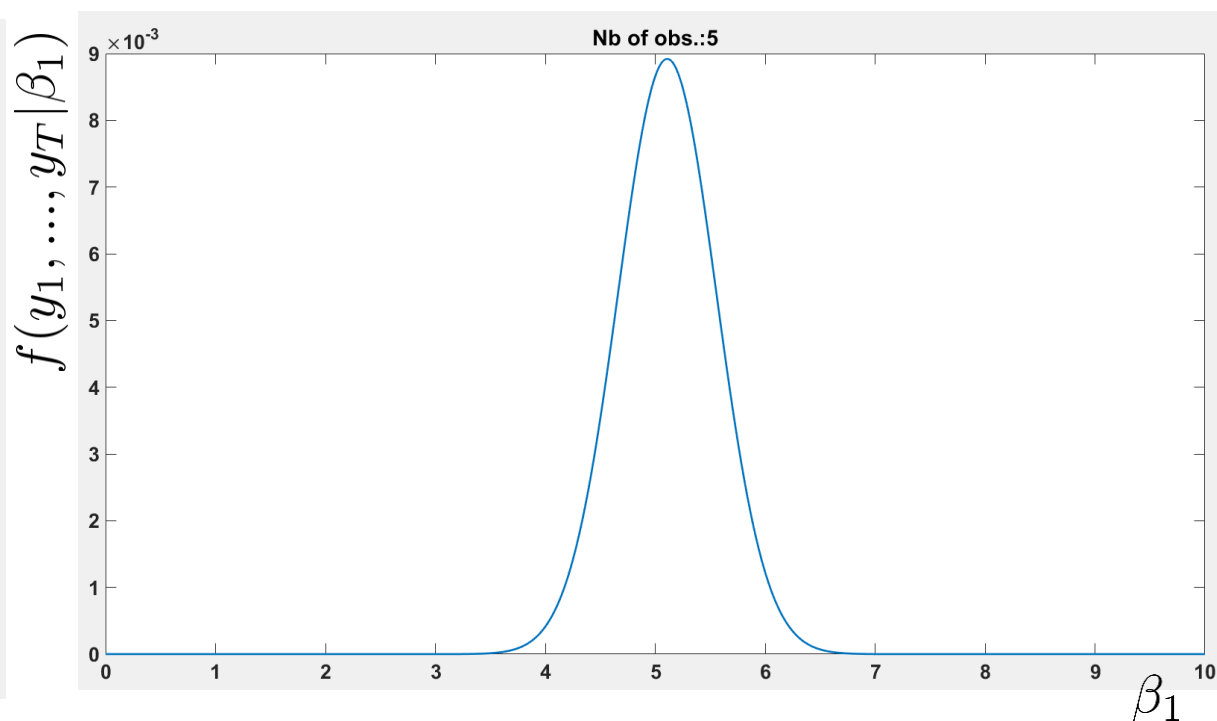
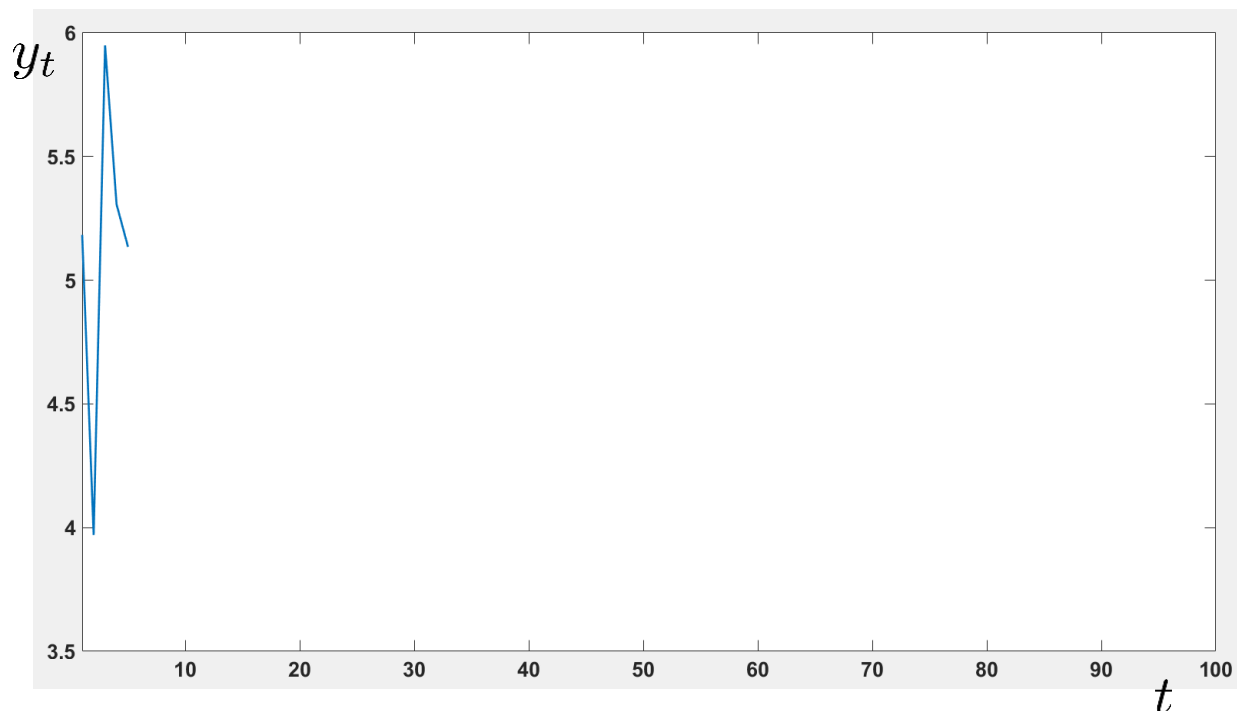
 **Model fully characterized:**  $y_t | \beta_1 \sim N(\beta_1, \sigma^2)$

