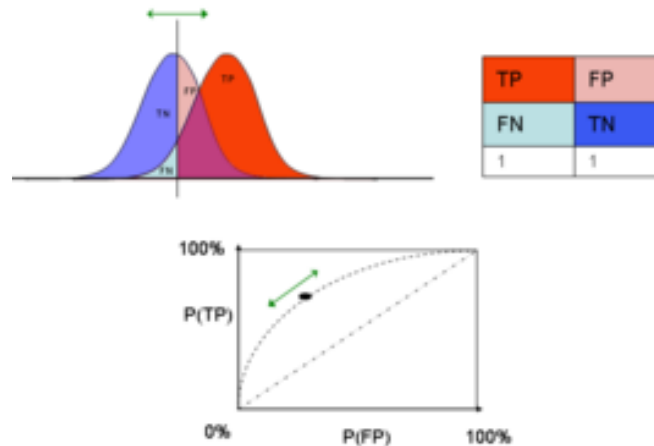


# Análise e Comparação de Desempenho

André E. Lazzaretti

UTFPR/CPGEI



# Matriz de Confusão

- Supondo que você esteja trabalhando com um problema de classificação de duas classes e tem em mãos  $\mathbf{y}_{\text{pred}}$  e  $\mathbf{y}_{\text{true}}$ :

		Predicted class	
		$P$	$N$
Actual Class	$P$	True Positives (TP)	False Negatives (FN)
	$N$	False Positives (FP)	True Negatives (TN)

true label	0	1
0	71	1
1	2	40
predicted label		

Erro, acurácia, taxas TP e FP,  
*precision, recall, f1:*

$$ACC = \frac{TP + TN}{FP + FN + TP + TN} = 1 - ERR$$

$$ERR = \frac{FP + FN}{FP + FN + TP + TN}$$

$$FPR = \frac{FP}{N} = \frac{FP}{FP + TN}$$

$$PRE = \frac{TP}{TP + FP}$$

**Tudo isso pode ser expandido para multiclass!**  
**Cuidado com a acurácia no caso desbalanceado!**

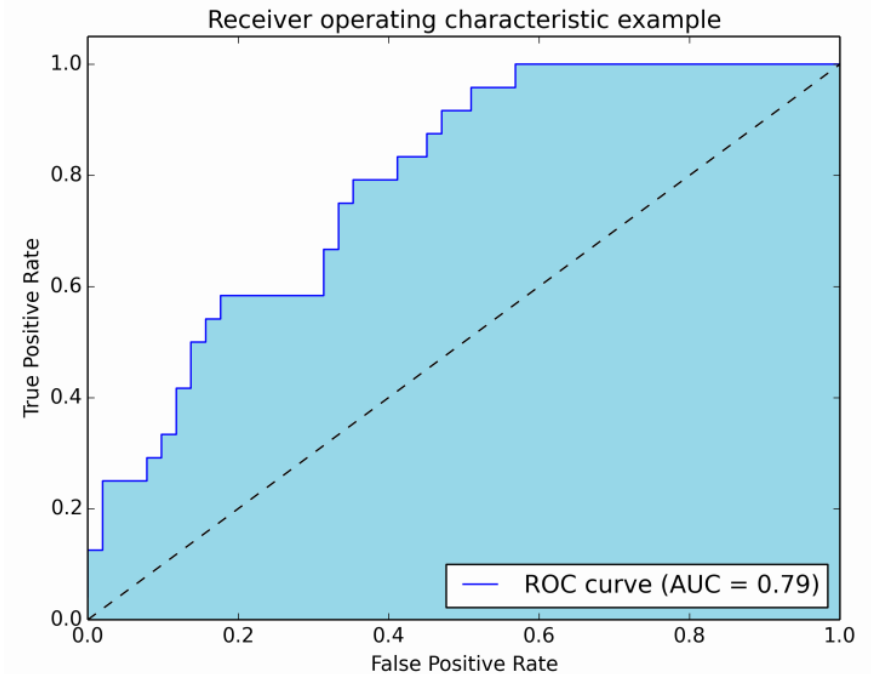
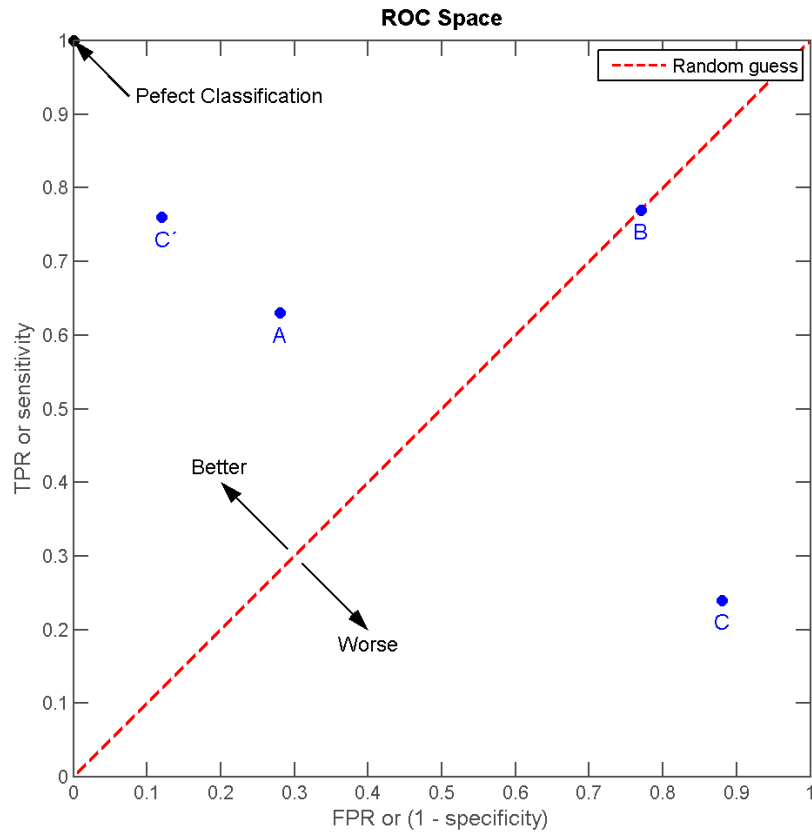
		Predicted class	
		<i>P</i>	<i>N</i>
Actual Class	<i>P</i>	True Positives (TP)	False Negatives (FN)
	<i>N</i>	False Positives (FP)	True Negatives (TN)

