

# Low default modelling: a comparison of techniques based on a real Brazilian corporate portfolio

MSc Guilherme Fernandes and MSc Carlos Rocha

Credit Scoring and Credit Control Conference XII August 2011

---

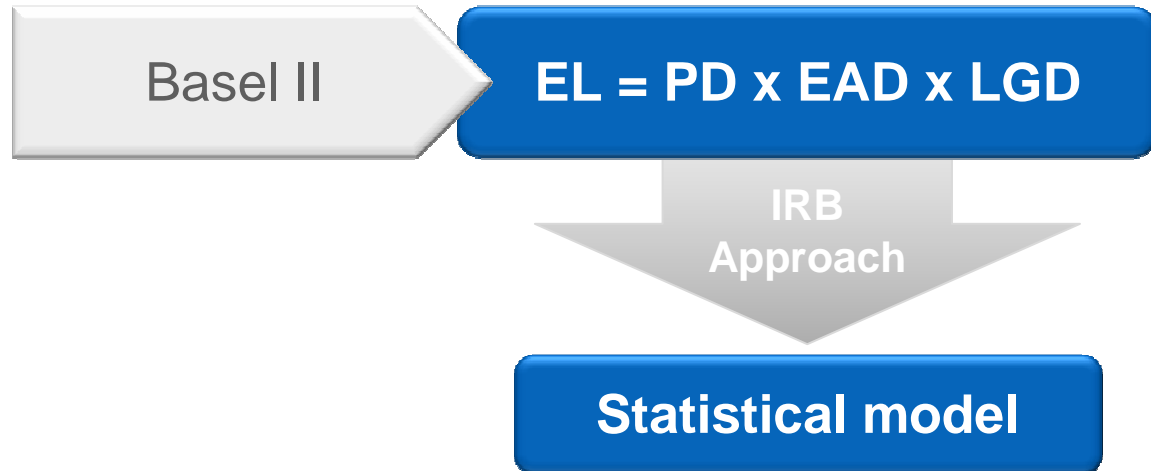


# Topics

- **The problem**
- Modelling techniques:
  - Classical logistic regression
  - Bayesian logistic regression
  - Limited logistic regression
  - Oversampling combined with correction
- Model validation
- Results
- Conclusion



## Previous work



- ❖ Hand and Henley (1997):
  - Statistical classification methods
  - Consumer credit scoring
  - Large data
- ❖ Pluto and Tasche (2006):
  - Low default portfolio
  - PD estimation

# The problem

## Common situation



- ❖ Vehicle financing
- ❖ Retail individuals
- ❖ Risk driver: Percentage of deposit

### Percentage of deposit

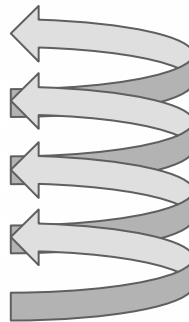
0%  
1%-30%  
31%-50%  
51%-80%  
81%-99%

### Default rate

30%  
20%  
10%  
5%  
1%

### Risk increase

1.5x  
2x  
2x  
5x



## LDP



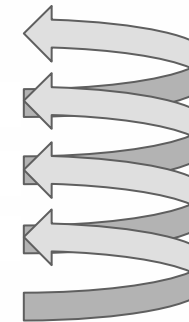
- ❖ Truck financing
- ❖ Middle companies
- ❖ Risk driver: Percentage of deposit

### Default rate

2.5%  
0.9%  
1.1%  
0.3%  
0.0%

### Risk increase

2.9x  
0.75x  
4.7x  
NA



# Topics

- Scenario and the problem
- **Modelling techniques:**
  - **Classical logistic regression**
  - Bayesian logistic regression
  - Limited logistic regression
  - Oversampling combined with correction
- Model validation
- Results
- Conclusion



# Classical logistic regression

**Model statement**

$$Y_i \sim \text{Bernoulli}(p_i)$$

**Target-covariates link**

$$p_i = \left( \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right)$$

**Likelihood function**

$$\ln(L(\beta | y)) = - \sum_{i=1}^n \ln(1 + \exp((1 - 2y_i)x_i' \beta))$$

Where:

$i$  = observation

$n$  = sample size

$y_i$  = response variable (1 : default; 0 : non default)

$x_i'$  = covariate's vector

$\beta$  = parameter vector

$p_i$  = probability of default (i - th observation)

# Topics

- Scenario and the problem
- **Modelling techniques:**
  - Classical logistic regression
  - **Bayesian logistic regression**
  - Limited logistic regression
  - Oversampling combined with correction
- Model validation
- Results
- Conclusion



# Bayesian logistic regression

**Model statement**

$$Y_i \sim \text{Bernoulli}(p_i)$$

**Target-covariates link**

$$p_i = \left( \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right) \quad \beta_k \sim \text{Normal}(\mu_k, \sigma_k^2)$$

**Posterior distribution**

$$\text{posterior}(\beta | y) \propto L(y | \beta) \cdot \text{prior}(\beta)$$

Onde:

$i$  = observation

$n$  = sample size

$y_i$  = response variable (1 : default; 0 : non default)

$\text{prior}(\beta_k) = \text{Normal}(\mu_k, \sigma_k^2)$

$x_i'$  = covariates vector

$\beta$  = parameters vector

$p_i$  = probability of default

$$L(\beta | y) = \exp\left(-\sum_{i=1}^n \ln(1 + \exp((1 - 2y_i)x_i' \beta))\right)$$

# Bayesian logistic regression

Model statement

$$Y_i \sim \text{Bernoulli}(p_i)$$

Target-covariates link

$$p_i = \left( \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right) \quad \beta_k \sim \text{Normal}(\mu_k, \sigma_k^2)$$

Posterior distribution

$$\text{posterior}(\beta | y) = \frac{L(y | \beta) \cdot \text{prior}(\beta)}{\int L(y | \beta) \cdot \text{prior}(\beta) d\beta}$$

