

Projetão da Disciplina

Germano C. Vasconcelos
Centro de Informática - UFPE

Objetivo



Realizar um projeto com base de dados reais em larga escala com modelos de redes neurais e outros classificadores



Motivações



- Possibilitar uma visão prática do uso de redes neurais na solução de problemas
- Consolidar os conhecimentos teóricos apresentados em sala de aula
- Permitir o contato com ferramentas do Github, Keras, Scikit-learn na Linguagem Python



- Classificação binária (2 classes)
 - Base real do mercado
 - Em larga escala: ~ 390 mil registros para treinamento, validação e teste
 - Problema: com base no perfil do cliente, decidir a quem conceder crédito (risco de inadimplência)

Descrição do Projeto



- Conjunto de classificadores disponíveis
 - Perceptron multicamadas (MLP) (obrigatório)
 - Máquina de Vetores de Suporte (opcional rodar 1 configuração)
 - Ensemble de MLPs (obrigatório)
 - Random Forest (usado para comparação)
 - Gradient Boosting (usado para comparação)
 - Ensemble de Classificadores (usado para comparação)
- Investigar diferentes topologias da rede e diferentes valores dos parâmetros (básico)
 - Número de camadas
 - Número de unidades intermediárias
 - Variação da taxa de aprendizagem
 - Função de ativação (logística, tangent hiperbolica, Relu)
 - Usar método de amostragem básica (repetitive oversampling)



Descrição do Projeto



- Parâmetros adicionais que podem ser explorados
 - Algoritmo de aprendizagem
 - Taxa de aprendizagem adaptativa
 - Outros



Preparação de Dados: (divisão e balanceamento)



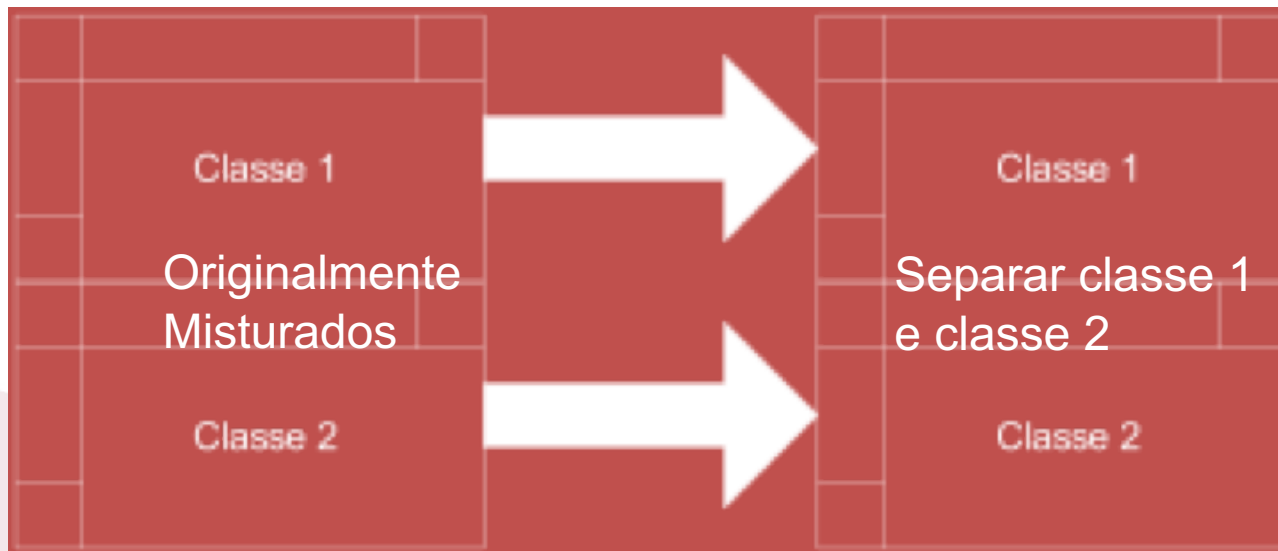
- Conjuntos de dados independentes
 - Treinamento (já está separado)
 - Validação (já está separado)
 - Teste (já está separado)
- Estatisticamente representativos e independentes
 - Não pode haver sobreposição