$$TPR = \frac{TP}{P} = \frac{TP}{FN + TP}$$

$$REC = TPR = \frac{TP}{P} = \frac{TP}{FN + TP}$$

$$F1 = 2 \frac{PRE \times REC}{PRE + REC}$$

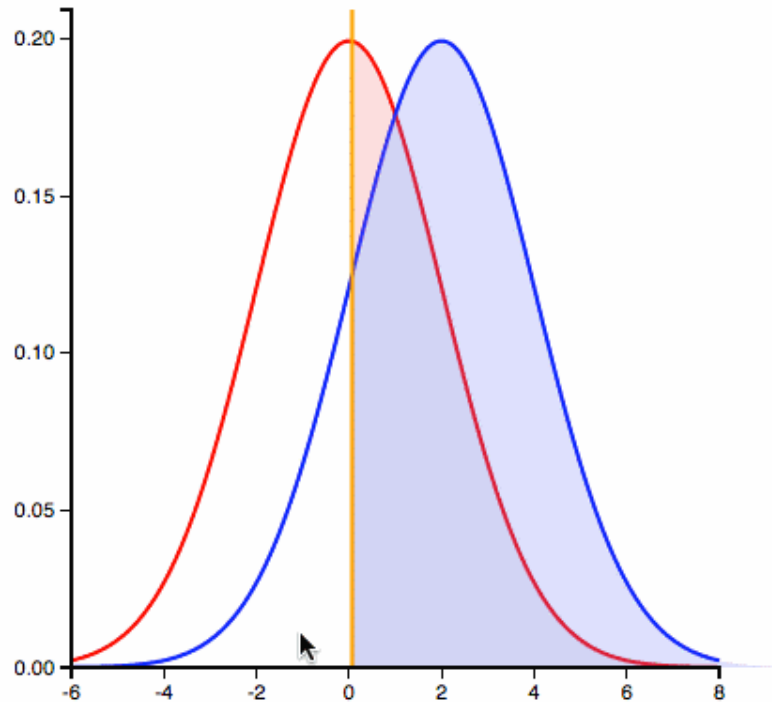
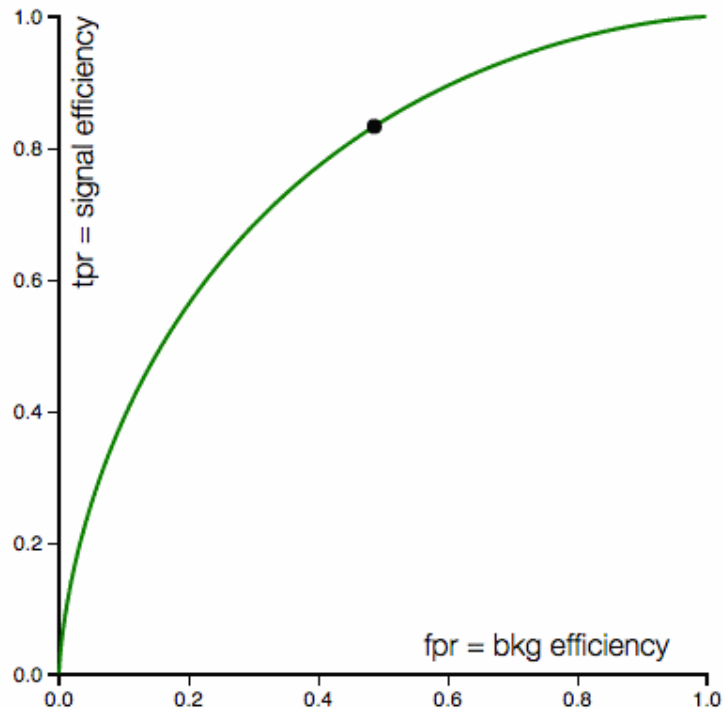
# *Receiver Operating Curve (ROC)*