Prior information of risk drivers

Data information

PDF of parameters for the model

Solved via Monte Carlo simulation (MCMC)

# Bayesian logistic regression

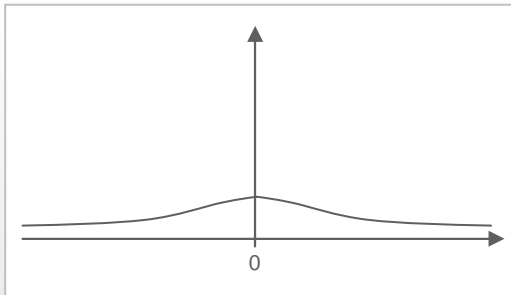
Posterior distribution

$$\text{posterior}(\beta | y) \propto L(y | \beta) \cdot \text{prior}(\beta)$$

Prior distributions

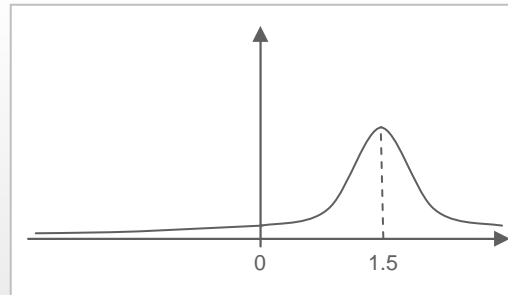
Non Informative

$$\text{prior}(\beta) = \text{Normal}(\mu = 0, \sigma^2 = 1000)$$



Informative

$$\text{prior}(\beta) = \text{Normal}(\mu = 1.5, \sigma^2 = 1.2)$$



Prior distribution

Knowledge about the parameters prior to data information

# Topics

- Scenario and the problem
- **Modelling techniques:**
  - Classical logistic regression
  - Bayesian logistic regression
  - **Limited logistic regression**
  - Oversampling combined with correction
- Model validation
- Results
- Conclusion



# Limited logistic regression

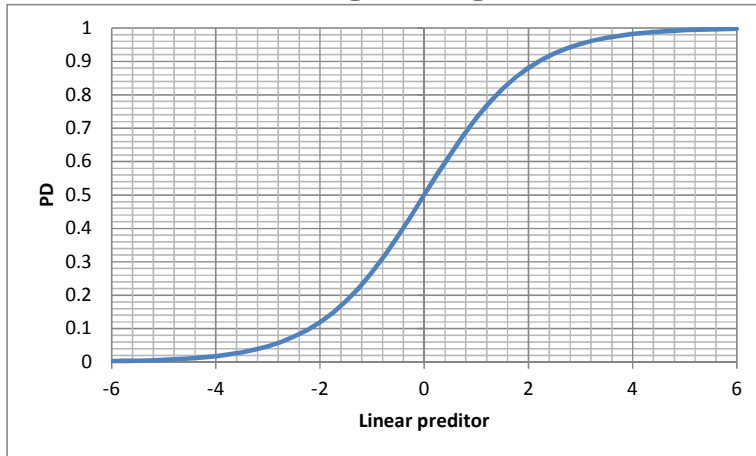
Model statement

$$Y_i \sim \text{Bernoulli}(p_i)$$

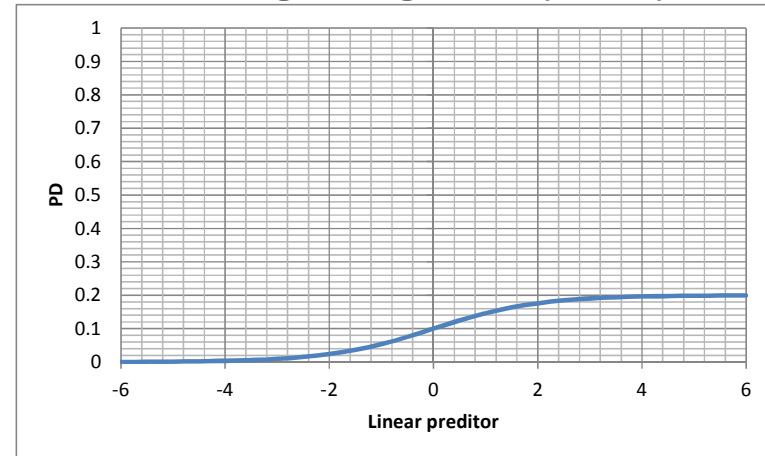
Target-covariates link

$$p_i = \omega \cdot \left( \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right)$$

Classical logistic regression



Limited logistic regression (w = 0.2)



# Limited logistic regression

**Model statement**

$$Y_i \sim \text{Bernoulli}(p_i)$$

**Target-covariates link**

$$p_i = \omega \cdot \left( \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right)$$

**Log-Likelihood function**

$$\ln(L(\beta | y)) = - \sum_{i=1}^n \left( y_i \ln \left( \omega \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right) + (1 - y_i) \ln \left( 1 - \omega \frac{\exp(x_i' \beta)}{1 + \exp(x_i' \beta)} \right) \right) I_{(0,1)}(\omega)$$

Onde:

$i$  = observation

$n$  = sample size

$y_i$  = response variable (1 : default; 0 : non default)

$\omega$  = upper bounding parameter ( $0 < \omega < 1$ )

