

MATH 60604A  
Statistical modelling  
§ 4h - Logistic model for proportions

HEC Montréal  
Department of Decision Sciences

# Logistic model for proportions

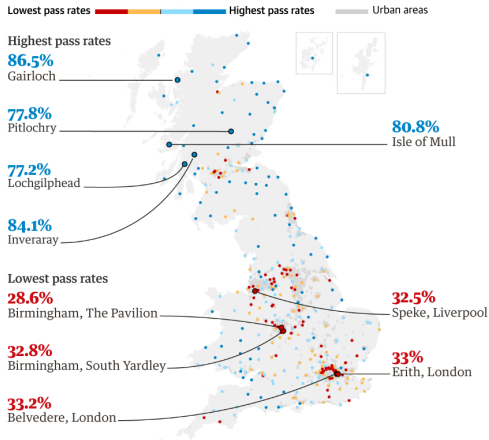
- Sometimes, we don't have access to individual records, but rather to aggregated counts such as the number of successes (out of  $m$  trials).
- We may use a binomial model instead by simply specifying the total number of trials associated to each number of successes.
- The parameter interpretation remains the same.

We consider the pass rate for all 346 Great-Britain driving license practical testing sites; the data are from 2018.

- 761 750 people succeeded in their exam out of 1 663 897 attempts.
- A news article from *The Guardian* hinted that exam takers in rural areas got an easy ride. Since we do not have a classification of urban/rural centers, we use the number of tests conducted as proxy.
- Other covariates are `sex` and the `region` for England; all of Scotland and Wales are pooled.

# Binomial model for driving license pass rate in Great-Britain

## Rural test centres tend to have higher pass rates than ones in cities



Source: The Guardian.

## SAS code to fit a logistic regression for binomial data

```
data gbdriving;
set statmod.gbdriving;
if(total < 500) then size="small";
else if (total < 1000) then size="medium";
else size = "large";
run;

proc logistic data=gbdriving;
class sex(ref="women") region(ref="London")
      size / param=glm;
model pass/total = sex region size /
      plr1 plcl expb;
run;
```

# Size of center per region

region	size		
	large	mediu	small
	N	N	N
East Midlands	40	3	3
East of England	54	.	.
London	48	4	6
North East England	29	5	8
North West England	63	3	2
Scotland	41	17	94
South East England	78	.	.
South West England	44	6	.
Wales	30	9	9
West Midlands	54	6	2
Yorkshire and the Hu	32	.	2

Scotland boasts the largest number of small centers (fewer than 500 exams per year).

# Model specification for Great-Britain driving licenses

---

## Model Information

<b>Data Set</b>	WORK.GBDRIVING
<b>Distribution</b>	Binomial
<b>Link Function</b>	Logit
<b>Response Variable (Events)</b>	pass
<b>Response Variable (Trials)</b>	total

---

<b>Number of Observations Read</b>	692
<b>Number of Observations Used</b>	692
<b>Number of Events</b>	761750
<b>Number of Trials</b>	1663897

---

---

## Model Fit Statistics

Criterion	Intercept and Covariates		
	Intercept Only	Log Likelihood	Full Log Likelihood
<b>AIC</b>	2294792.5	2278217.4	26619.303
<b>SC</b>	2294804.8	2278390.0	26791.848
<b>-2 Log L</b>	2294790.5	2278189.4	26591.303

---

---

## Type 3 Analysis of Effects

Effect	DF	Wald	
		Chi-Square	Pr > ChiSq
<b>sex</b>	1	8510.4974	<.0001
<b>region</b>	10	5565.9869	<.0001
<b>size</b>	2	1537.2919	<.0001

---

# Odds estimates for Great-Britain driving licenses data

---

<b>Odds Ratio Estimates and Profile-Likelihood Confidence Intervals</b>				
<b>Effect</b>	<b>Unit</b>	<b>Estimate</b>	<b>95% Confidence Limits</b>	
<b>sex men vs women</b>	1.0000	1.335	1.327	1.343
<b>region East Midlands vs London</b>	1.0000	1.279	1.262	1.297
<b>region East of England vs London</b>	1.0000	1.241	1.225	1.257
<b>region North East England vs London</b>	1.0000	1.500	1.475	1.524
<b>region North West England vs London</b>	1.0000	1.231	1.216	1.246
<b>region Scotland vs London</b>	1.0000	1.261	1.243	1.280
<b>region South East England vs London</b>	1.0000	1.257	1.243	1.271
<b>region South West England vs London</b>	1.0000	1.405	1.385	1.425
<b>region Wales vs London</b>	1.0000	1.447	1.423	1.472
<b>region West Midlands vs London</b>	1.0000	1.046	1.033	1.060
<b>region Yorkshire and the Hu vs London</b>	1.0000	1.094	1.078	1.110
<b>size large vs small</b>	1.0000	0.614	0.597	0.631
<b>size mediu vs small</b>	1.0000	0.766	0.741	0.792

---

All other things being constant,

- The odds of men are 33% higher than women of obtaining a driver license;
- Greater London is the region with the lowest success rate after accounting for the site volume; the odds of success are 50% higher in North East England and 44.7% higher in Wales, etc.
- The odds of success are 63% higher in small center than in large centers (1/0.614).
- All parameters are statistically significant.



## Remark on models for Bernoulli/binomial data

- While the deviance and Pearson  $\chi^2$  statistics are reported for logistic binomial model, their distribution depends on the unknown parameter vector  $\beta$ .
- As such, the deviance is approximately  $\chi_{n-p-1}^2$  only when the number of trials  $m$  is in the several thousands.
- Comparisons of deviance, which amount to likelihood ratio tests, are however valid.

# Revisiting the US road casualties example

We can fit a binomial model for the crash where the "event" is death.

Parameter Estimates and Profile-Likelihood Confidence Intervals				
Parameter		Estimate	95% Confidence Limits	
Intercept		-10.8702	-10.8913	-10.8495
time	night	0.2593	0.2372	0.2815
year	2018	0.2322	0.2101	0.2544

Odds Ratio Estimates and Profile-Likelihood Confidence Intervals				
Effect	Unit	Estimate	95% Confidence Limits	
time night vs day	1.0000	1.296	1.268	1.325
year 2018 vs 2010	1.0000	1.261	1.234	1.290

- The estimated rate of death dying on the road during the day in 2010 is  $\hat{\pi} = \exp(\hat{\beta}_0) / \{1 + \exp(\hat{\beta}_0)\} = 0.000019016$ , so a death rate of 1.9 per 100000 inhabitants. This estimate is slightly higher than the one from the negative binomial model.
- The odds of dying during nighttime (relative to daytime) increase by 29.6%, whereas the odds for 2018 (relative to 2010) increase by 26.1%.