OLS criterion:  $\hat{\beta}_1 = \operatorname{argmin}_{\beta_1} \sum_{t=1}^T \epsilon_t^2$

OLS estimator:  $\hat{\beta}_1 = \frac{1}{T} \sum_{t=1}^T y_t$

Large sample prop:  $\hat{\beta}_1 \sim N(\beta_1, \frac{\sigma^2}{T})$

ML criterion - Likelihood function:  $f(y_1, \dots, y_T | \beta_1) = \prod_{t=1}^T f(y_t | \beta_1) = \prod_{t=1}^T \frac{\exp\left(-\frac{(y_t - \beta_1)^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}$





# Maximum likelihood theory: Example

Linear regression with a constant:  $y_t = \beta_1 + \epsilon_t$

ML criterion - Likelihood function:


$$f(y_1, \dots, y_T | \beta_1) = \prod_{t=1}^T f(y_t | \beta_1) = \prod_{t=1}^T \frac{\exp\left(-\frac{(y_t - \beta_1)^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}$$

ML estimator :  $\hat{\beta}_1 = \operatorname{argmax}_{\beta_1} \ln f(y_1, \dots, y_T | \beta_1) = \operatorname{argmax}_{\beta_1} \sum_{t=1}^T \left(-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(y_t - \beta_1)^2}{2\sigma^2}\right)$

  $\hat{\beta}_1 = \operatorname{argmax}_{\beta_1} -\frac{1}{2\sigma^2} \sum_{t=1}^T (y_t - \beta_1)^2$    $\hat{\beta}_1 = \frac{1}{T} \sum_{t=1}^T y_t$

Same estimator as the OLS estimator

Large sample prop:  $\hat{\beta}_1 \sim N\left(\beta_1, \left(\frac{-d^2 \ln f(y_1, \dots, y_T | \beta_1)}{d\beta_1^2}\right)^{-1}\right) = N\left(\beta_1, \frac{\sigma^2}{T}\right)$

 
$$\begin{aligned} \frac{d^2 \ln f(y_1, \dots, y_T | \beta_1)}{d\beta_1^2} &= \frac{d^2}{d\beta_1^2} \left[ \sum_{t=1}^T \left( -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(y_t - \beta_1)^2}{2\sigma^2} \right) \right], \\ &= -\sum_{t=1}^T \frac{1}{\sigma^2} \quad \left( = -\frac{T}{\sigma^2} \right) \end{aligned}$$

OLS criterion:  $\hat{\beta}_1 = \operatorname{argmin}_{\beta_1} \sum_{t=1}^T \epsilon_t^2$

OLS estimator:  $\hat{\beta}_1 = \frac{1}{T} \sum_{t=1}^T y_t$

Large sample prop:  $\hat{\beta}_1 \sim N\left(\beta_1, \frac{\sigma^2}{T}\right)$

# Summary

## 1. Fully Characterized your model:

### Linear regression:

$$y_t|x_t \sim N(\beta_1 + \beta_2 x_t, \sigma^2)$$

### Probit model:

$$P[C_t = 1|x_t] = \Phi(\beta_1 + \beta_2 x_t)$$

## 2. Write down the likelihood function:

$$\prod_{t=1}^T \frac{\exp\left(-\frac{(y_t - \beta_1 - \beta_2 x_t)^2}{2\sigma^2}\right)}{\sqrt{2\pi\sigma^2}}$$

$$\prod_{t=1}^T \Phi(\beta_1 + \beta_2 x_t)^{c_t} (1 - \Phi(\beta_1 + \beta_2 x_t))^{1-c_t}$$

