



# Modelagem preditiva, Navalha de Occam e Processos Gaussianos para Machine Learning.

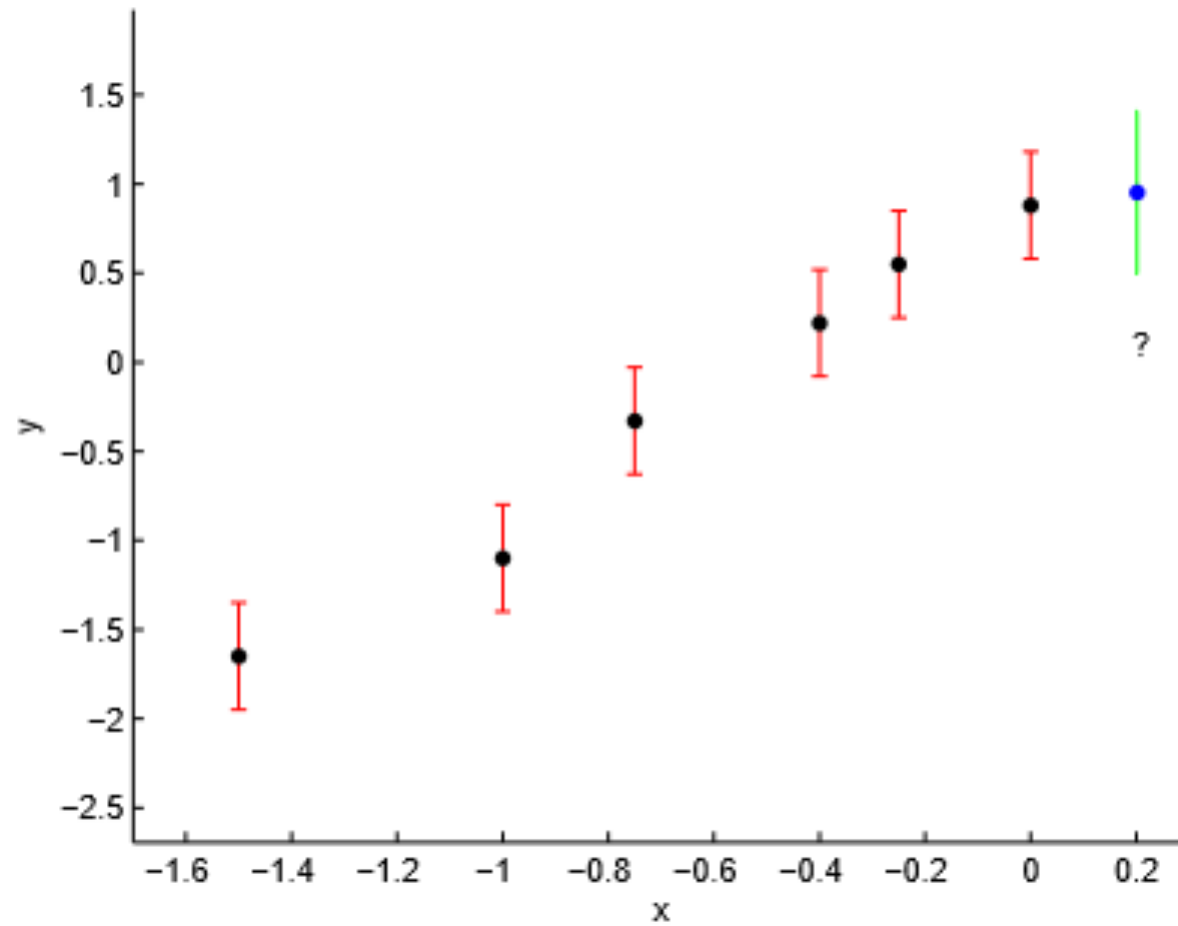
Fabio Ramos

Departamento de Matemática Aplicada

Instituto de Matemática

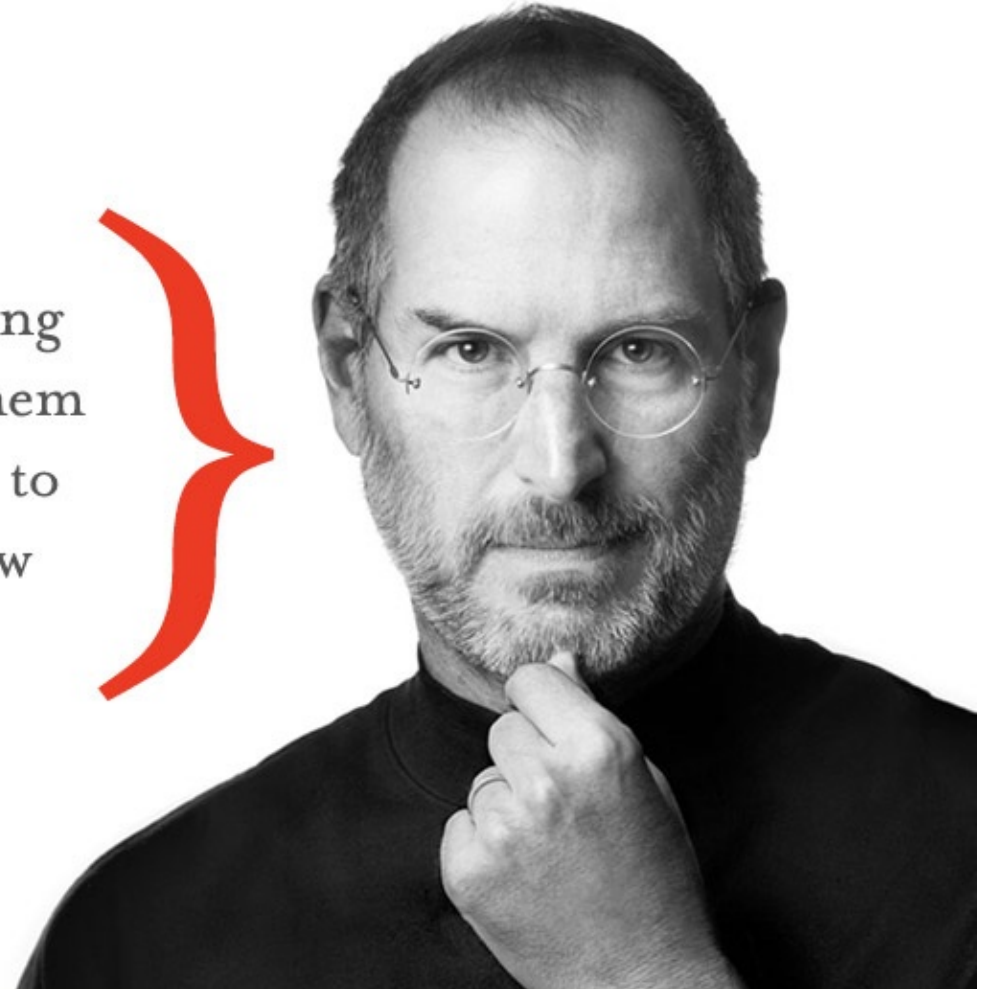
UFRJ

# O OBJETIVO

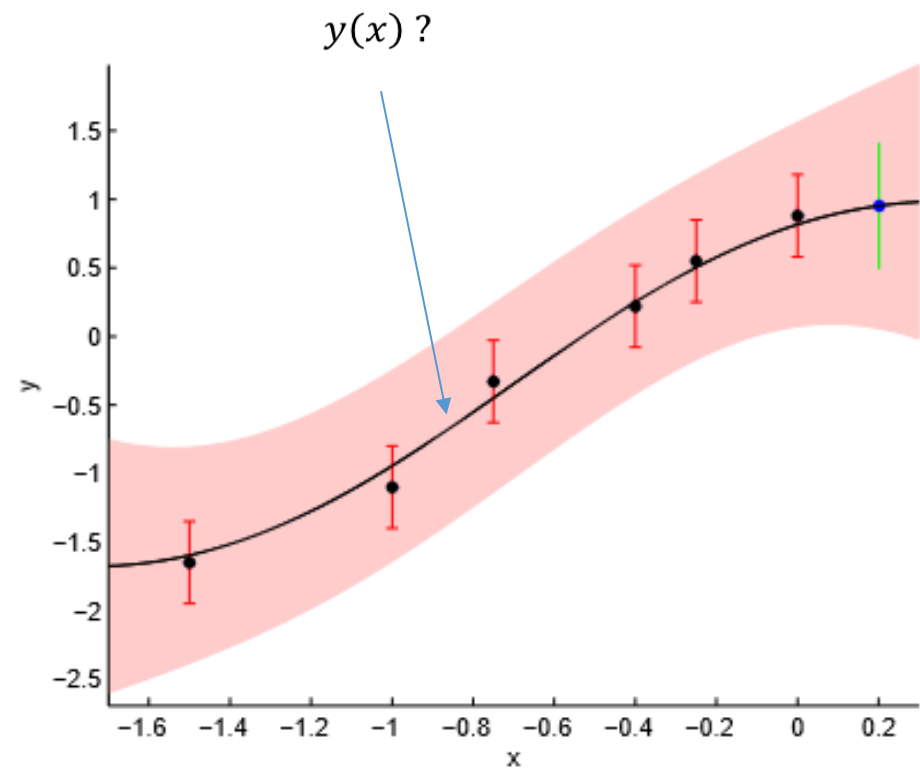
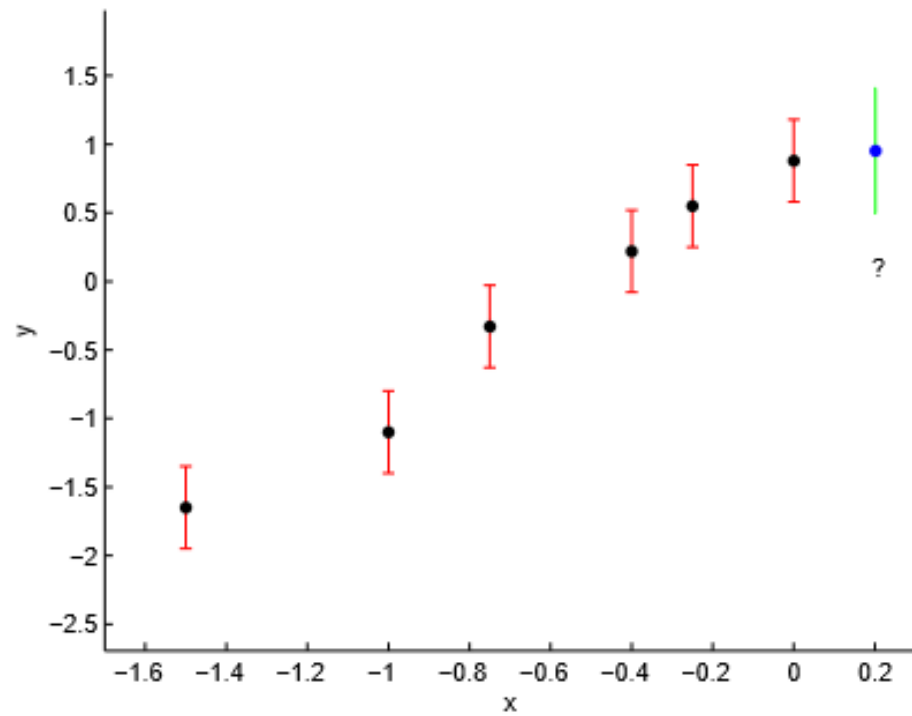


You can't connect the dots looking forward; you can only connect them looking backwards. So you have to trust that the dots will somehow connect in your future.