# *Receiver Operating Curve (ROC)*

## ROC curve demo

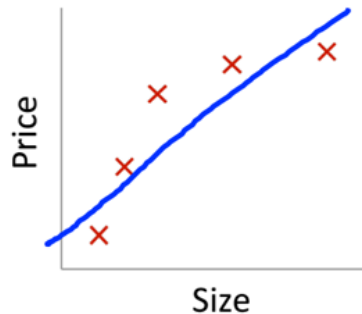
mean #1:  mean #2:  variance #1:  variance #2:



# Mínimos Quadrados - Regressão

- Regressão: minimizar o seguinte funcional em relação à  $\theta$ :

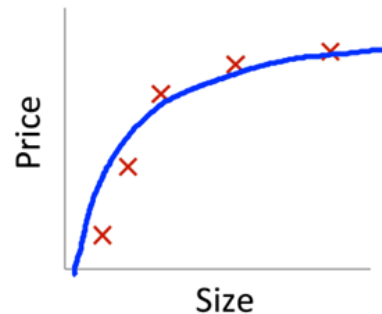
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2$$



$h_{\theta}$

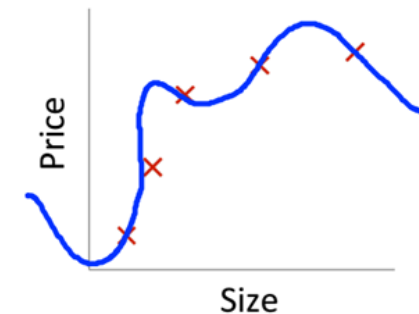
Size  
 $\theta_0 + \theta_1 x$

High bias  
(underfit)



Size  
 $\theta_0 + \theta_1 x + \theta_2 x^2$

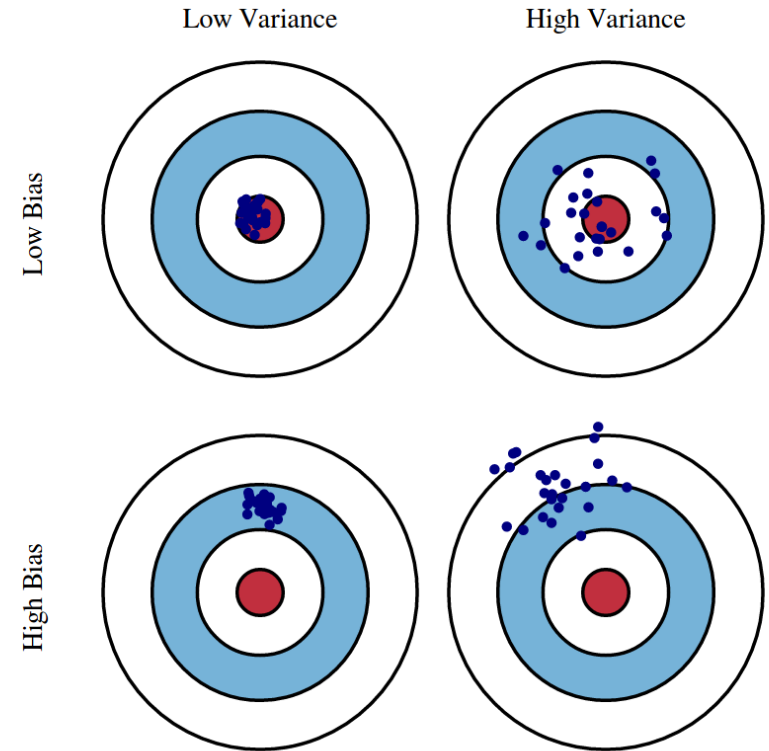
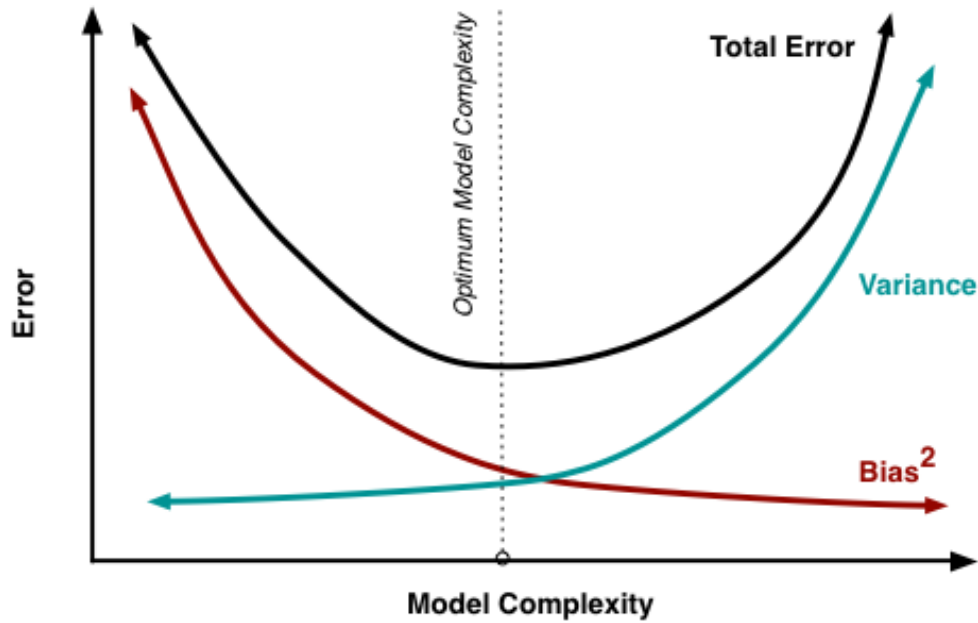
“Just right”