Preparação de Dados: (divisão e balanceamento)

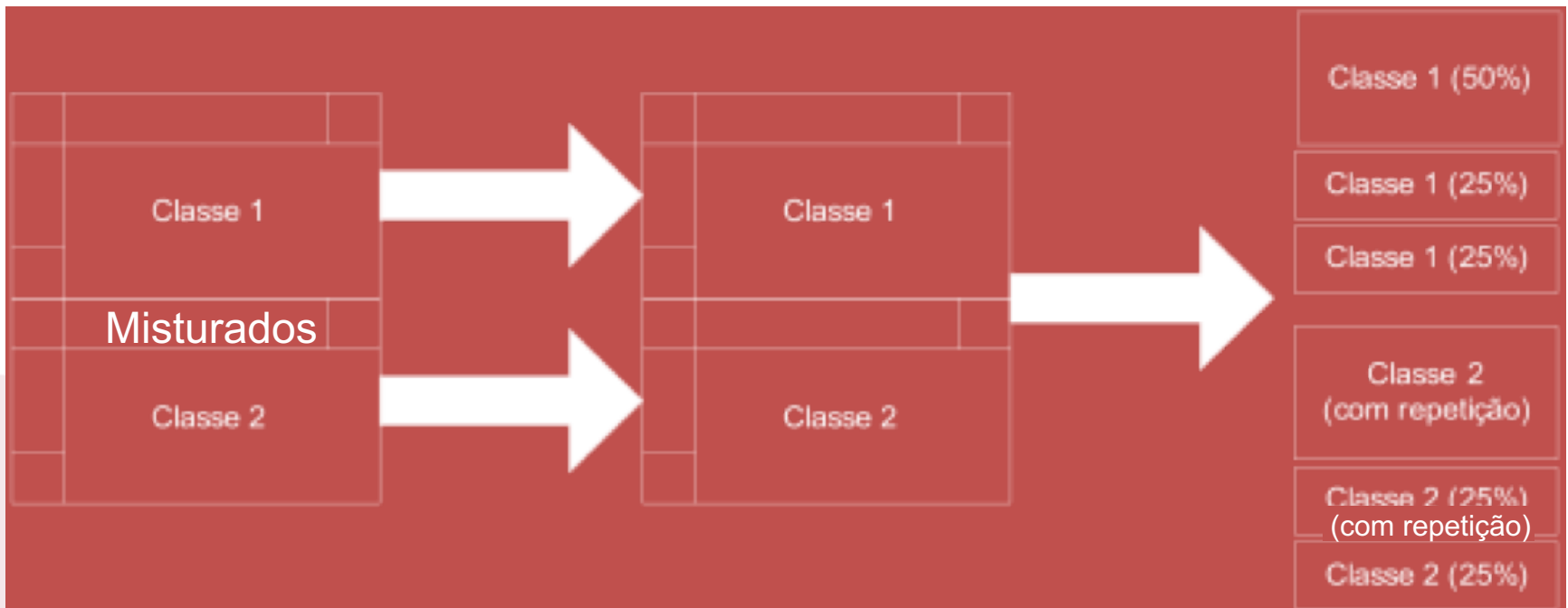


Particionamento dos Dados – Primeira etapa



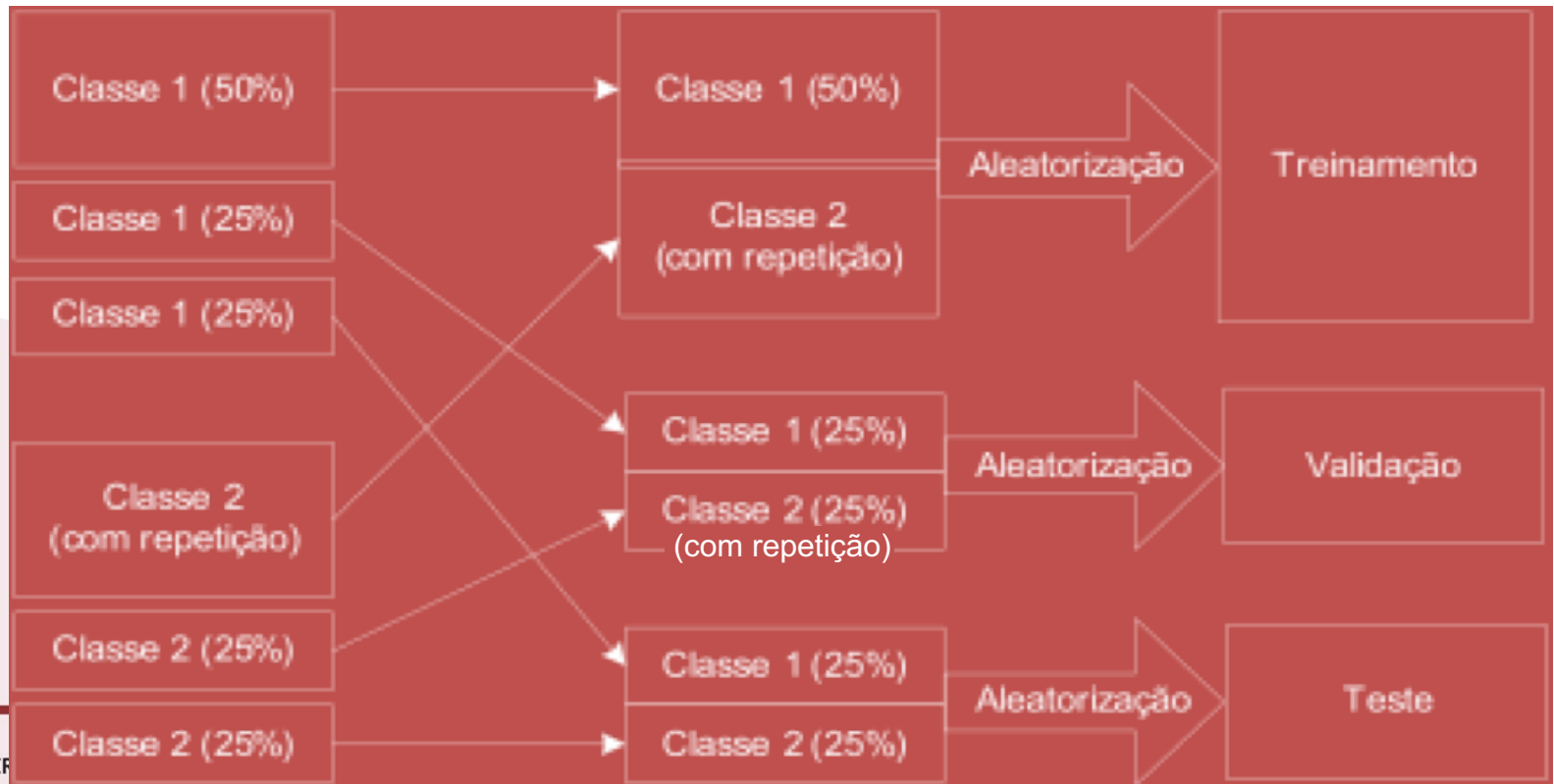
Preparação de Dados: (divisão e balanceamento)

Particionamento dos Dados – Segunda etapa



Preparação de Dados: (divisão e balanceamento)

Particionamento dos Dados – Terceira etapa



■ Classificação

- Teste estatístico Kolmogorov-Smirnov -KS (**principal**)
- MSE (erro médio quadrado)
- Matriz de confusão
- Auroc (Área sob a Curva Roc)
- Recall, Precision e F-Measure