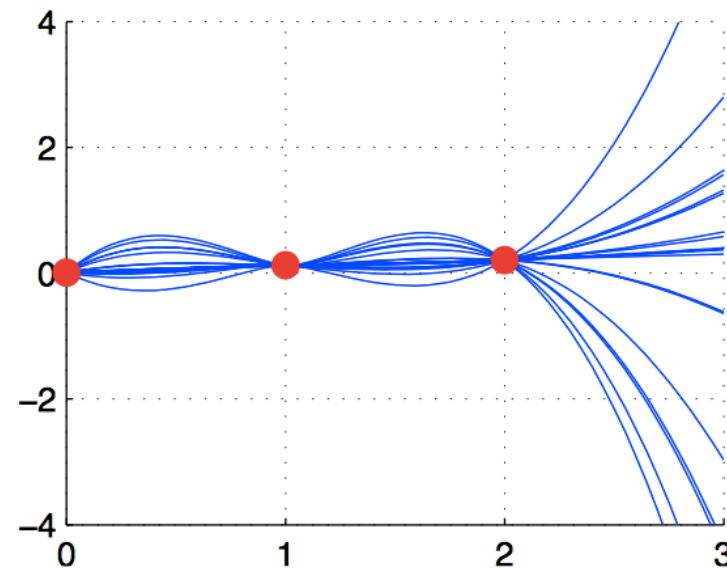
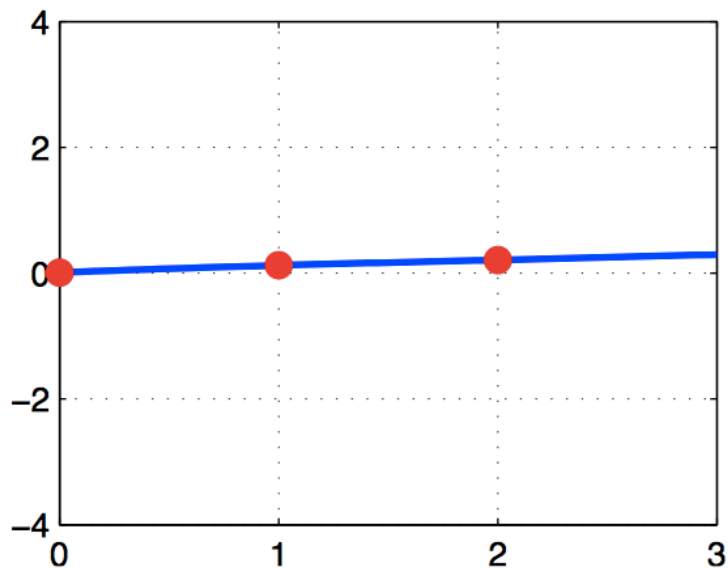
*-Steve Jobs*



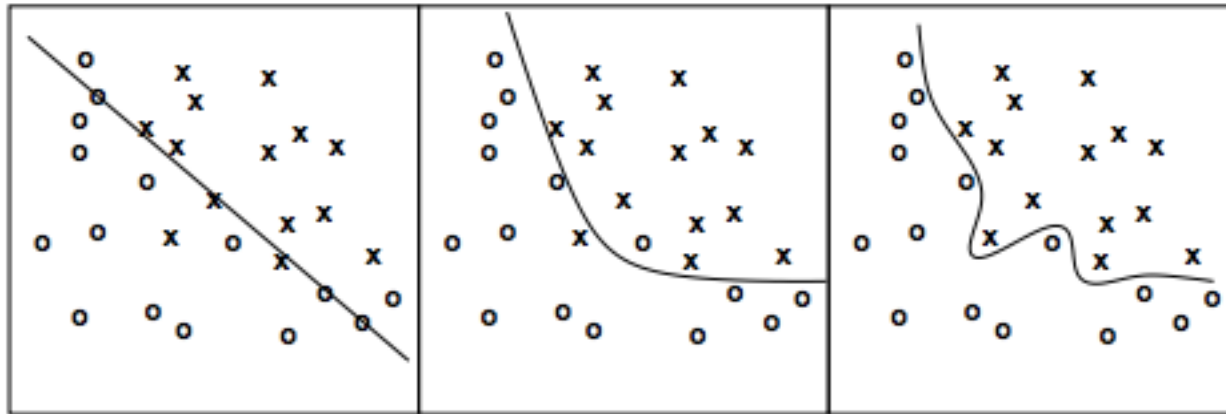
# O OBJETIVO



# Em qual modelo acreditar ?



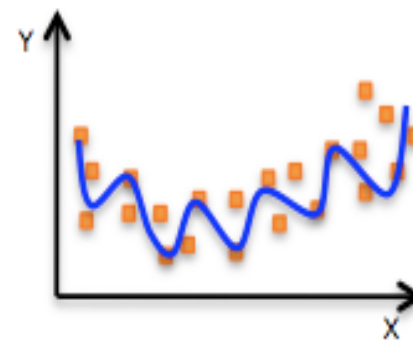
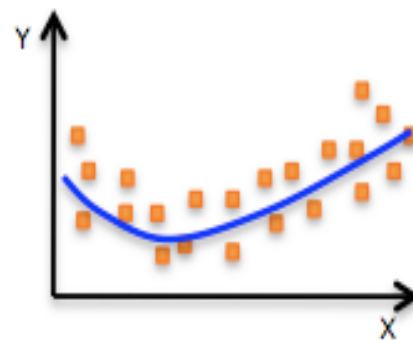
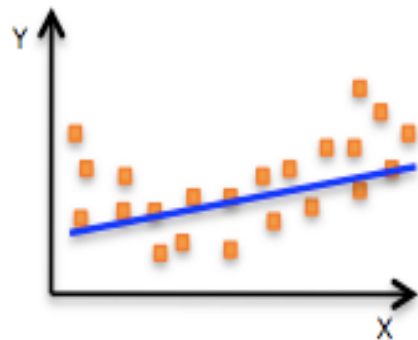
*Navalha de Occam* - A hipótese mais provável é a mais simples e consistente com os dados