## 3. Compute the log-likelihood function:

$$\sum_{t=1}^T \left( -\frac{1}{2} \ln(2\pi\sigma^2) - \frac{(y_t - \beta_1 - \beta_2 x_t)^2}{2\sigma^2} \right)$$

$$\sum_{t=1}^T c_t \ln[\Phi(\beta_1 + \beta_2 x_t)] + (1 - c_t) \ln(1 - \Phi(\beta_1 + \beta_2 x_t))$$

## 4. Maximize the log-likelihood function with respect to the parameters

$$\hat{\beta}_1 = \bar{y} - \hat{\beta}_2 \bar{x}, \hat{\beta}_2 = \frac{\sum_{t=1}^T (x_t - \bar{x})(y_t - \bar{y})}{\sum_{t=1}^T (x_t - \bar{x})^2}, \hat{\sigma}^2 = \frac{\sum_{t=1}^T \hat{\epsilon}_t^2}{T}$$

No analytical formula

## 5. Compute the Hessian at the maximum likelihood estimate:

$$H(\hat{\theta}) = \frac{d^2 \ln f(y_1, \dots, y_T)}{d\theta d\theta}$$

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\sigma}^2 \end{pmatrix} \sim N\left( \begin{pmatrix} \beta_1 \\ \beta_2 \\ \sigma^2 \end{pmatrix}, (-H(\hat{\beta}_1, \hat{\beta}_2, \hat{\sigma}^2))^{-1} \right)$$

**Large T**

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \sim N\left( \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, (-H(\hat{\beta}_1, \hat{\beta}_2))^{-1} \right)$$

*Limited dependent variables*  
**Model comparison**

# Likelihood ratio test

$$\begin{aligned} m &= \# \text{ of restrictions} \\ &= K_U - K_R \end{aligned}$$

## How to compare nested models ?

**Likelihood ratio test (equivalent of the Fisher test in a ML context).**

**Full model (or unrestricted model):**

$$P[C_t = 1|x_t] = \Phi(\beta_1 + \beta_2 x_t)$$



$$\begin{aligned} \ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) = \\ \sum_{t=1}^T c_t \log[\Phi(\hat{\beta}_1 + \hat{\beta}_2 x_t)] + (1 - c_t) \log(1 - \Phi(\hat{\beta}_1 + \hat{\beta}_2 x_t)) \end{aligned}$$

**Nested model (or restricted model):**

$H_0$  : Restrictions on the model hold

$$P[C_t = 1|x_t] = \Phi(\beta_1) \Rightarrow H_0 : \beta_2 = 0$$



$$\begin{aligned} \ln f(c_1, \dots, c_T | \tilde{\beta}_1) = \\ \sum_{t=1}^T c_t \log[\Phi(\tilde{\beta}_1)] + (1 - c_t) \log(1 - \Phi(\tilde{\beta}_1)) \end{aligned}$$

### Likelihood ratio test

$$2[\ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) - \ln f(c_1, \dots, c_T | \tilde{\beta}_1)] \sim \chi^2(m) \text{ under } H_0$$

# Information Criteria

In general, how to compare models with the same dependent variable ?

Information criteria have been proposed for general comparisons.

Advantage: Simple and intuitive

Drawback: Multiple criteria exist and can lead to select different models.

## Example:

- Multiple models to compare: Model 1 to Model 10.
- Compute the information criterion for each model.
- The best model exhibits the highest value of the information criterion.

Information criteria (IC): Likelihood function evaluated at the MLE - penalty

Differs from one IC to another

Two popular information criteria: AIC and BIC

# Information Criteria

- Akaike Information criterion (**AIC**) :

$$AIC = \ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) - \underbrace{K}_{\text{penalty}_{AIC}} \quad \text{where } K \text{ denotes } \# \text{ parameters}$$

- Bayesian Information criterion (**BIC**) :

$$BIC = \ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) - \underbrace{\frac{K}{2} \log T}_{\text{penalty}_{BIC}}$$

**AIC penalty is smaller than the BIC penalty**



**BIC leads to select smaller models than AIC**

# AIC and BIC justifications

- Akaike Information criterion (**AIC**) :  $AIC = \ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) - K$

## Asymptotic justification:

- Selected model minimizes the Kullback-Leibler divergence (with respect to the true model)

➡ No claim that the selected model is the good one.

➡ The selected model should be good at predicting the data

- Bayesian Information criterion (**BIC**) :  $BIC = \ln f(c_1, \dots, c_T | \hat{\beta}, x_{1:T}) - \frac{K}{2} \log T$

## Asymptotic justification:

- BIC is proportional to the marginal likelihood of the Bayesian model with uninformative priors.

➡ Selected model is the true model!

➡ The true model must be among the models in competition.

**BIC is favoured in financial research**