$x_i'$  = covariates vector

$\beta$  = parameters vector

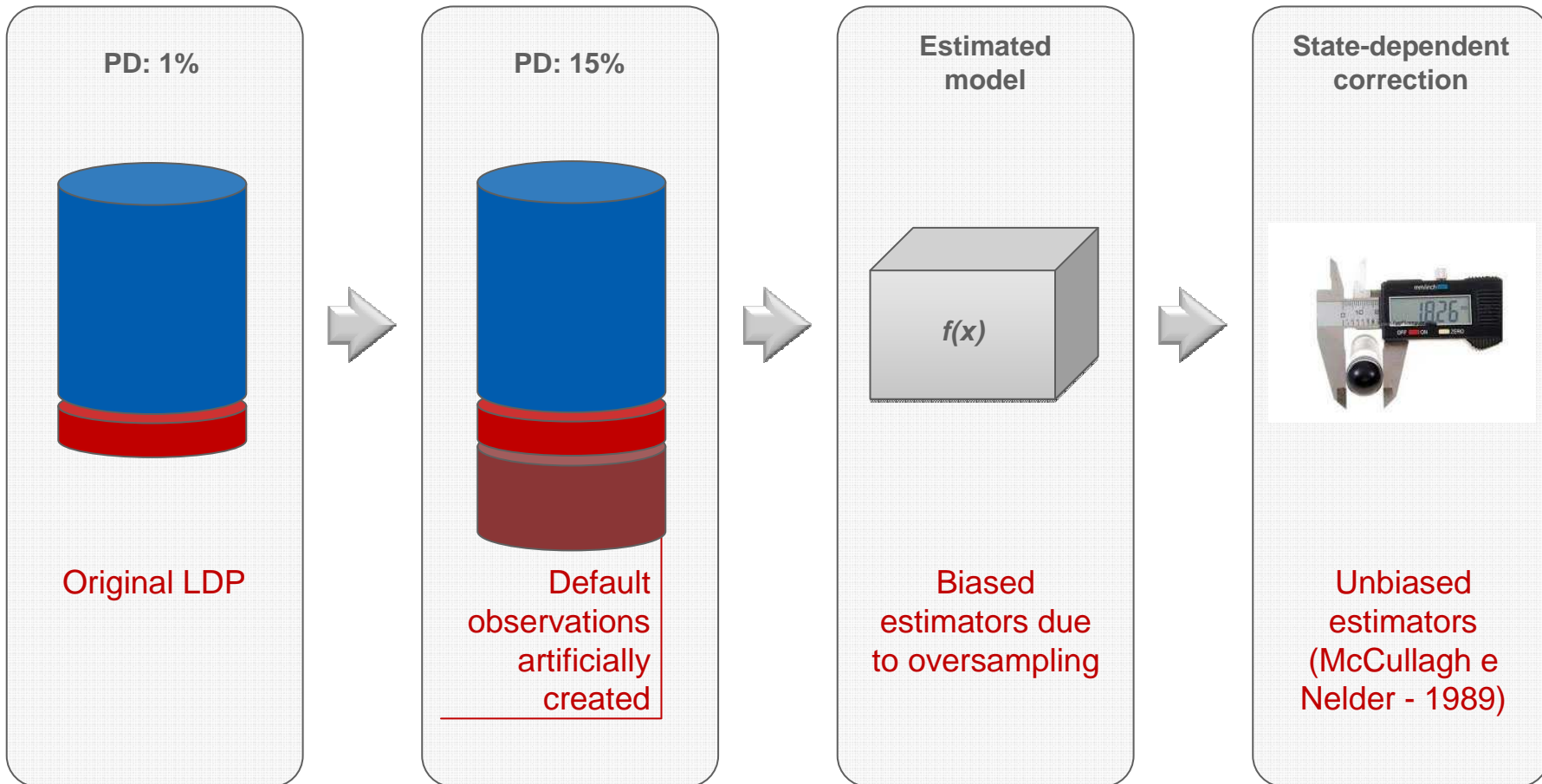
$p_i$  = probability of default

# Topics

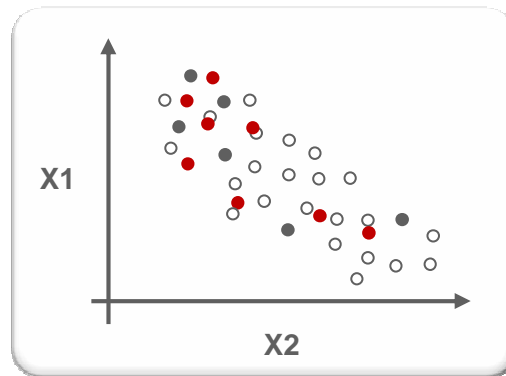
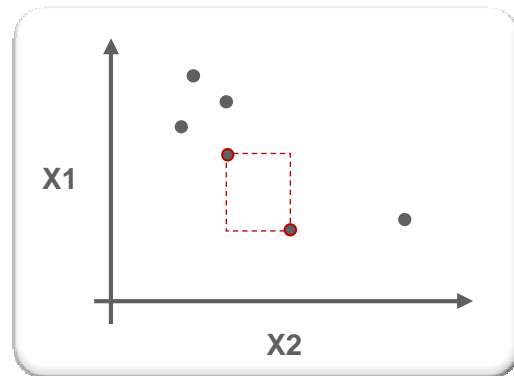
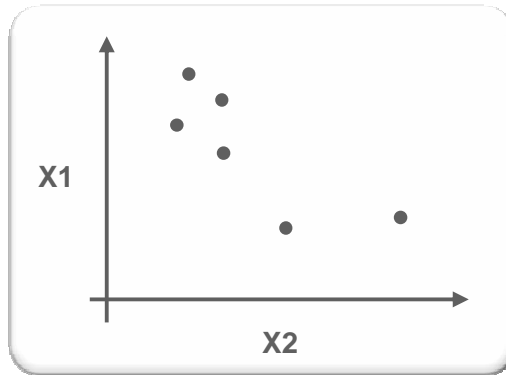
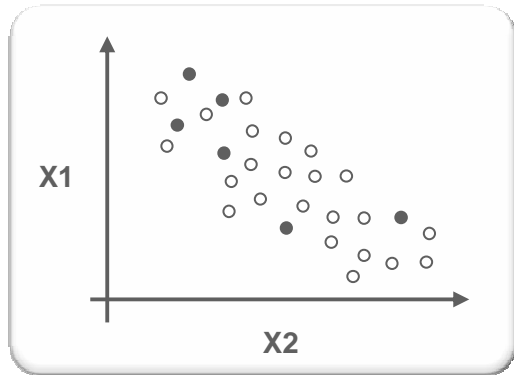
- Scenario and the problem
- **Modelling techniques:**
  - Classical logistic regression
  - Bayesian logistic regression
  - Limited logistic regression
  - **Oversampling combined with correction**
- Model validation
- Results
- Conclusion



