

# MINI PROJETO 1

Modelos probabilísticos e Dados

# IDENTIFICAÇÃO DE DISTRIBUIÇÕES

#### **OBJETIVO:**

O objetivo deste projeto é identificar quais distribuições (funções de densidade de probabilidade - no caso contínuo, ou funções de probabilidade - no caso discreto) descrevem bem variáveis quantitativas extraídas de um *datasets*.

O resultado final esperado é um relatório que identifique, com bons argumentos, a escolha de um ou mais modelos probabilísticos para ajuste de uma variável quantitativa extraída de *dataset*.

Este projeto é estritamente individual.

#### **O QUE DEVE SER FEITO:**

Você precisa escolher uma variável quantitativa em *datasets* públicos de sua escolha. A sua variável pode ser discreta ou contínua.

Limpe e prepare os dados para processamento (tratando valores NaN ou N/A, por exemplo). Fique atento ao dicionário de dados (se houver) para identificar quais colunas do *dataset* de fato são quantitativas e eventualmente remover valores inválidos.

A seguir, estude a variável escolhida e procure identificar uma função adequada que descreva as probabilidades de ocorrência dos valores que essa variável pode assumir.

Sugerimos fortemente que o trabalho siga as seguintes fases:

- 1. Seleção de um dataset e escolha uma variável quantitativa adequada.
  - Não há restrições em relação à base de dados a utilizar, desde que não seja as mesmas bases da PNAD já usadas na disciplina. Aconselha-se evitar variáveis de bases com pequeno tamanho amostral.
  - Tornamos disponível uma <u>Lista de datasets</u> que pode ajudar nesta fase do trabalho.
    Atenção: nem todas as bases de dados desta lista têm variáveis quantitativas, analise com cuidado. Você não precisa ficar restrito a esta lista
  - Indique o dataset e a variável que escolheu no piazza. IMPORTANTE!!!



- 2. Limpeza da variável escolhida, se necessário.
- 3. Inspeção visual da distribuição dos valores da variável escolhida usando um histograma, por exemplo.
- 4. Formulação de hipóteses sobre o formato da distribuição dos dados (simetria, assimetria positiva e assimetria negativa) e escolha PELO MENOS DUAS DISTRIBUIÇÕES TEÓRICAS DIFERENTES PARA MODELAR SUA ÚNICA VARIÁVEL QUANTITATIVA definida no item 1. Justifique por que escolheu suas distribuições teóricas.
- 5. Tentativa de estimar os parâmetros da família de distribuições escolhida no item acima a partir dos dados.
- 6. As distribuições do pacote scipy.stats têm uma função chamada fit () que procura estimar os parâmetros a partir do conjunto de dados.
  - Use o fit() para fazer estimativa dos parâmetros da família de distribuições escolhida no item 4.
  - Compare os parâmetros estimados a partir do fit() com os parâmetros estimados por você no item 5. Para cada uma das suas distribuições teóricas, opte por um ajuste: o do item 5 ou o obtido pelo comando fit().
- 7. Construa o histograma dos dados junto com a fdp de cada distribuição teórica e analise.
- 8. Construa o QQ-Plot (quantil amostral vs quantil teórico) e analise. **Dica:** veja Exemplo 6.8 do Magalhães e Lima (7ª. edição) de como obter as frequências relativas acumuladas a partir de uma amostra de tamanho n e de como obter os quantis teóricos.
- Construa um gráfico com a frequência relativa acumulada (a partir dos dados) vs a função de distribuição acumulada e analise.
- 10. Faça um teste de aderência para a distribuição (veja o arquivo MiniProjetol Aderencia Numpy Pseudocodigo.ipynb). Teste de aderência é útil para mensurar a qualidade do ajuste do modelo teórico aos dados.
- 11. Elabore uma tabela que contrasta sua variável com as distribuições teóricas escolhidas e a qualidade do ajuste em cada caso. Analise essa tabela e faça a escolha da melhor das distribuições teóricas para o ajuste dos dados.



### **ENTREGÁVEIS ESPERADOS E DATAS:**

#### Turmas A, B e C:

Item	Data	Descrição		
Indicação de dataset	19/09/2016	Indicar <i>dataset</i> e variável de interesse em post no <b>Piazza</b> de Ciência dos dados		
Entrega intermediária (check)	20/09/2016	Histogramas das variáveis candidatas e possíveis distribuições adequadas (Itens 1 a 4 completos).		
Relatório final	23/09/2016	Relatório enviado na pasta MiniProjeto1 no Github.		

#### **FÓRUM DE DISCUSSÃO:**

Um fórum de discussão foi criado no <u>Piazza</u> - procure participar para tirar suas dúvidas e ajudar seus colegas: (<a href="https://piazza.com/insper.edu.br/fall2016/cd2016">https://piazza.com/insper.edu.br/fall2016/cd2016</a> 2)

Não aceitaremos mesma variável quantitativa analisada por dois alunos da mesma sala ou de turmas diferentes. Assim, aproveite o fórum para descrever as variáveis quantitativas que irá trabalhar. Uma vez publicadas, um outro aluno não poderá mais utilizá-las neste projeto. Esse fórum será único para as três turmas, fazendo com que isso seja válido para todas as turmas.

# Engenharia Ciência dos Dados

# Insper

## RUBRICS DE AVALIAÇÃO DO OBJETIVO DE APRENDIZADO

Objetivo de aprendizado	Insatisfatório (I)	Em desenvolvimento (D)	Essencial (C)	Proficiente (B)	Avançado (A)
Especificar as	_ ·	Conseguiu fazer a leitura	Para a variável	Realizou os	Realizou os comportamentos
distribuições de	insuficientes ou	dos dados mas não	quantitativa escolhida:	comportamentos de C de	de B e C de maneira excelente
probabilidades 	atrasadas	avançou na análise		maneira excelente e:	e:
adequadas para as					
variáveis			- Leu os dados		
		Escolheu um conjunto de	adequadamente	Traçou adequadamente	- Avaliou entre pelo menos
		dados que já tinha sido		as fdp's ou fp's junto aos	duas distribuições alternativas
		escolhido pelo colega	- Traçou um histograma	histogramas.	usando um teste de aderência
			considerando densidade		e formulou uma conclusão
			no eixo y (normed=True).	- Construiu e avaliou	coerente em relação a escolha
				adequadamente o ajuste	da distribuição que gera o
		Não indicou	- Elegeu pelo menos duas	dos dados usando QQ-plot.	melhor ajuste.
		adequadamente a URL	distribuições teóricas e		
		do dataset escolhido ou	justificou escolhas.	- Construiu e avaliou	
		os nomes específicos das		adequadamente função de	
		variáveis	- Estimou os parâmetros	distribuição acumulada e	
			das distribuições teóricas a	frequência relativa	
			partir dos dados ou	acumulada.	
			utilizou comando fit(), mas	Analiaau anéficaa	
			justificou no item 6	- Analisou gráficos	
			escolha final para	impecavelmente.	
			estimativas.	(Até item 9!)	
			(Até item 6!)	price reciti 5.7	