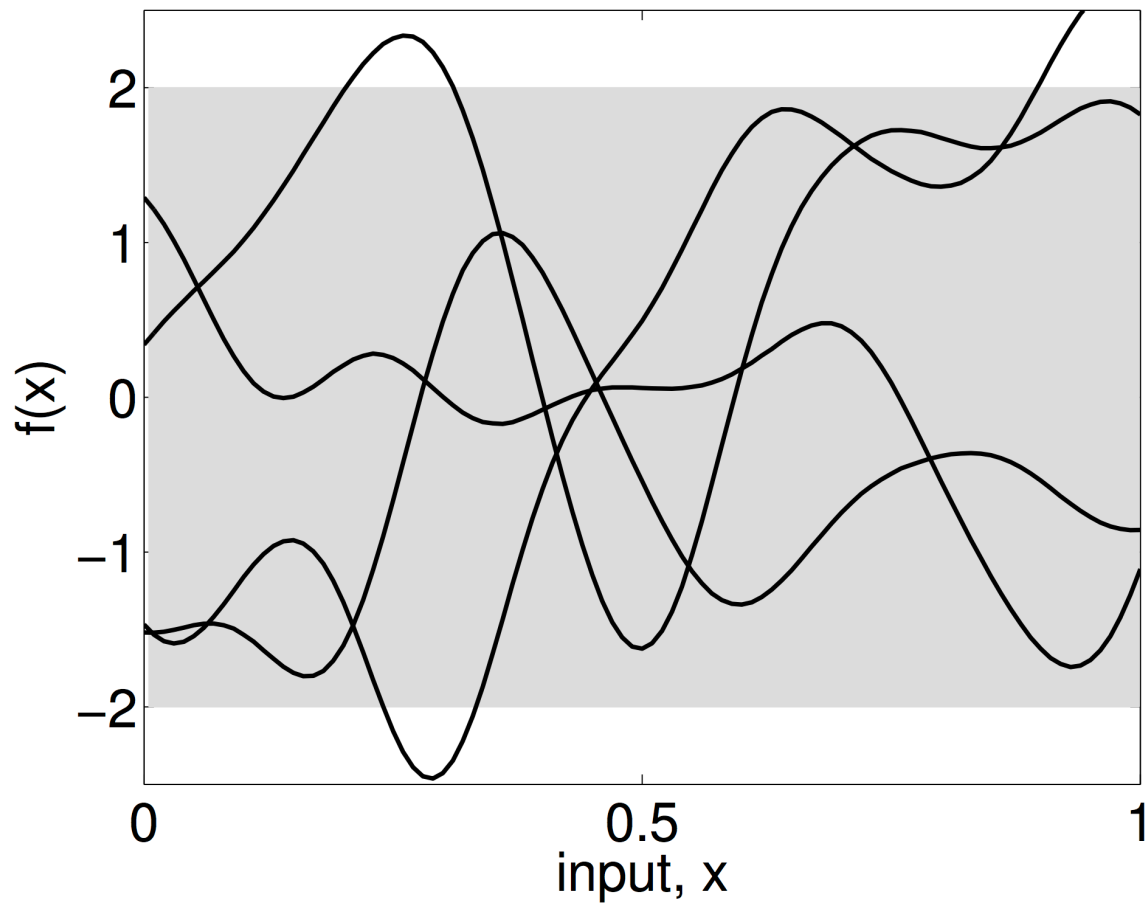
Inadequado

Solução de  
Compromisso

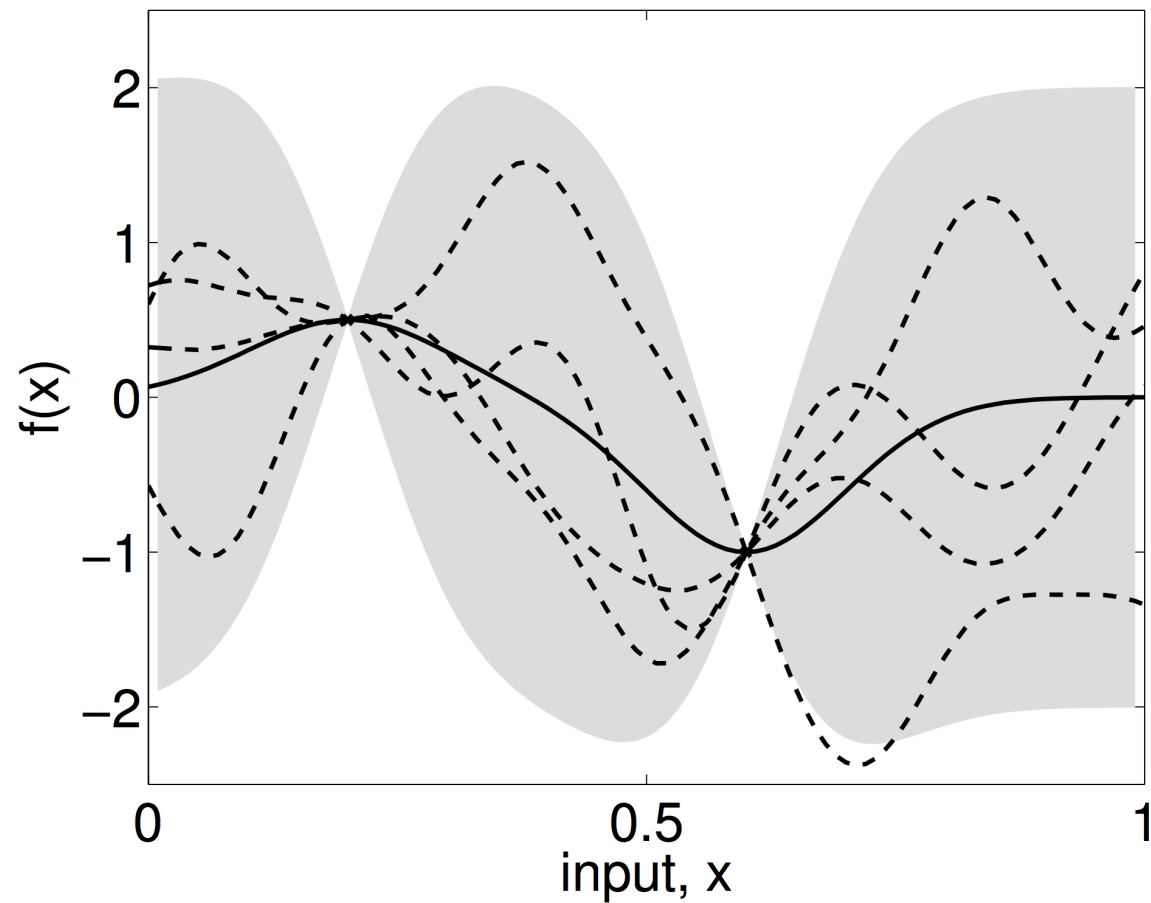
Overfitting



# Gaussian Processes - Priors over functions



(a), prior

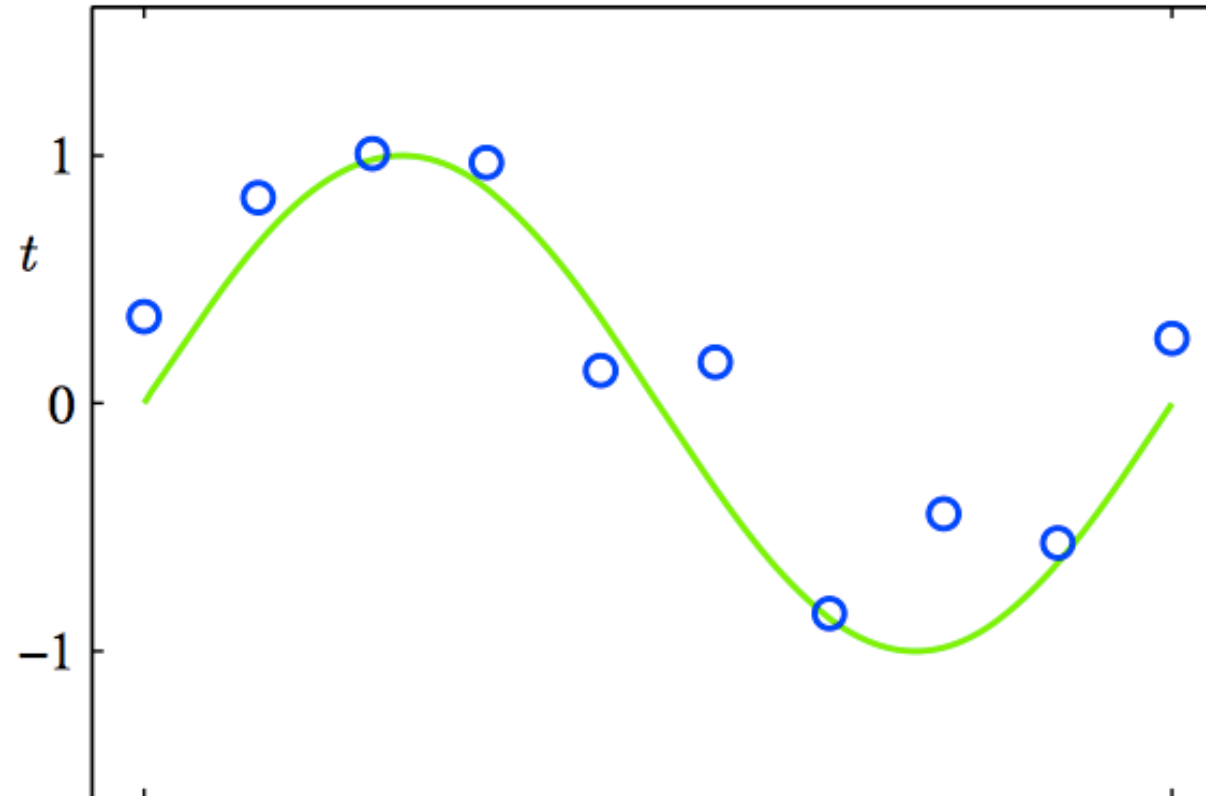


(b), posterior

# Interpolação polinomial

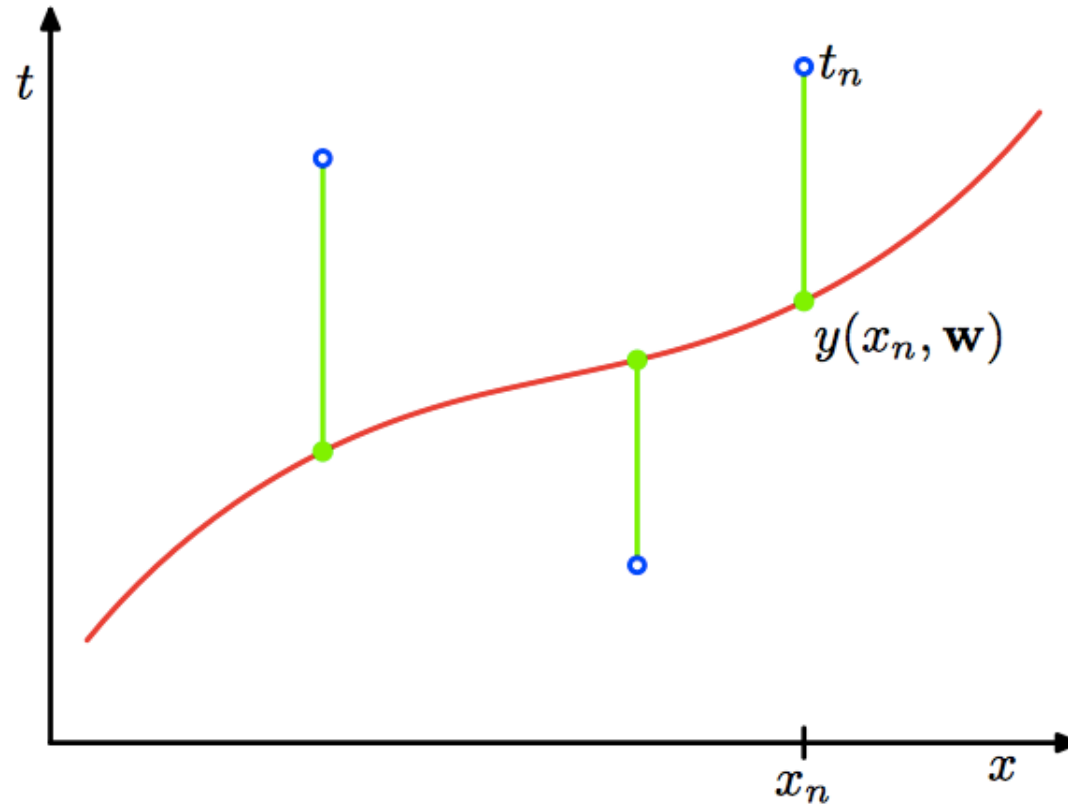


# SINAL (senoidal) + RUÍDO



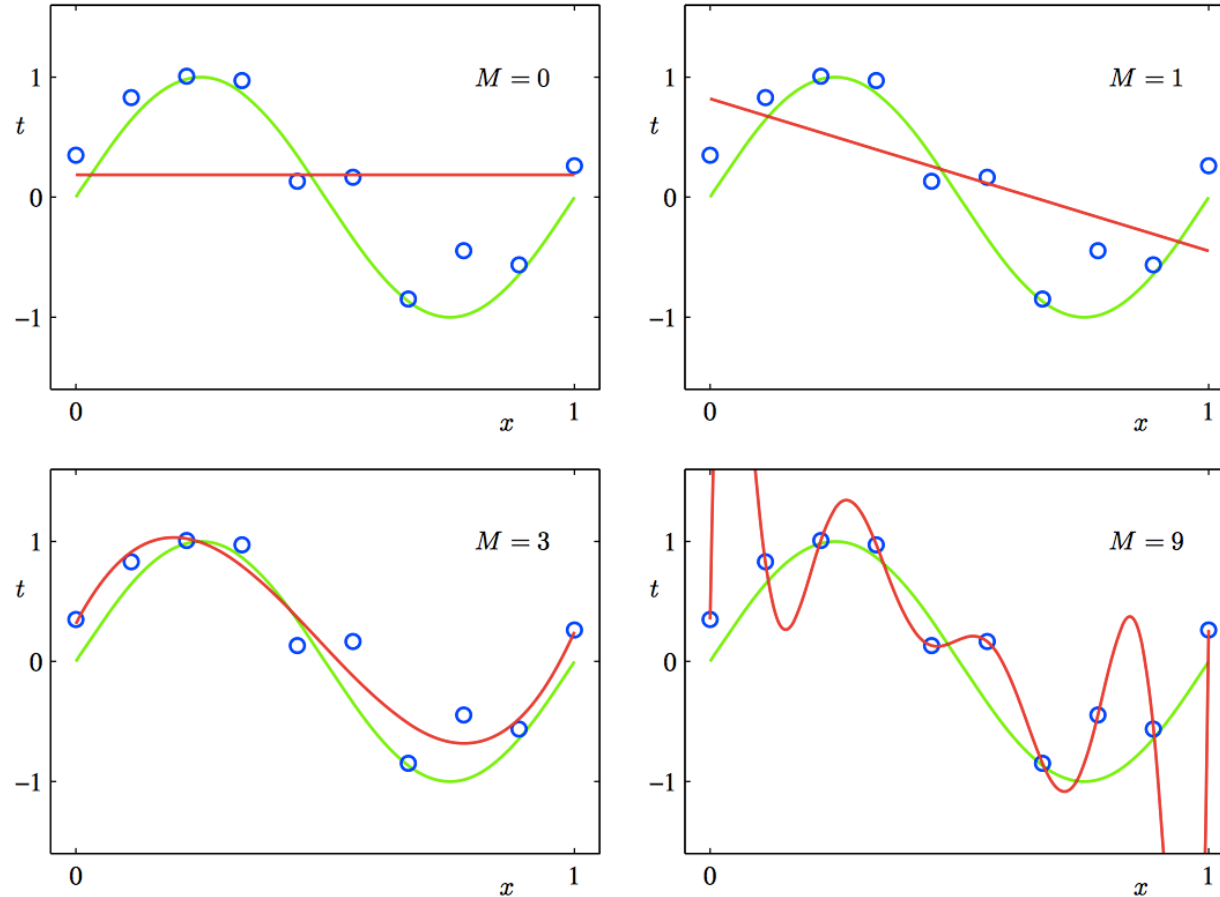
$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

# Ajustando: erro médio quadrático



$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \xrightarrow{\text{minimizador } \mathbf{w}^*} E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