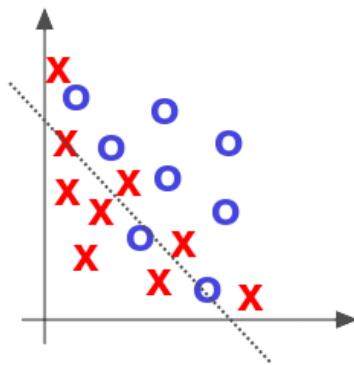
Size  
 $\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$

High variance  
(overfit)

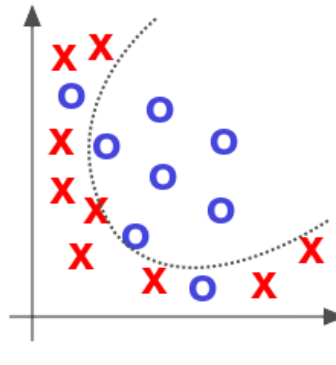
# Bias x Variância



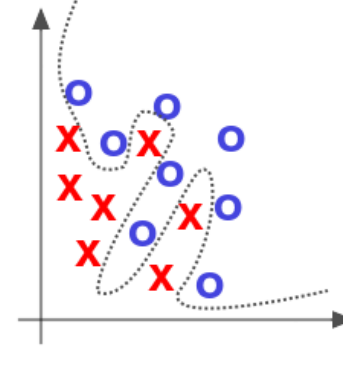
# Overfitting - Classificação



Under Fit

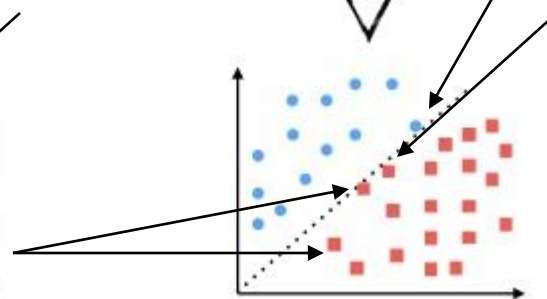
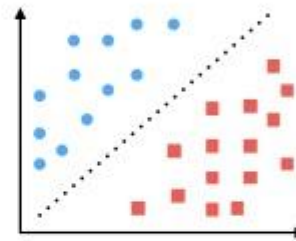
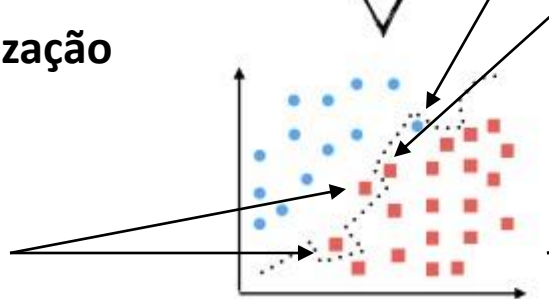
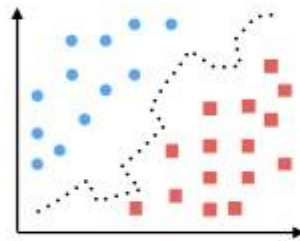


