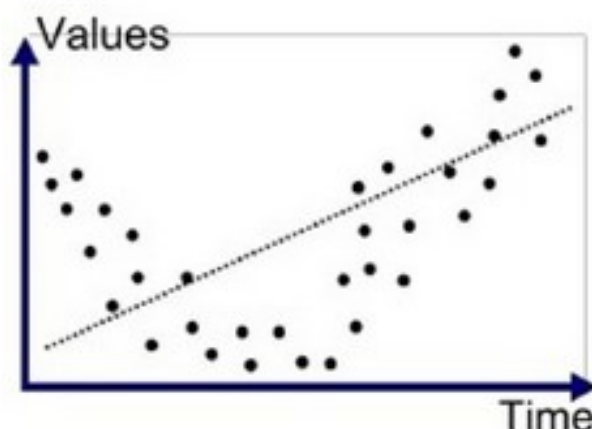


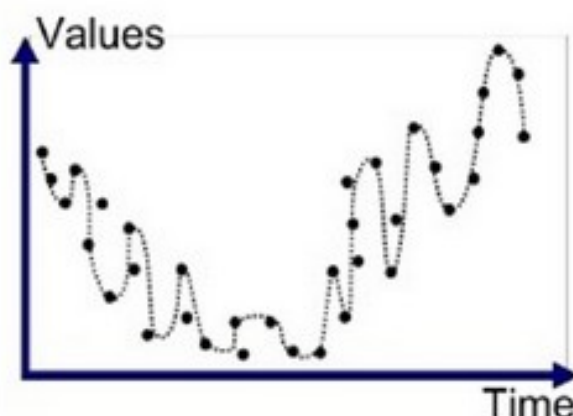
PIBIC - atividade 4

Overfitting e Underfitting

Overfitting e underfitting são dois conceitos que podem ser considerados “erros” na área de machine learning, ambos ocorrem na fase de treinamento e podem comprometer os resultados finais da predição. Como podemos ver pelos nomes, underfitting é quando nossa rede neural não consegue se adaptar aos dados de treinamento, consequentemente errando mais nos dados de teste, geralmente um gráfico de underfitting é assim:

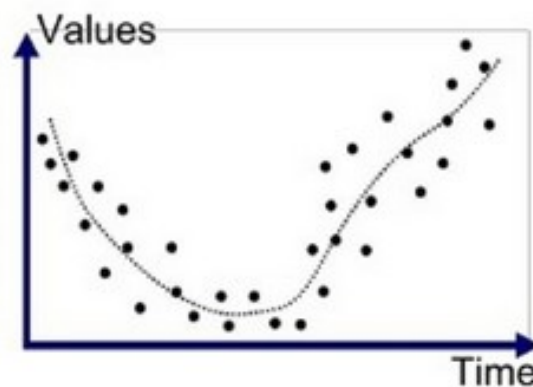


Perceba que a predição ocorreu de maneira linear, com uma reta, e não conseguiu se adaptar nem um pouco aos dados de treinamento, podendo ser considerada um exemplo de “underfitting”. O overfitting é exatamente o contrário, ele ocorre quando o modelo se adapta até demais aos nossos dados de treinamento, como podemos ver no gráfico abaixo:



O modelo pode parecer muito bem treinado a primeira vista, e ele está! O problema é que o nosso modelo se adaptou bem mais do que deveria, e se nós adicionarmos um dado que se distancie minimamente dos dados de treinamento, a precisão de nossa predição será mínima.