# Melhor Interpolação Polinomial

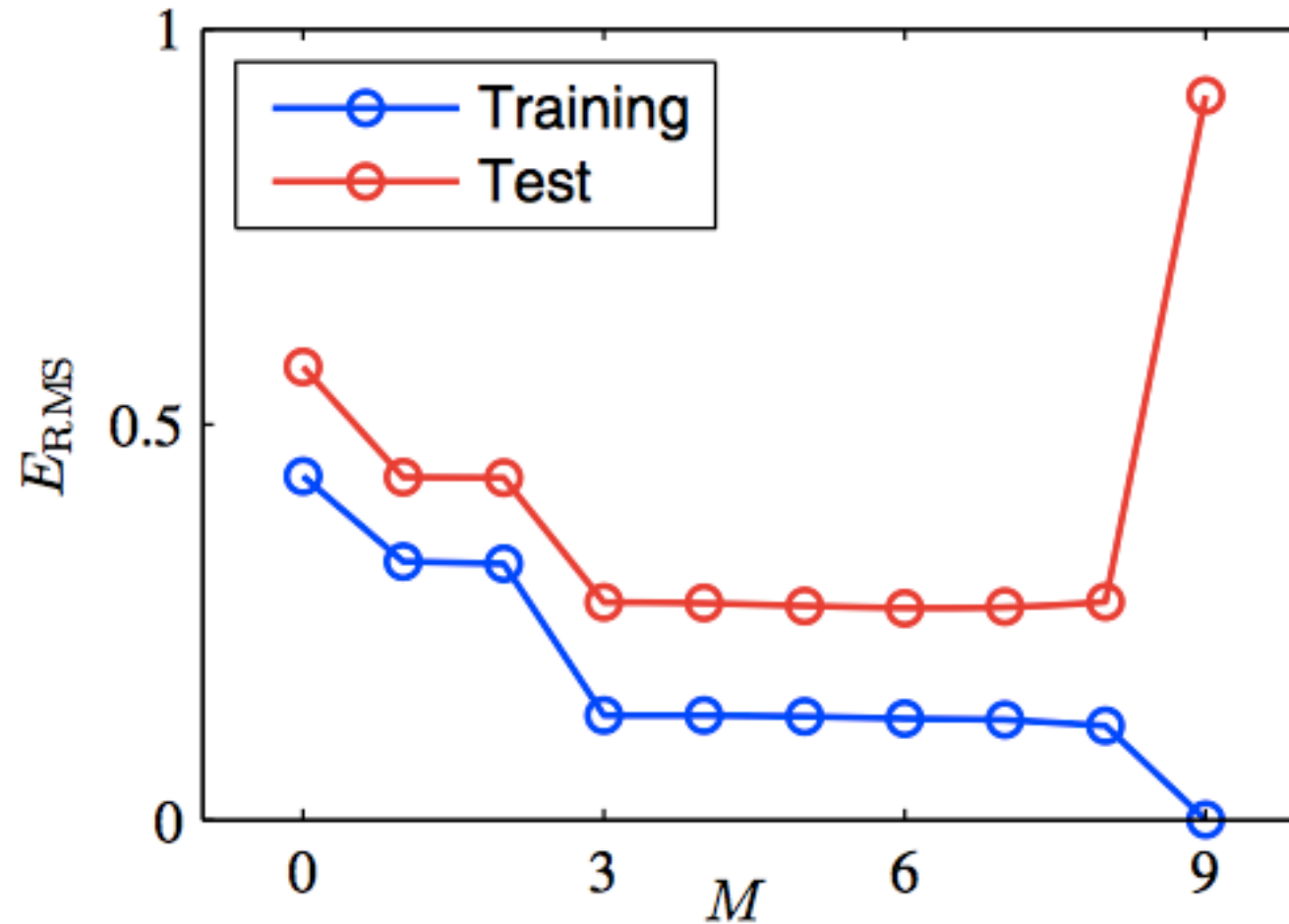


$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_j x^j$$

# Coeficientes dos polinômios interpoladores

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	7.99	232.37
$w_2^*$			-25.43	-5321.83
$w_3^*$			17.37	48568.31
$w_4^*$				-231639.30
$w_5^*$				640042.26
$w_6^*$				-1061800.52
$w_7^*$				1042400.18
$w_8^*$				-557682.99
$w_9^*$				125201.43

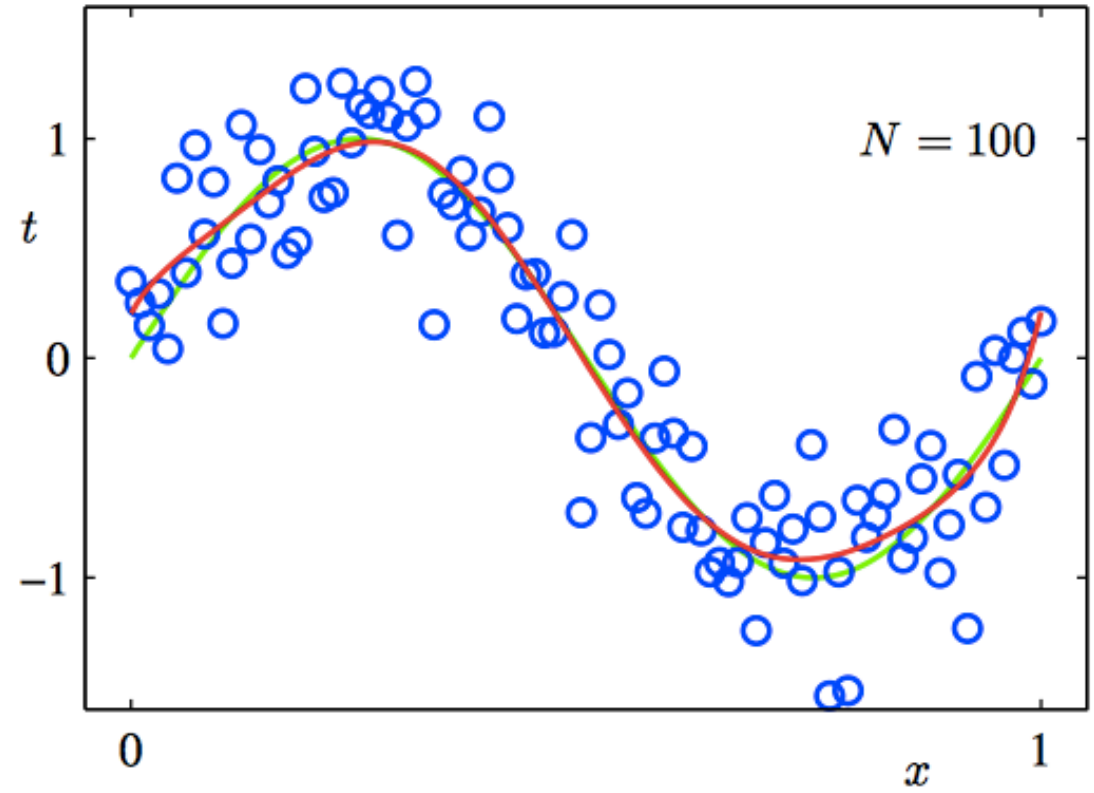
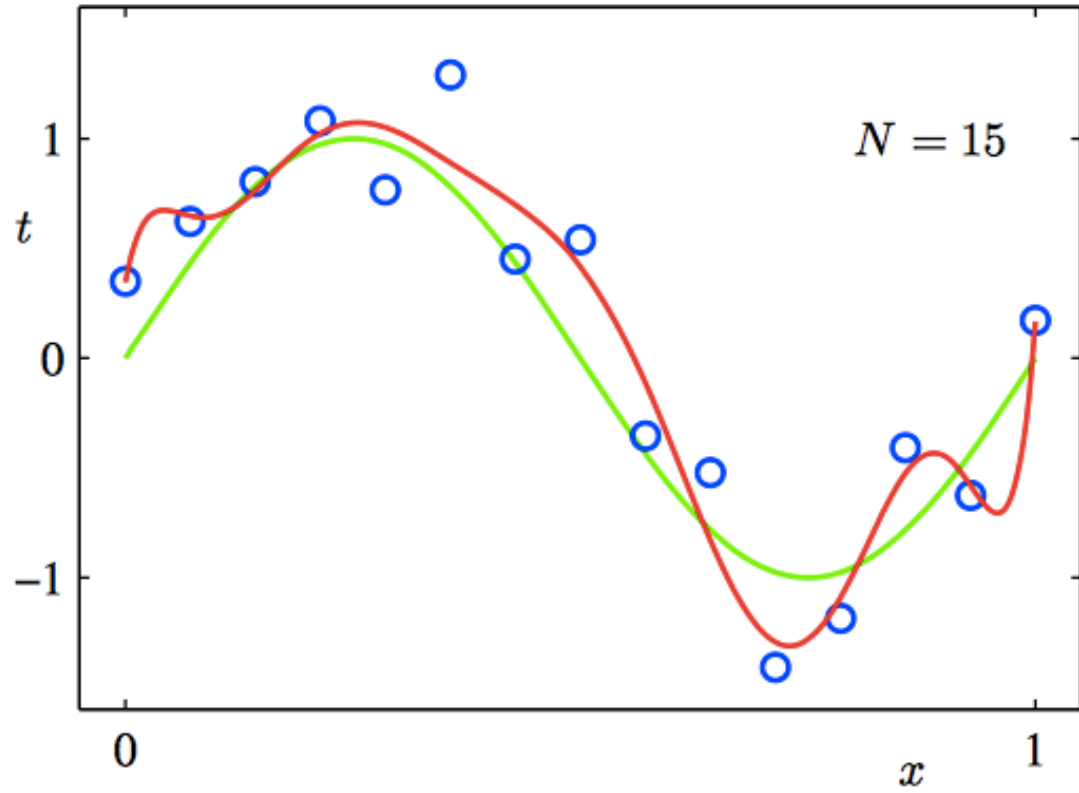
# Overfitting (Treino é treino, jogo é jogo)