Appropriate



Over Fit

**Capacidade  
de  
Generalização**

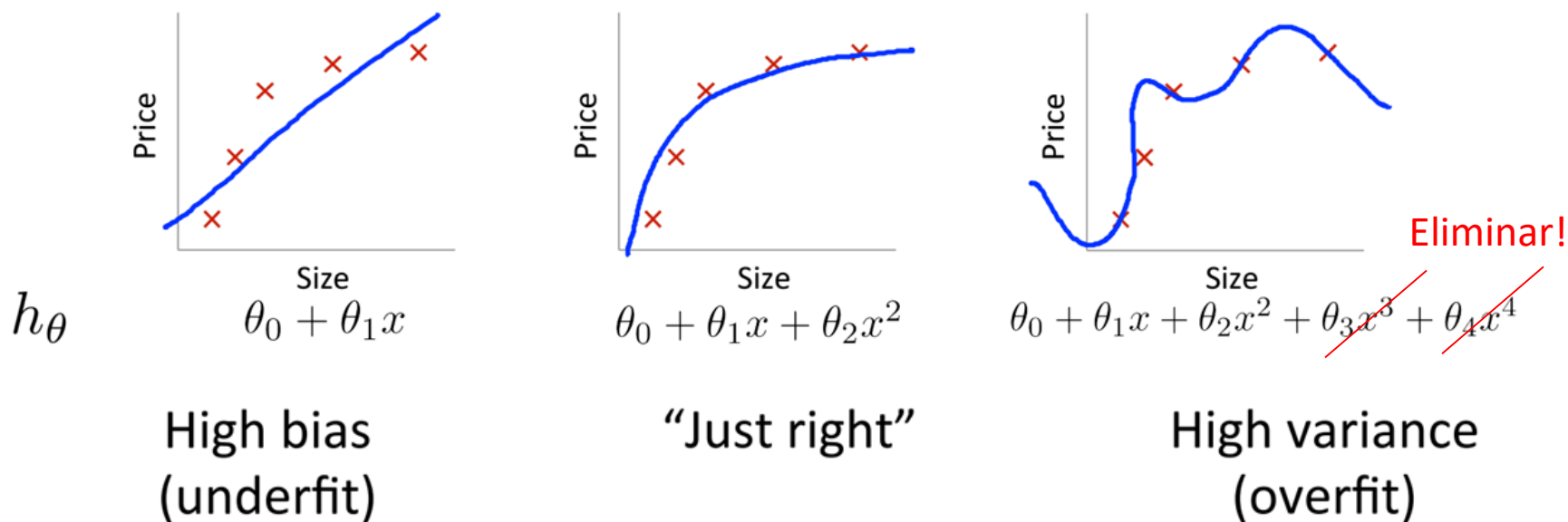




# Regularização

- Sem conhecimento prévio do que se classifica como “verdadeiro”, é possível julgar qual modelo é melhor?
- Na prática não é possível, porém pode-se utilizar um modelo de penalização de modelos excessivamente complexos.
- Isso normalmente se chama **regularização**.

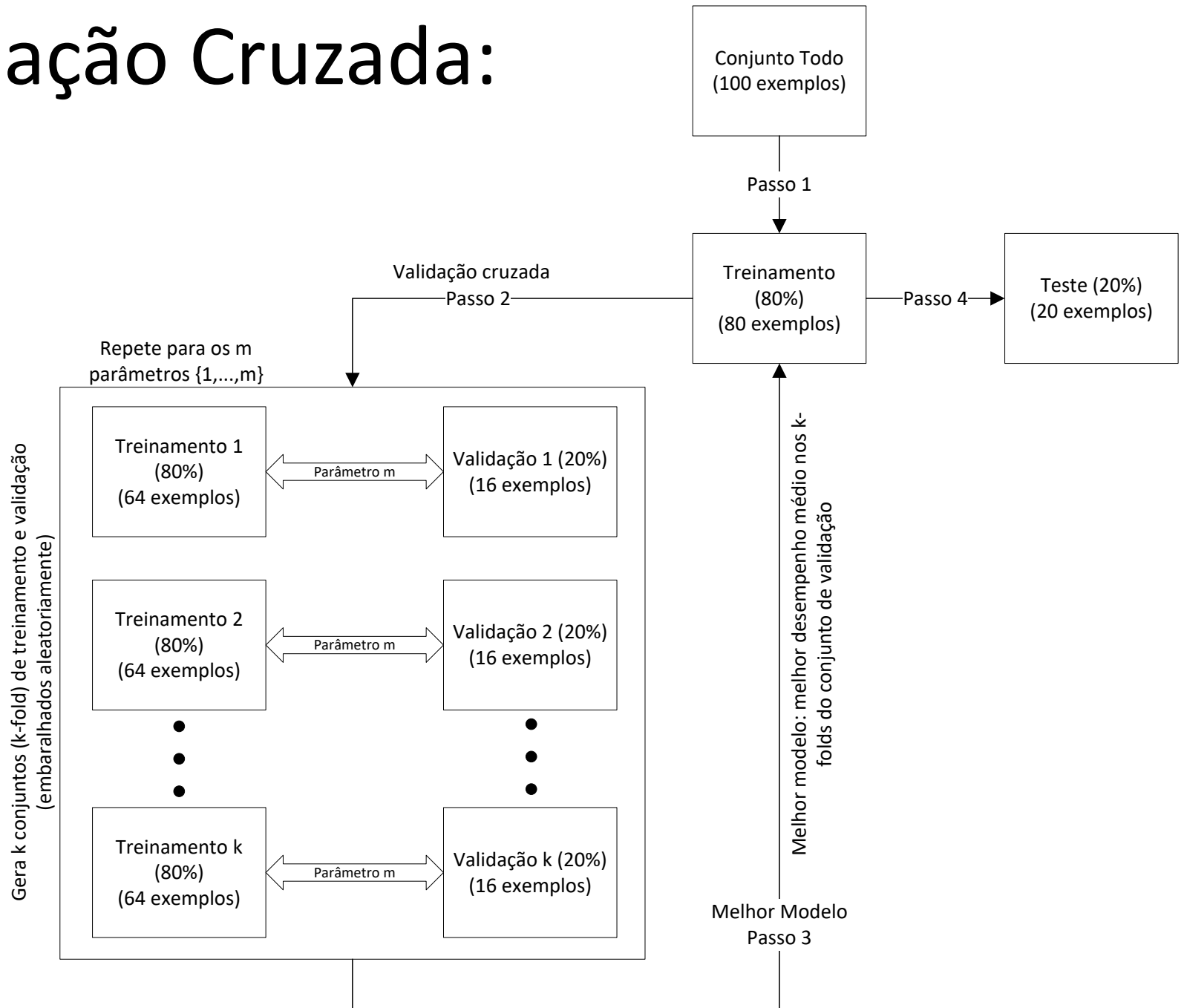
# Regularização



Alternativa: alterar a minimização do seguinte funcional em relação à  $\theta$ :