# Oversampling and state-dependent correction



# Oversampling technique: SMOTE\*



➤ Artificial observations

1. Select defaulted observations
2. Randomly choose 2 observations
3. Randomize one point within the space defined
4. Estimate model with the inflated default rate database

# State-dependent sample

## Weighted Log-likelihood function

$$\begin{aligned} \ln(L_w(\beta | y)) &= \omega_1 \sum_{[Y_i=1]} \ln(p_i) + \omega_0 \sum_{[Y_i=0]} \ln(1 - p_i) \\ &= - \sum_{i=1}^n \omega_i \ln(1 + \exp((1 - 2y_i)x_i' \beta)) \end{aligned}$$

Parameters ( $\beta$ ) estimated by WMLE are biased, even on large data. McCullagh e Nelder (1989) present the correction so the model will be unbiased.

Onde:  $i$  = observation  
 $n$  = sample size  
 $x_i'$  = covariates vector

$p_i$  = probability of default  
 $\omega_i$  = weight of  $i^{\text{th}}$  observation =  $\omega_1 Y_i + \omega_0 (1 - Y_i)$   
 $\tau$  = default rate in population  
 $y_i$  = response variable (1: default; 0: non default)

$\beta$  = parameters vector  
 $\omega_1 = \tau / \bar{y}$   
 $\omega_0 = (1 - \tau) / (1 - \bar{y})$   
 $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$

$$\begin{aligned} \xi_i &= 0.5 Q_{ii} [(1 + \omega_1) \pi_i - \omega_1] \\ Q_{ii} &= \text{elemento}_{ii} [X(X'WX)^{-1}X'] \end{aligned}$$

$$\text{viés}(\hat{\beta}) = (X'WX)^{-1} X'W\xi_i$$

$$\tilde{\beta} = \hat{\beta} - \text{viés}(\hat{\beta})$$

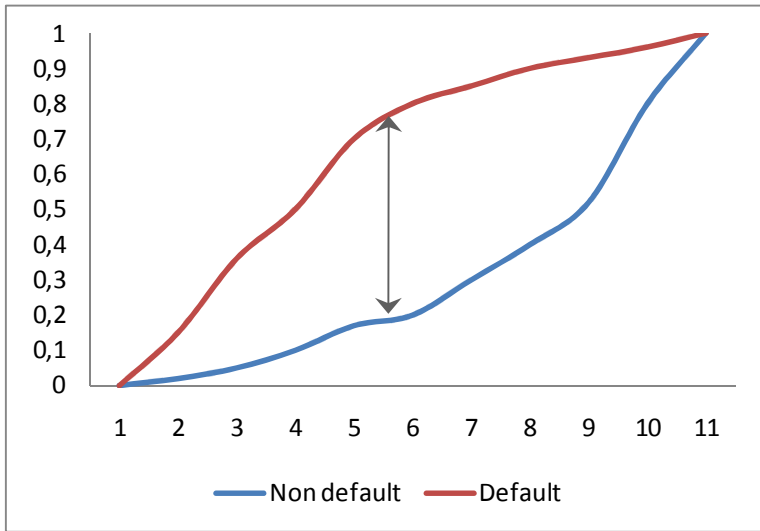
# Topics

- Scenario and the problem
- Modelling techniques:
  - Classical logistic regression
  - Bayesian logistic regression
  - Limited logistic regression
  - Oversampling combined with correction
- **Model validation**
- Results
- Conclusion