Conjunto de 100 pontos de teste.

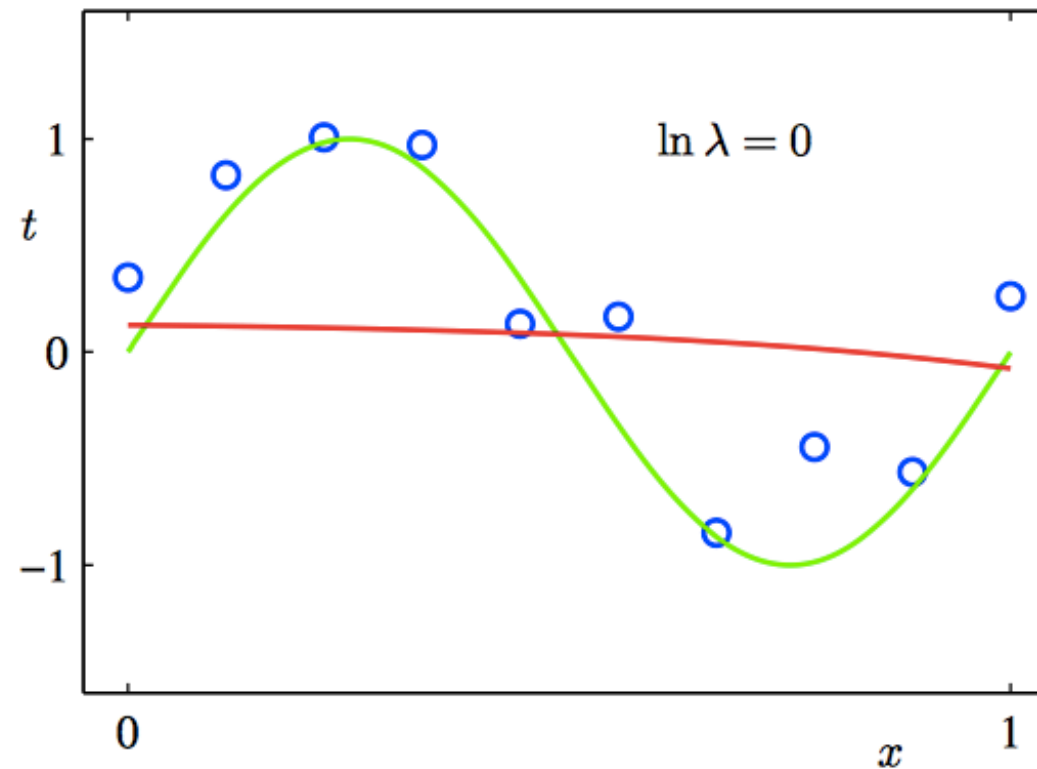
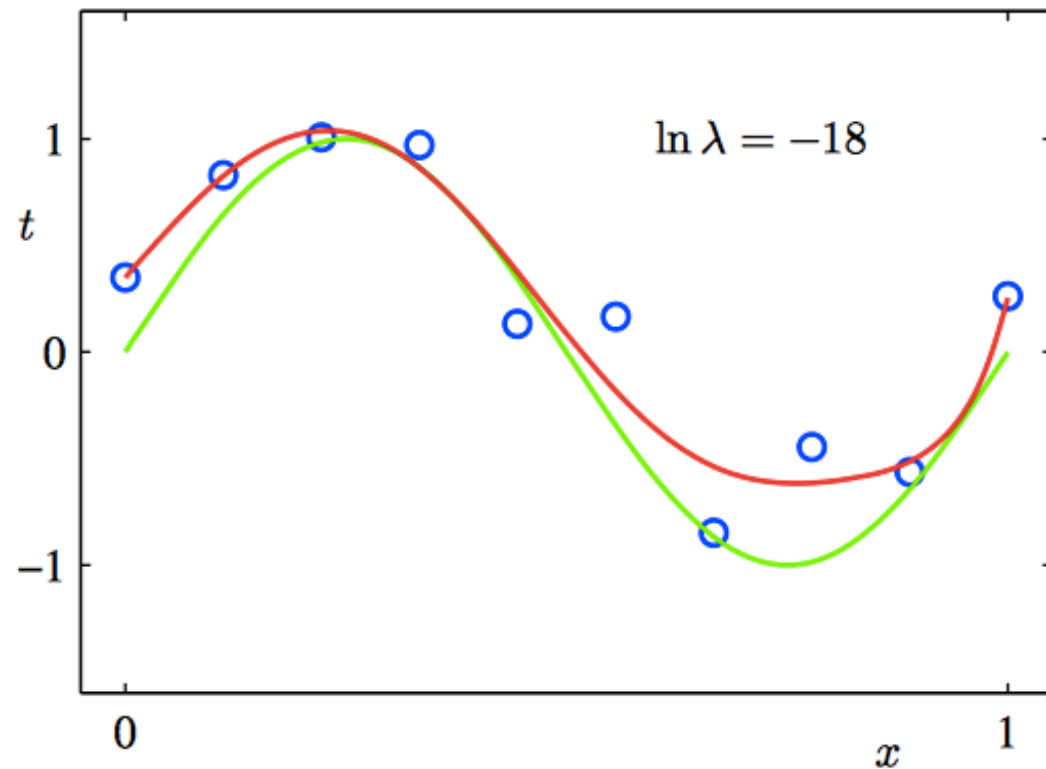
# Treinando com mais dados

$M=9$



# Penalização e Regularização

M=9



$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

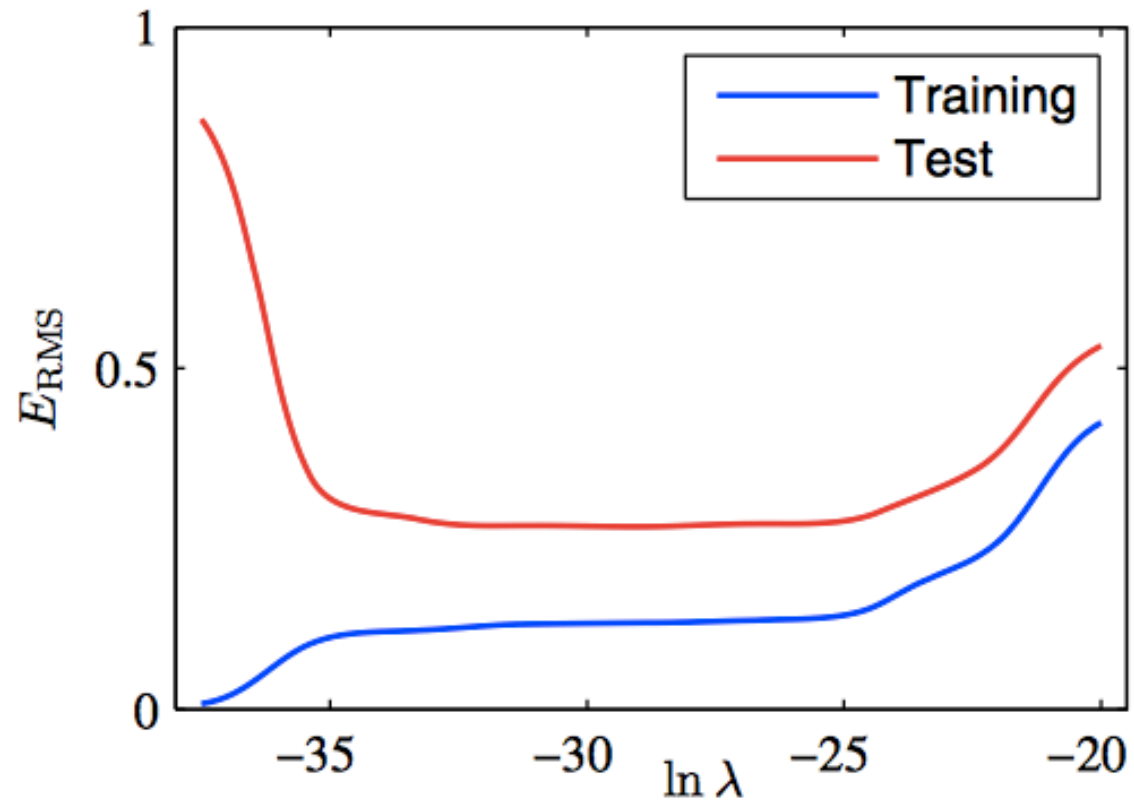
# Coeficientes dos polinômios interpoladores - o efeito regularizador

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232.37	4.74	-0.05
$w_2^*$	-5321.83	-0.77	-0.06
$w_3^*$	48568.31	-31.97	-0.05
$w_4^*$	-231639.30	-3.89	-0.03
$w_5^*$	640042.26	55.28	-0.02
$w_6^*$	-1061800.52	41.32	-0.01
$w_7^*$	1042400.18	-45.95	-0.00
$w_8^*$	-557682.99	-91.53	0.00
$w_9^*$	125201.43	72.68	0.01