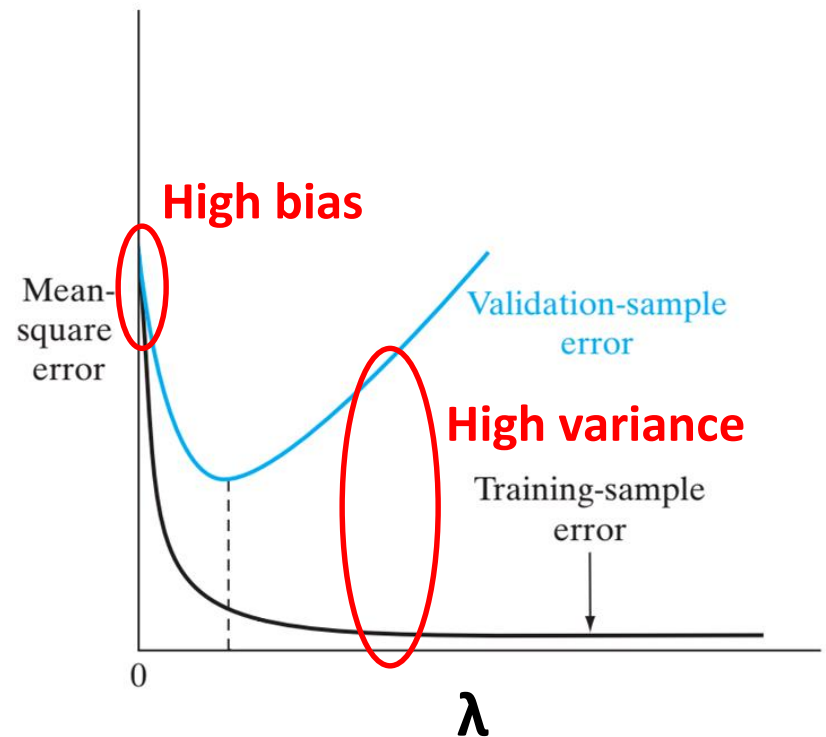
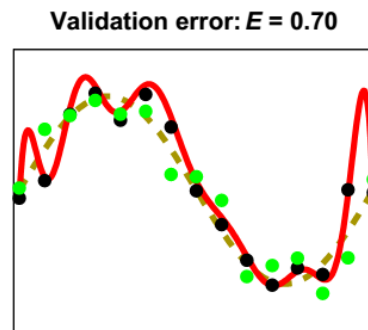
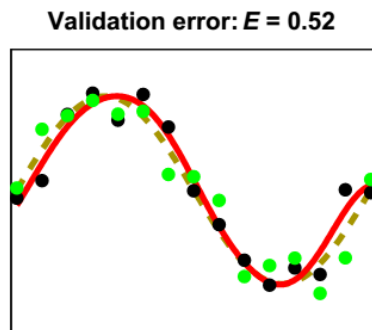
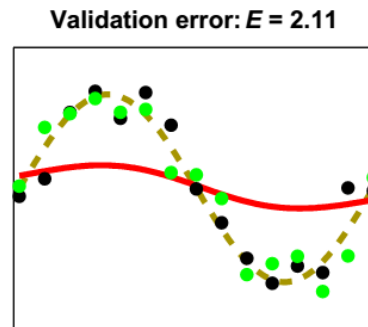
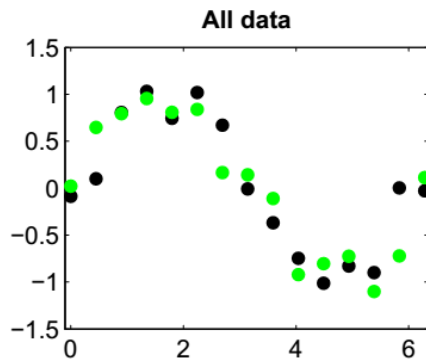
$$J(\theta) = \frac{1}{2m} \left[ \sum_{i=1}^m (h_\theta(x_i) - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

# Validação Cruzada:



# Regularização

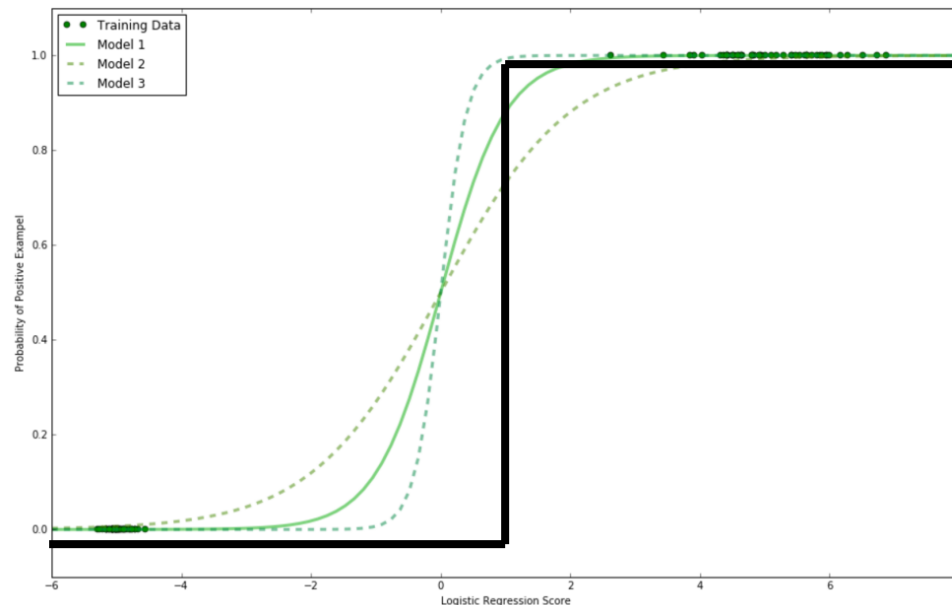
Selecionar  $\lambda$  com o menor erro no conjunto de validação:



Complexidade (ordem do modelo)

# Regularização – Logistic Regression

- Por ser um problema de otimização baseado em likelihood – overfitting!
- Caso linear: no limite pode ser uma função degrau, localizada em diferentes pontos (limiares).