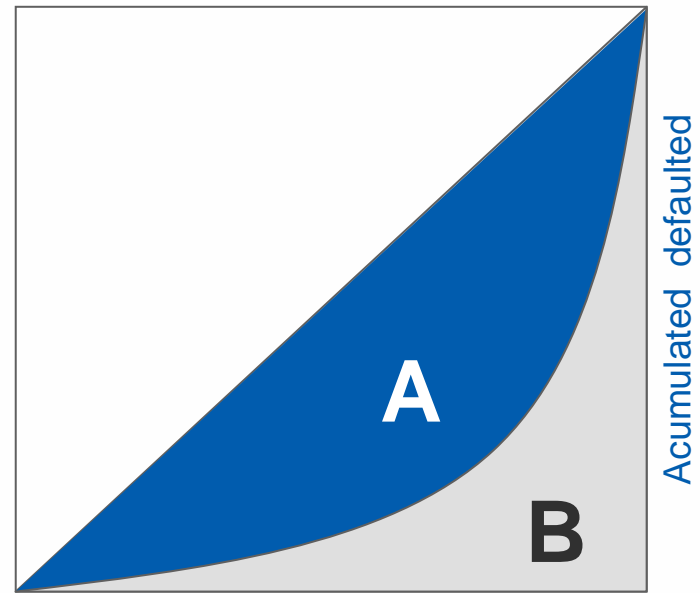
# Model validation

## KS



$$KS = \text{Max}( |NDAcum - Dacum| )$$

## Gini



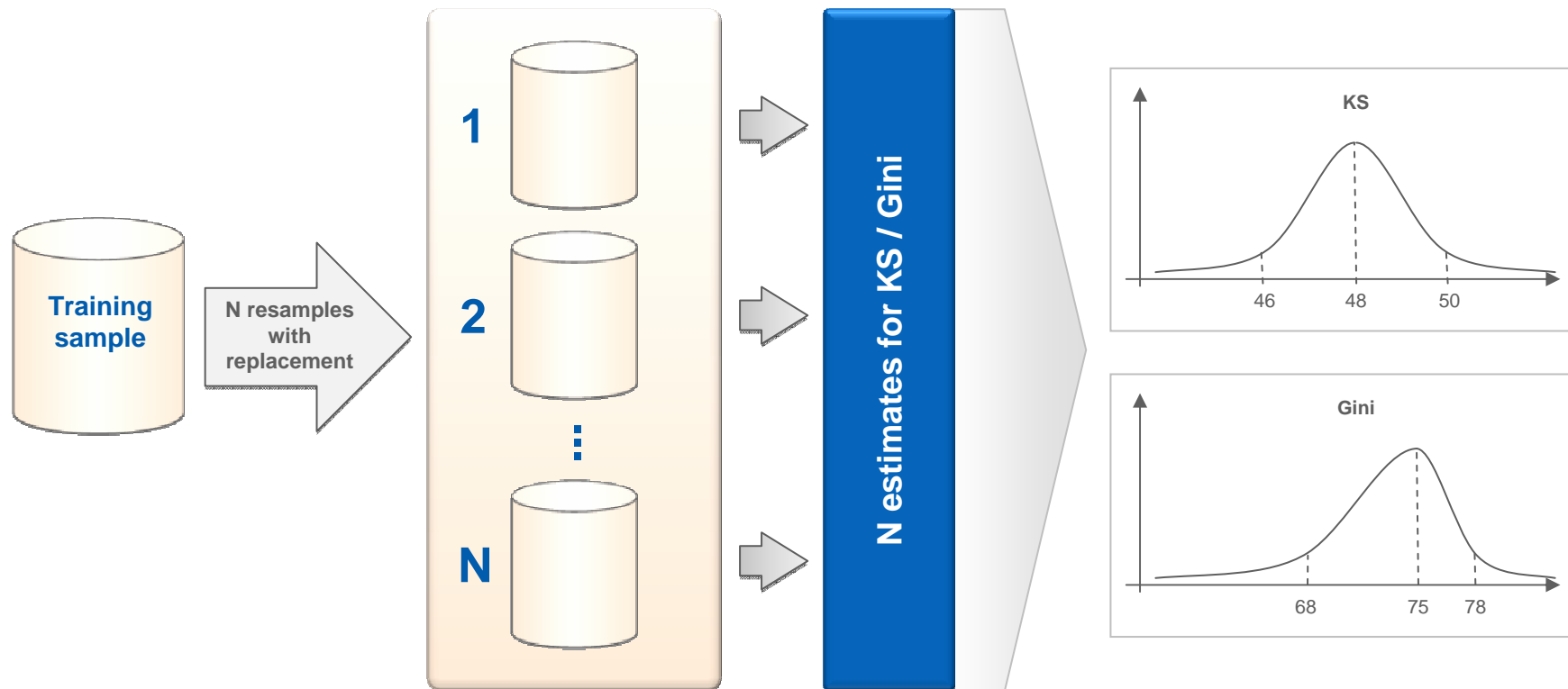
Acumulated non-defaulted

Acumulated defaulted

$$Gini = A / (A+B)$$

# Model validation

❖ Alternative method for out-of-time sample



# Topics

- Scenario and the problem
- Modelling techniques:
  - Classical logistic regression
  - Bayesian logistic regression
  - Limited logistic regression
  - Oversampling combined with correction
- Model validation
- **Results**
- Conclusion