Estes são dois gráficos de modelos que deram errado, um modelo ideal de rede neural precisa ser capaz de generalizar os dados de treino para poder encaixar os dados de teste, então seu gráfico seria assim:



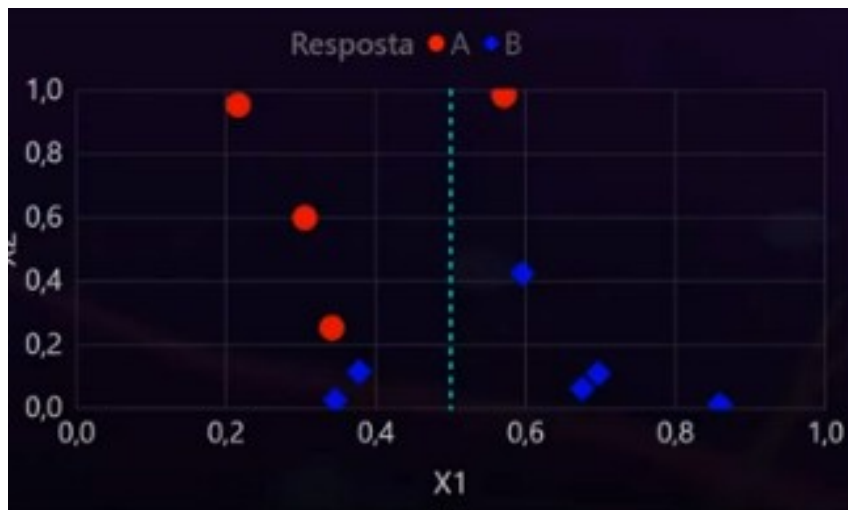
Desta maneira, quando encaixarmos um novo dado, a predição de nossa rede neural seria mais generalizada, logo sendo mais precisa quando comparado com os dois modelos passados.

Acurácia, precisão e recall (f1, precision and recall)

Estes 3 conceitos são o que usamos para poder ver quão boa é a nossa rede neural, para podermos entender estes 3 conceitos precisamos saber o que são os: falsos positivos (FP), falsos negativos (FN), verdadeiros positivos (TP) e verdadeiros negativos (TN). Primeiramente é importante deixar claro que quando uso a nomenclatura “positivo” e “negativo”, estou me referindo a afirmação que diz se o dado é ou não da classe que estamos buscando.

		Valor Previsto	
		Positivo	Negativo
Valor Verdadeiro	Positivo	TP Verdadeiro Positivo	FN Falso Negativo
	Negativo	FP Falso Positivo	TN Verdadeiro Negativo

Esta é uma clássica tabela de falsos positivos e negativos, ela é muito utilizada para podermos organizar as predições de nosso modelo, vou utilizar o seguinte modelo para exemplificar:



Vamos tratar como positivos os pontos vermelhos (A), e como negativos os pontos azuis (B), neste modelo nós possuímos: 2FP, 1FN, 3TP, 4TN, e é com estes valores que nós iremos classificar o nosso modelo como sendo “bom” ou “ruim”.

Vamos começar com a precisão, ela mede a proporção de positivos classificados corretamente, podemos obtê-la com a seguinte fórmula:

$$\frac{\sum \text{True Positive}}{\sum \text{True Positive} + \sum \text{False Positive}}$$

De acordo com essa fórmula, a precisão do nosso modelo seria: $3/(3+2) = 0.6$ ou 60% de precisão. Agora, vamos medir a acurácia de nosso modelo, que é apenas a proporção entre dados corretos e o total dos dados analisados pelo nosso modelo, ela é definida pela seguinte fórmula:

$$\frac{\sum \text{True Positive} + \sum \text{True Negative}}{\text{Test Data Size}}$$

Dito isto, a acurácia de nosso modelo-exemplo seria: $(3+4)/10 = 0.7$ ou 70% de acurácia. E por fim, vamos ver o recall, que mede a proporção de positivos identificados corretamente, ou seja, o quão bom a nossa rede neural é em detectar positivos, sua fórmula é a seguinte:

$$\frac{\sum \text{True Positive}}{\sum \text{True Positive} + \sum \text{False Negative}}$$

O recall do nosso modelo é igual a: $3/(3+1) = 0.75$ ou 75% de recall.

Método holdout