- Regularização quadrática:

$$f'(\mathbf{w}) = \text{NLL}(\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

$$\mathbf{g}'(\mathbf{w}) = \mathbf{g}(\mathbf{w}) + \lambda \mathbf{w}$$

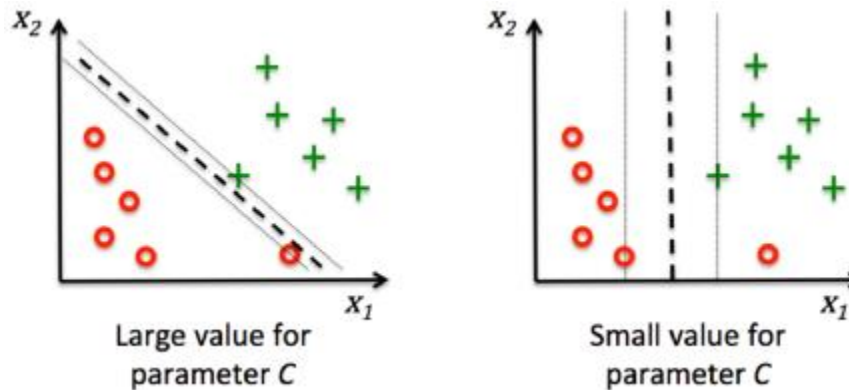
$$\mathbf{H}'(\mathbf{w}) = \mathbf{H}(\mathbf{w}) + \lambda \mathbf{I}$$

- Ou Bayesian Logistic Regression!

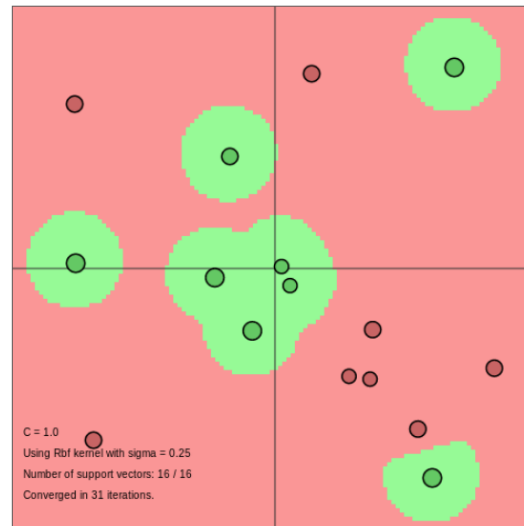
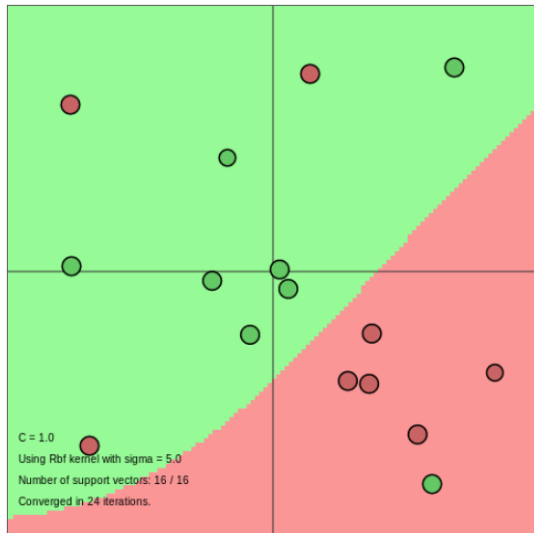
# Regularização - SVM

- Da aula anterior: minimize  $J(\mathbf{w}, \mathbf{w}_0, \boldsymbol{\xi}) = \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$

Influência do  
parâmetro C:



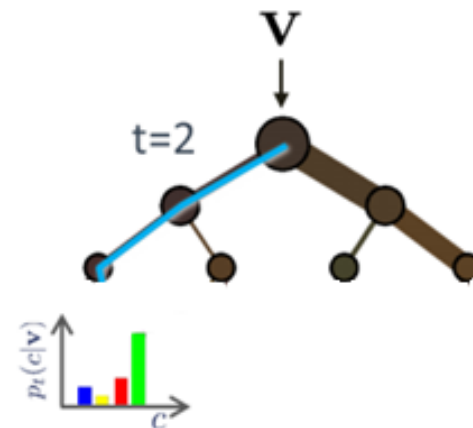
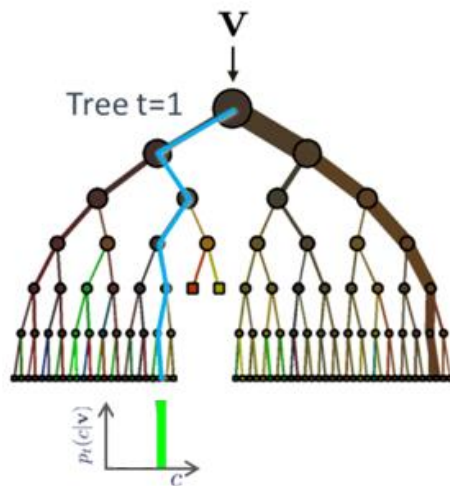
O que isso  
significa?  
Relação com  
regularização!