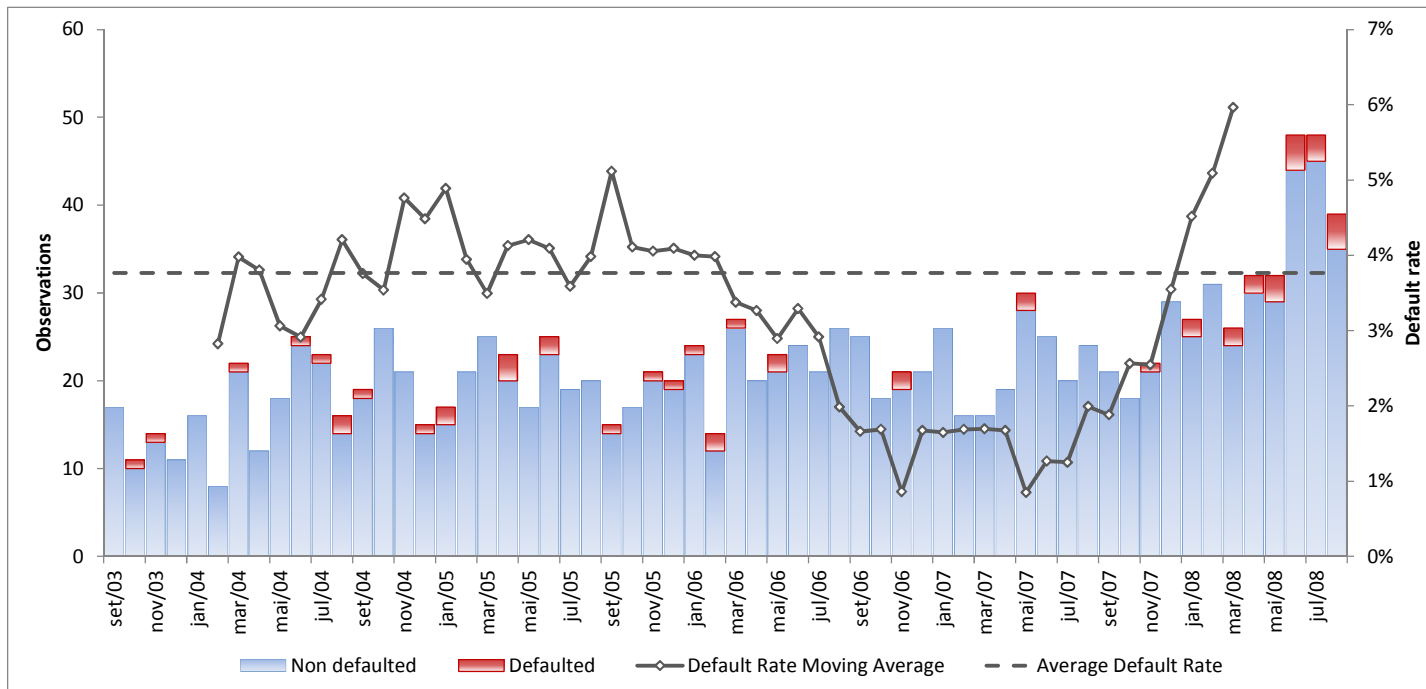


# Results

Data set

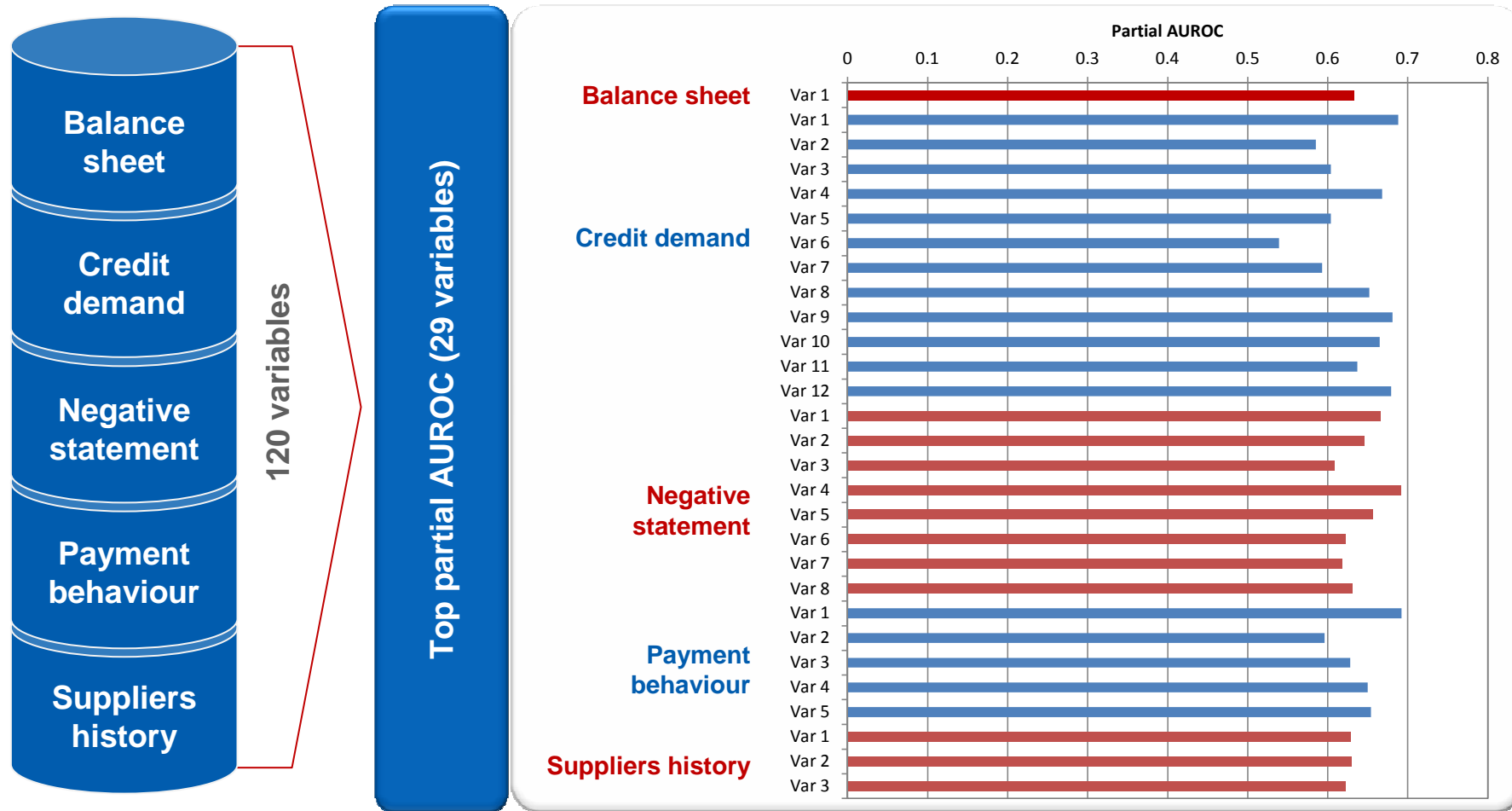


- Revenue > US\$ 130 per year
- 1.327 obs
- Period: 2003 – 2008
- 50 defaulted observation