Experimentos



- Base já processada
- Importante
 - Registrar o desempenho de forma evolutiva, a cada etapa



Experimentos



- Sugestão:

- Iniciar com um modelo MLP e um modelo Random Forest
- Após bom desempenho com esses modelos, outros classificadores podem ser investigados
- ensembles de MLPs,
- (opcional) gradient boosting, ensembles mistos, votação, meta-classificadores). (Opcional)
- Esquemas de ensembles: Voting classifier, Bagging classifier



Experimentos



Parametros que podem ser variados: MLPs

- # camadas (1 ou 2)
- # neurônios (iniciar pequeno e aumentar na necessidade)
- Taxa de aprendizagem
- Função de ativação (logística, tangent, Relu)
- Otimizadores (adadelta, adam, RMs prop, SGD)
- Drop out
- Regularização
- # Epocas: 10.000 (parar aprendizagem pelo overfitting)
- Patience (Max fail): 10 (se parando ainda precoce aumentar para 20)



Experimentos

Parametros sugestivos: Random Forest

- # estimadores
- Max depth
- Max features
- Min_sample_leaf

Experimentos

Parametros sugestivos: Gradient Boosting, Xgboost

- Loss: deviance
- Learning rate
- # estimators
- Subsample
- Criterion: Friedman_mse
- Min_samples_leaf
- Max depth