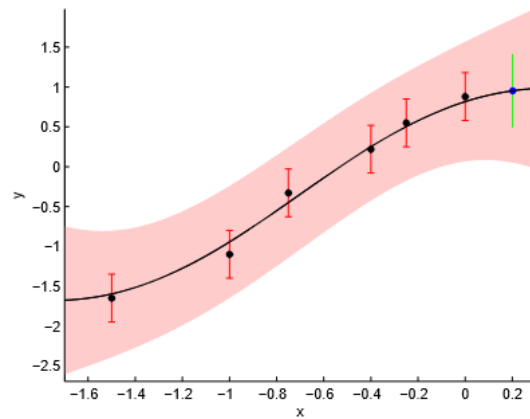
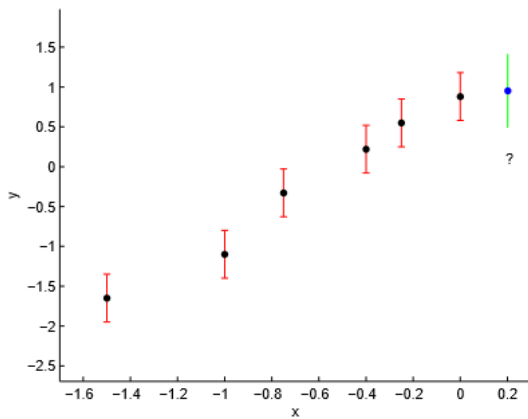
# Domando o Overfitting (Equilibrando Treino e Jogo)

M=9



# Uma abordagem bayesiana

”Teoria da probabilidade não é nada além de senso comum reduzido ao cálculo” - Laplace



”Senso comum não é nada comum” -  
Voltaire

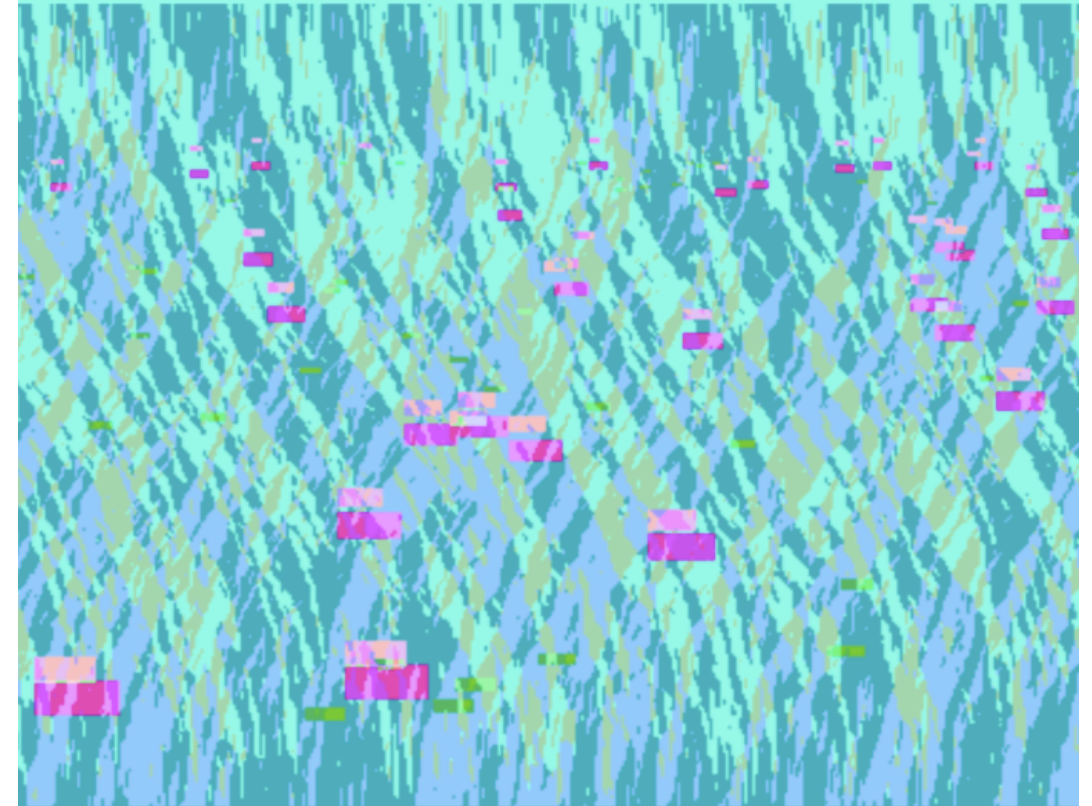
## Sir David John Cameron MacKay



(22 de Abril 1967 – 14 de Abril 2016)