# Results

## Variables available and selection



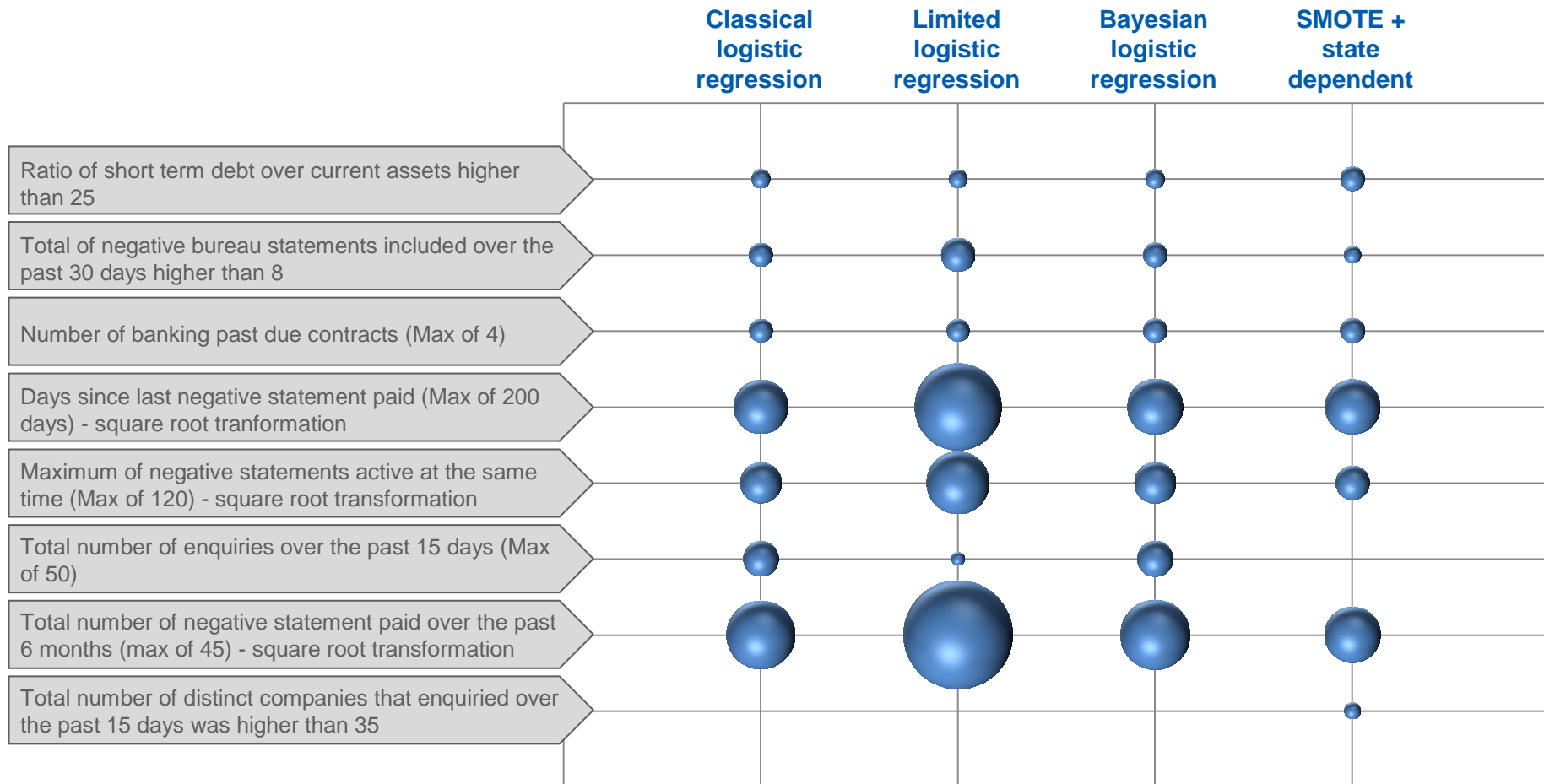
# Results

## Parameters per model

Variable	Domain	Classical model	Limited logistic	Bayesian model	SMOTE + state dep
Intercept		-3.62280		-3.69740	-3.13471
w (upper bound limit)			0.14641		
Ratio of short term debt over current assets higher than 25	(0;1); Dummy for values higher than 25	0.84510	0.96934	0.86510	1.09750
Total of negative bureau statements included over the past 30 days higher than 8	(0;1); Dummy for values higher than 8	1.07010	2.12650	1.09010	0.77100
Total number of distinct companies that enquired over the past 15 days was higher than 35	(0;1); Dummy for values higher than 35				0.73950
Number of banking past due contracts (Max of 4)	Integer; limited in 4	0.26380	0.30292	0.26860	0.27870
Days since last negative statement paid (Max of 200 days) - square root tranformation	Real; limited in 14.14	-0.17440	-0.33405	-0.17930	-0.17714
Maximum of negative statements active at the same time (Max of 120) - square root transformation	Real; limited in 10.95	0.17130	0.33579	0.17270	0.14201
Total number of enquiries over the past 15 days (Max of 50)	Integer; limited in 50	0.03200	0.01352	0.03240	
Total number of negative statement paid over the past 6 months (max of 45) - square root transformation	Real; limited in 6.70	-0.46370	-0.86208	-0.47180	-0.37800