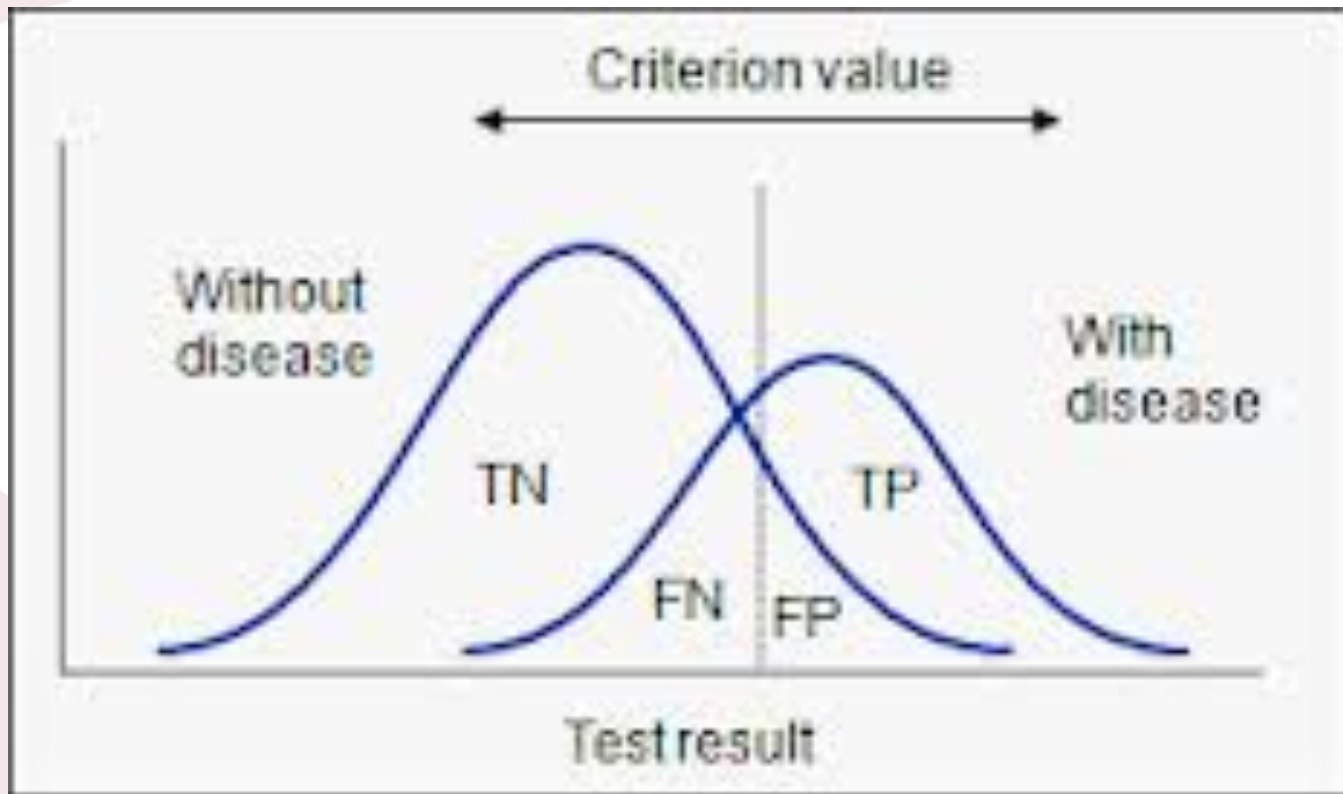
Ferramentas Úteis



- Google Colab
- Keras
- Scikit Learn
- Pandas + Imbalanced learn
- Optuna: variação de parâmetros



Avaliação (Desempenho e Resultados)



Avaliação (Desempenho e Resultados)

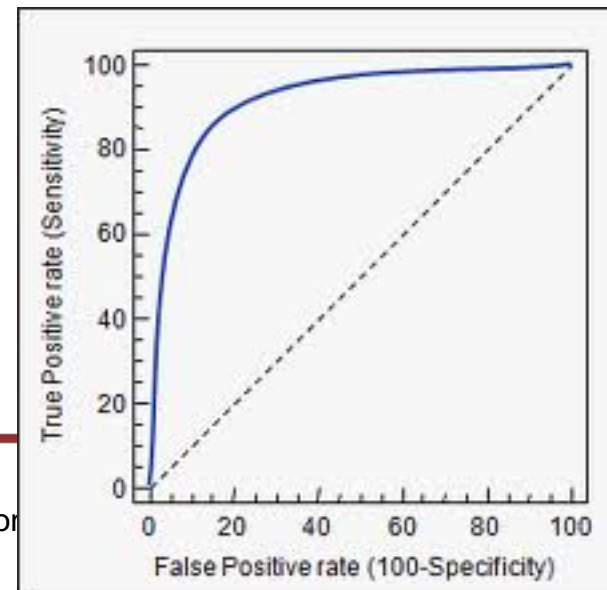
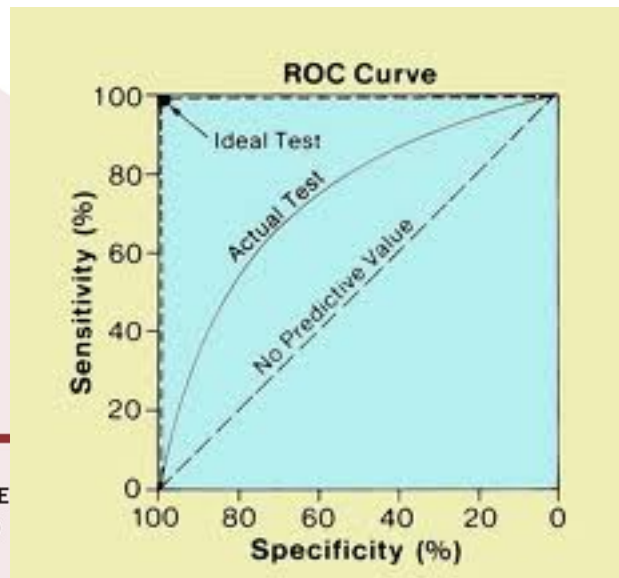
Matriz de Confusão

		Actual classification	
		positive	negative
Hypothesis	positive	true positive (tp)	false positive (fp)
	negative	false negative (fn)	true negative (tn)