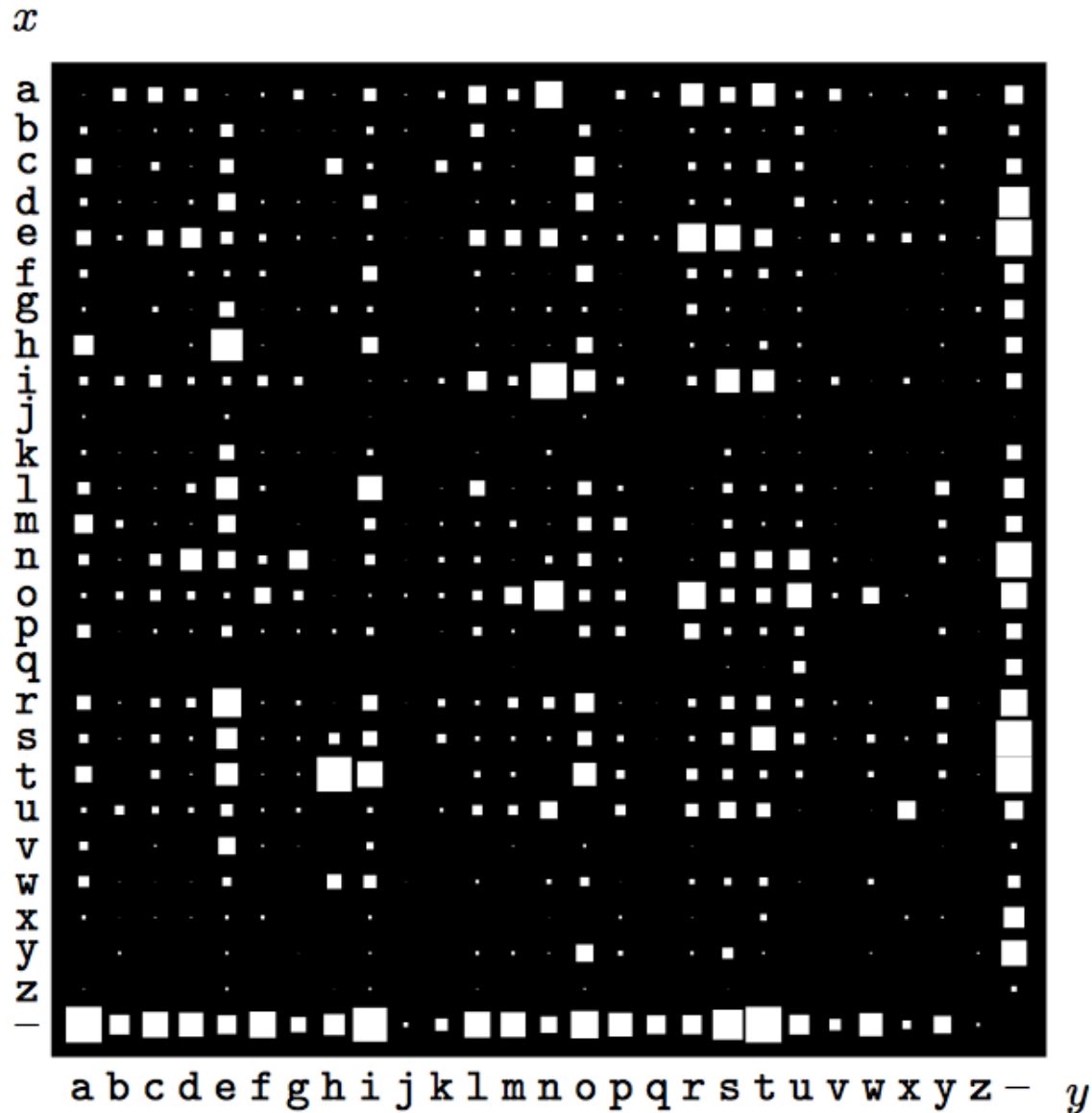
David J.C. MacKay

## Information Theory, Inference, and Learning Algorithms



Cambridge University Press, 2003

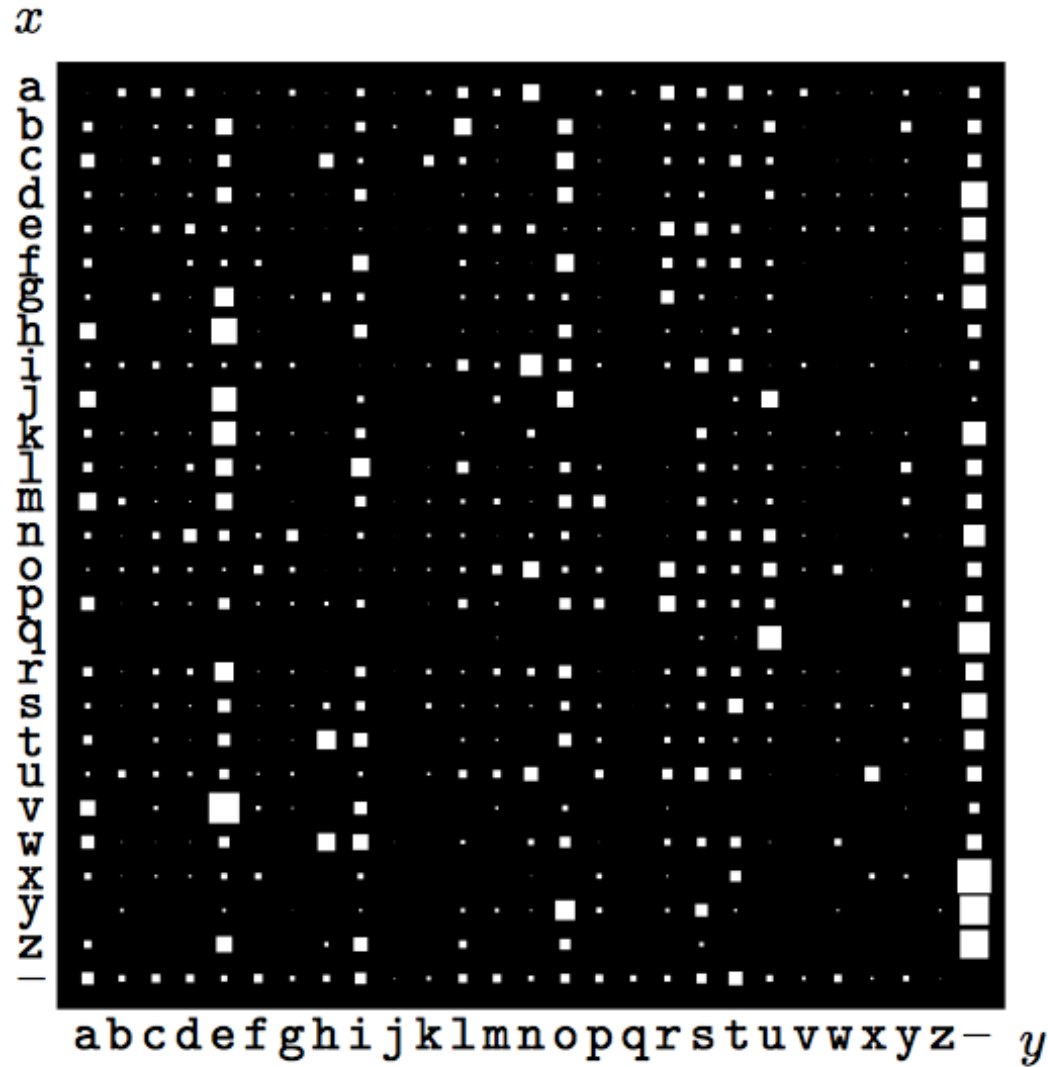
# Teorema de Bayes



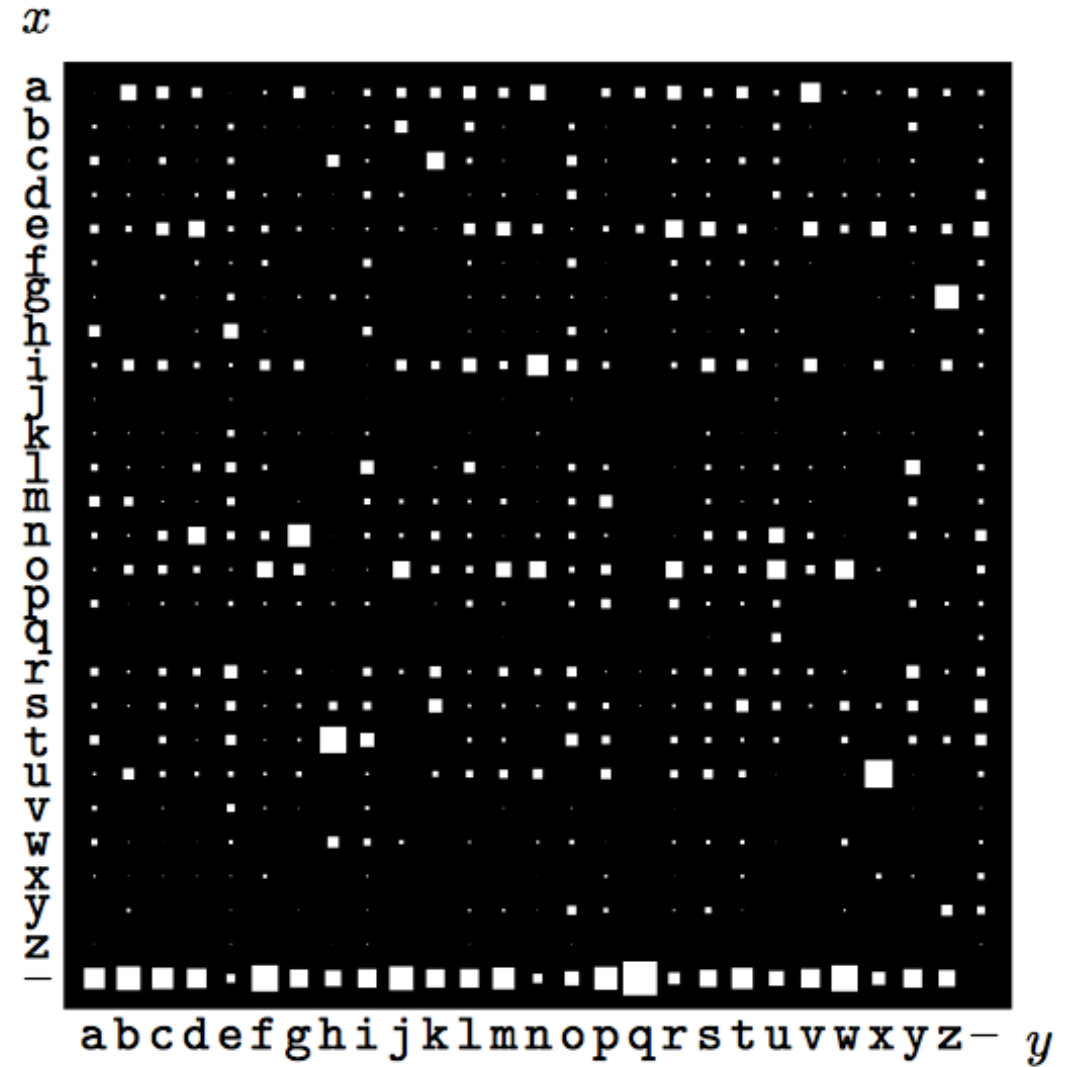
# Busca por padrões – FAQ Linux

$i$	$a_i$	$p_i$	
1	a	0.0575	a
2	b	0.0128	b
3	c	0.0263	c
4	d	0.0285	d
5	e	0.0913	e
6	f	0.0173	f
7	g	0.0133	g
8	h	0.0313	h
9	i	0.0599	i
10	j	0.0006	j
11	k	0.0084	k
12	l	0.0335	l
13	m	0.0235	m
14	n	0.0596	n
15	o	0.0689	o
16	p	0.0192	p
17	q	0.0008	q
18	r	0.0508	r
19	s	0.0567	s
20	t	0.0706	t
21	u	0.0334	u
22	v	0.0069	v
23	w	0.0119	w
24	x	0.0073	x
25	y	0.0164	y
26	z	0.0007	z
27	-	0.1928	-

# Teorema de Bayes

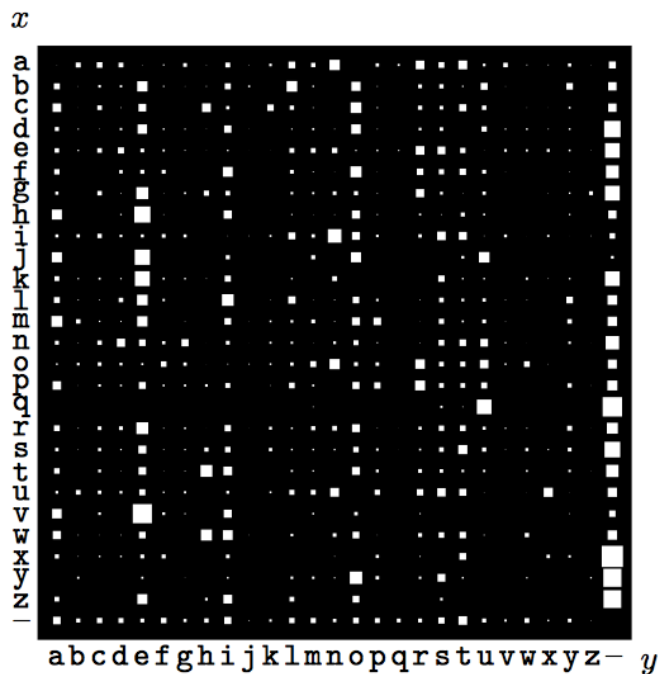


(a)  $P(y|x)$

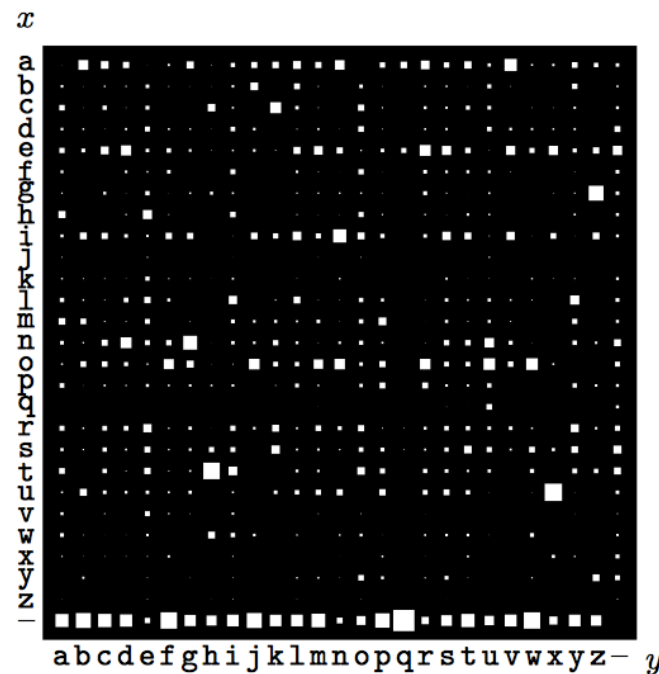


(b)  $P(x|y)$

$$\begin{aligned}
 P(y | x, \mathcal{H}) &= \frac{P(x | y, \mathcal{H})P(y | \mathcal{H})}{P(x | \mathcal{H})} \\
 &= \frac{P(x | y, \mathcal{H})P(y | \mathcal{H})}{\sum_{y'} P(x | y', \mathcal{H})P(y' | \mathcal{H})}.
 \end{aligned}$$