# Results

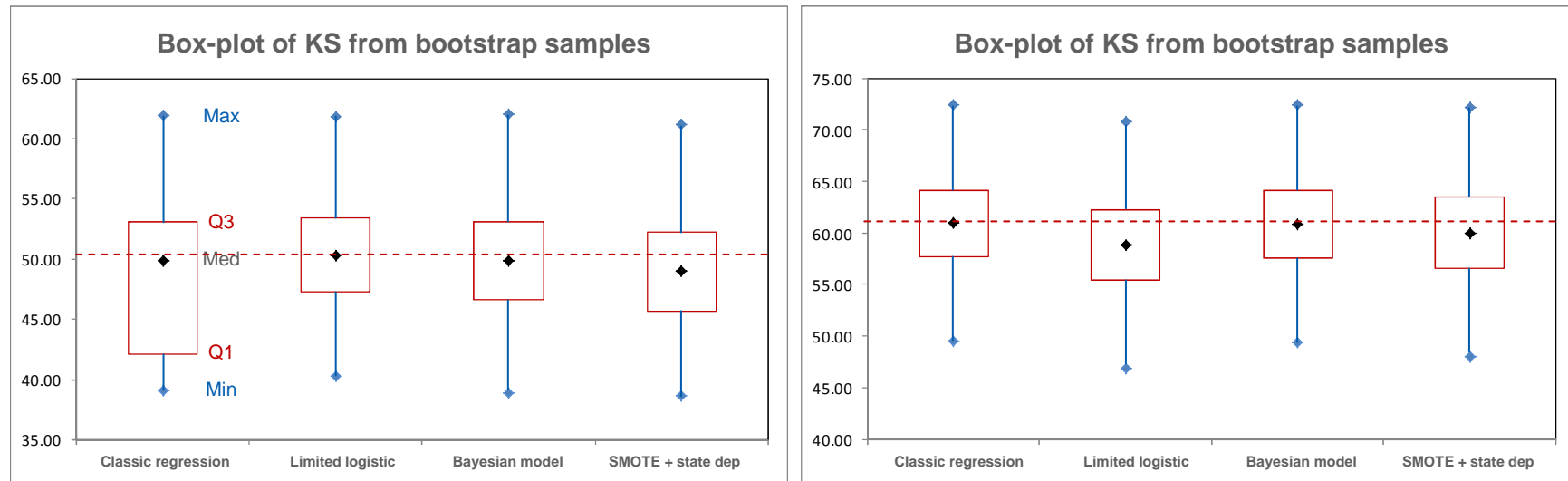
## Importance of variables



# Results

## Performance measures

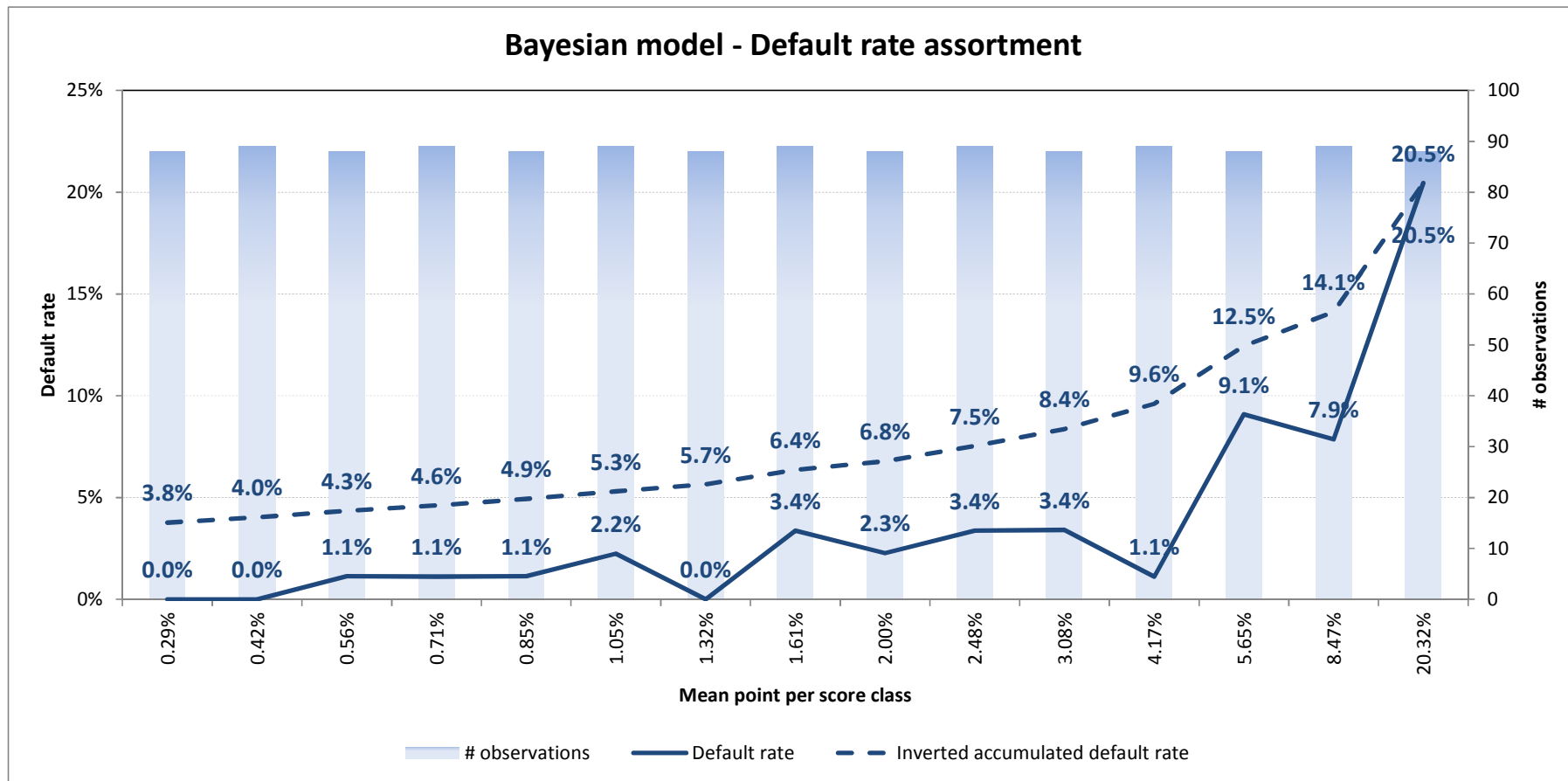
10.000 bootstrap samples



**Bayesian model: High performance and low variability**

# Results

## Default rate assortment



# Conclusion

---

- Pluto & Tasche: estimates PD for a pre-existing rating grade
- Approaches presented:
  - Limited logistic regression: Best KS, worst Gini;
  - SMOTE + State dependent: incorporates different variable;
  - Classic logistic regression: High KS and Gini, but higher variability on bootstrap;
  - Bayesian logistic regression: High KS and Gini.
- Bayesian model gives reasonably assorted default rate even on LDP
- Future research: how to incorporate informative prior using specialist's knowledge and bureau information

**Guilherme B Fernandes**

*guilherme.fernandes@br.experian.com*

**Carlos A. Rocha**

*carlos.rocha@br.experian.com*

## Questions

---

Serasa  Experian