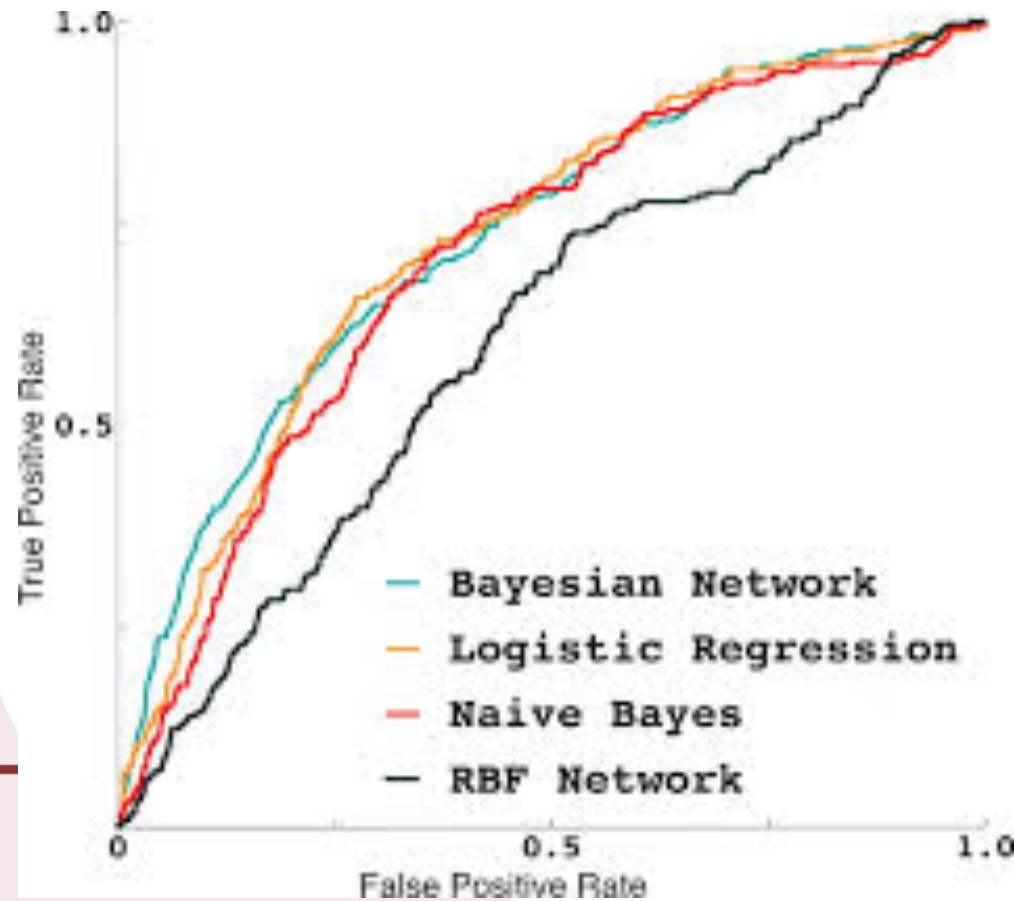
Avaliação (Desempenho e Resultados)

		condition		
		+	-	
test outcome	+	True positive	False positive (type I error, p-value)	→ positive predictive value
	-	False negative (type II error)	True negative	→ negative predictive value
		↓	↓	
		sensitivity	specificity	

Curvas ROC



Curvas ROC: Exemplo



Ferramentas para o Projeto



- Código em Python
 - <https://github.com/RomeroBarata/IF702-redes-neurais>
- Conjuntos de dados do problema
 - <http://www.cin.ufpe.br/~gcv/web Ici/intro.html>



Resultados do Projeto



- Apresentação com todos do grupo com descrição do problema, divisão dos dados, estrutura experimental e interpretação dos resultados
- Entrega no final do semestre