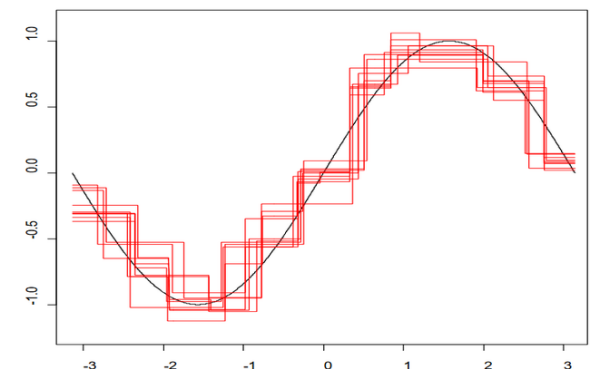
- Maior o sigma → mais “suave” o classificador → reduz overfitting → classificador global
- Menor o sigma → mais “sharp” o classificador → overfitting → classificador local

# Regularização – *Decision Tree*

- Limitando máximo comprimento das árvores:



- Ensembles – Random Forest:

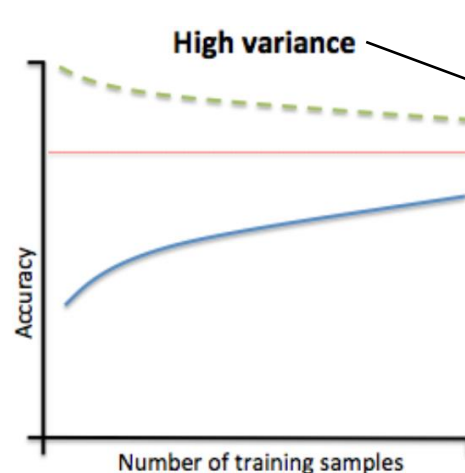
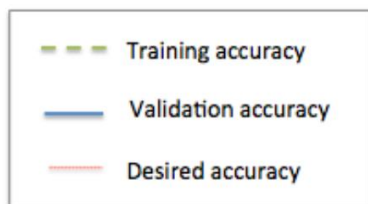
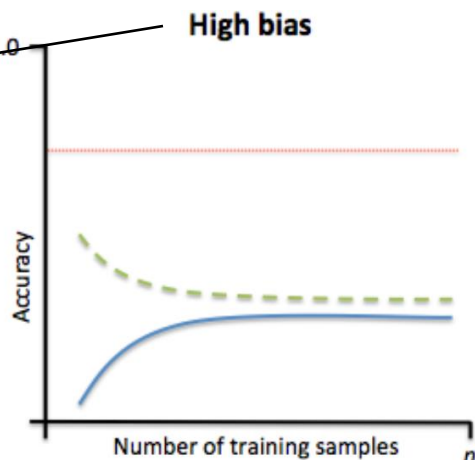


# Análise de Curvas

Mais dados de treinamento provavelmente não vão ajudar significativamente

Tente adicionar novas *features*!

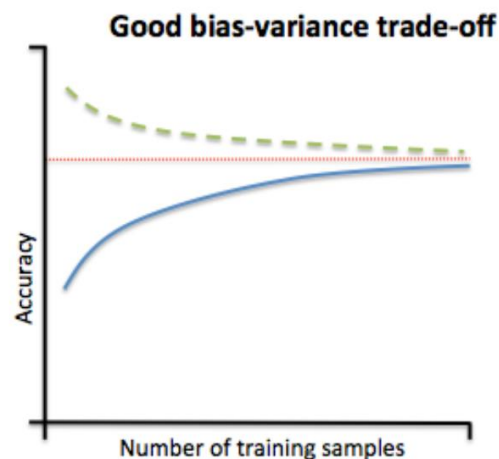
Verificar a validação cruzada. Aumentar *range* das variáveis.



**Overfitting!**

Mais dados de treinamento provavelmente vão ajudar!

Tente reduzir núm. de *features*



Verificar a validação cruzada. Diminuir *range* das variáveis.



# Outras Dicas

- Se você é um especialista no problema, faça uma inspeção visual dos dados e verifique se você consegue diferenciar o que o classificador está errando;
- Checar a anotação (rótulos) dos dados;
- Verifique os erros principais da matriz de confusão;

# Para o trabalho final

- Muito importante analisar os pontos discutidos aqui:
  - Análise de curvas
  - Curva ROC ou Matriz de confusão
  - O que você fez (como agiu) para contornar um problema

# Referências

- Aulas Prof. Andrew Ng
- Aulas Prof. David Tax
- Livro Python Machine Learning - Sebastian Raschka - Capítulo 6