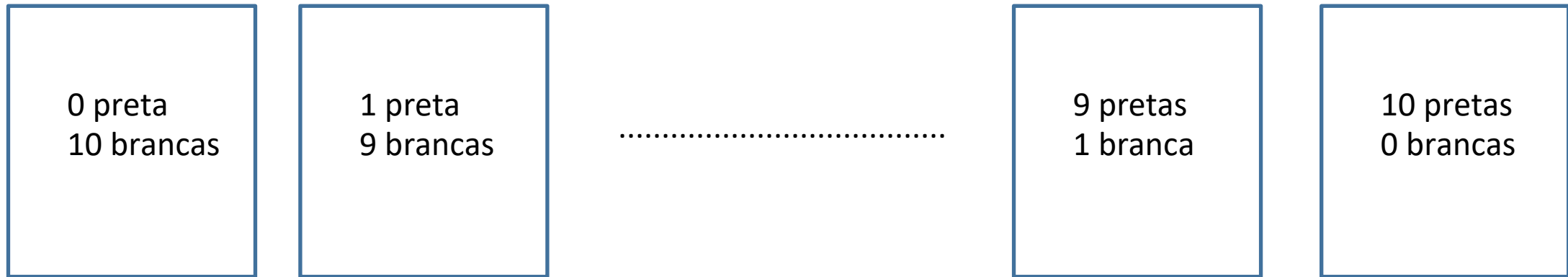
(a)  $P(y|x)$



(b)  $P(x|y)$



# Previsão Bayesiana



Urna  $u$  contém  $u$  bolas pretas e  $10-u$  bolas brancas

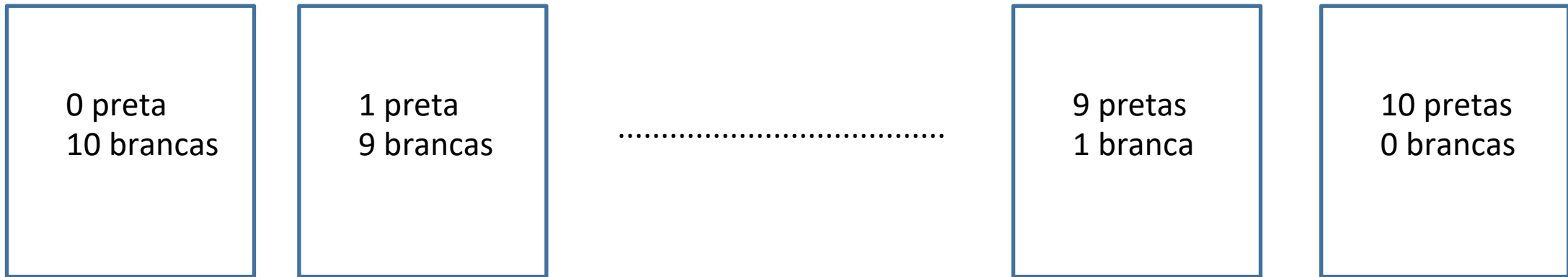
Carlos seleciona uma urna ao acaso e retira  $N$  vezes com reposição

$N=10$  retiradas e  $n_b = 3$  pretas foram observadas

Qual a probabilidade que Carlos esteja usando a urna  $u$  ?



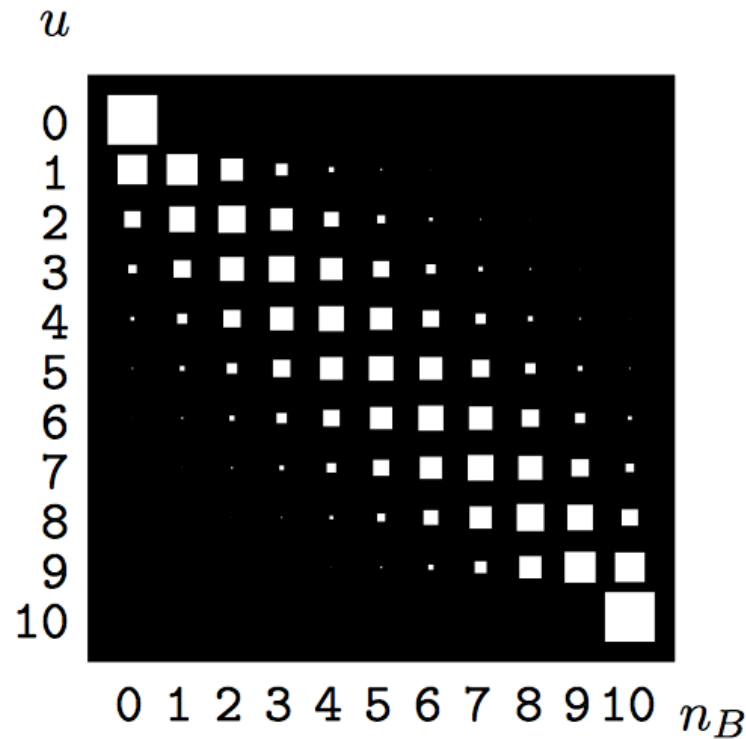
# Inferência



$N=10$  retiradas e  $n_b = 3$  pretas foram observadas

$$\begin{aligned} P(u | n_B, N) &= \frac{P(u, n_B | N)}{P(n_B | N)} \\ &= \frac{P(n_B | u, N)P(u)}{P(n_B | N)} \end{aligned}$$

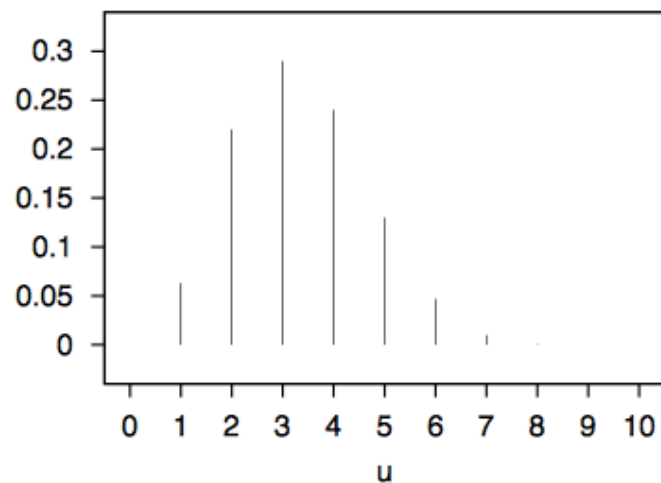
*Você não pode fazer inferência sem fazer hipóteses*



$$P(u) = \frac{1}{11} \quad \text{Conhecimento a priori}$$

$$P(n_B | u, N) = \binom{N}{n_B} f_u^{n_B} (1 - f_u)^{N - n_B}$$

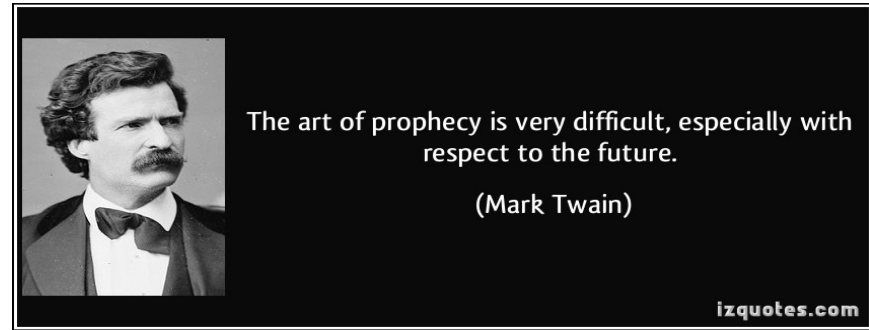
$$P(n_B | N) = \sum_u P(u, n_B | N) = \sum_u P(u) P(n_B | u, N)$$



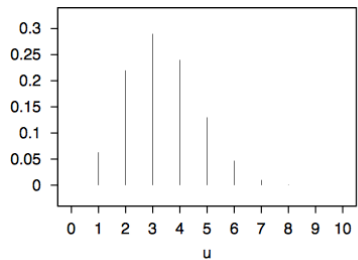
$u$	$P(u   n_B = 3, N)$
0	0
1	0.063
2	0.22
3	0.29
4	0.24
5	0.13
6	0.047
7	0.0099
8	0.00086
9	0.0000096
10	0

$$\begin{aligned}
 P(u | n_B, N) &= \frac{P(u)P(n_B | u, N)}{P(n_B | N)} \\
 &= \frac{1}{P(n_B | N)} \frac{1}{11} \binom{N}{n_B} f_u^{n_B} (1 - f_u)^{N - n_B}
 \end{aligned}$$

# Prediction



$$P(\text{ball}_{N+1} \text{ is black} \mid n_B, N) = \sum_u P(\text{ball}_{N+1} \text{ is black} \mid u, n_B, N)P(u \mid n_B, N)$$



$u$	$P(u \mid n_B = 3, N)$
0	0
1	0.063
2	0.22
3	0.29
4	0.24
5	0.13
6	0.047
7	0.0099
8	0.00086
9	0.0000096
10	0

$$P(\text{ball}_{N+1} \text{ is black} \mid n_B = 3, N = 10) = 0.333$$

# Previsão Bayesiana

$$p(y_* | \mathcal{D}) = \int p(y_* | \mathcal{D}, \theta) \underbrace{p(\theta | \mathcal{D})}_{\text{posterior}} d\theta$$

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{\int p(\mathcal{D} | \theta') p(\theta') d\theta'}$$

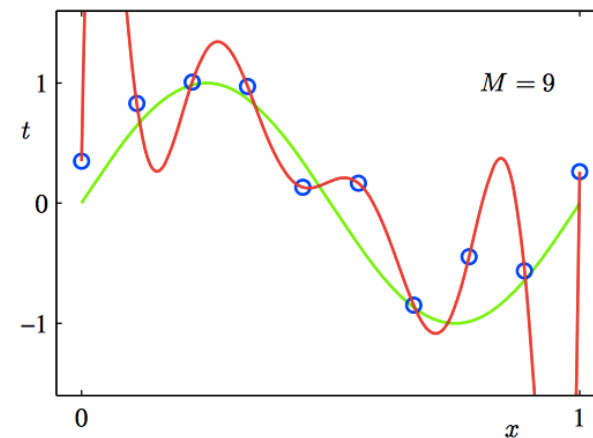
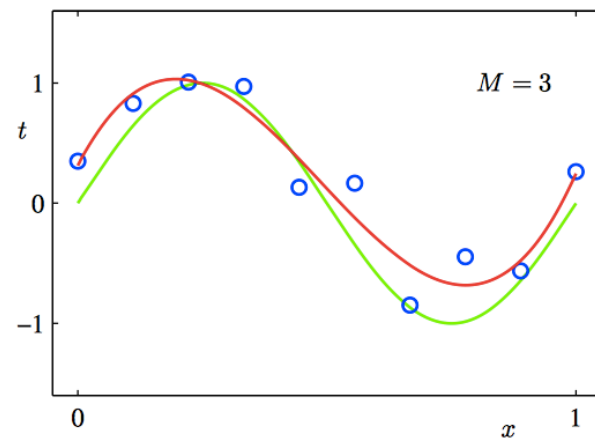
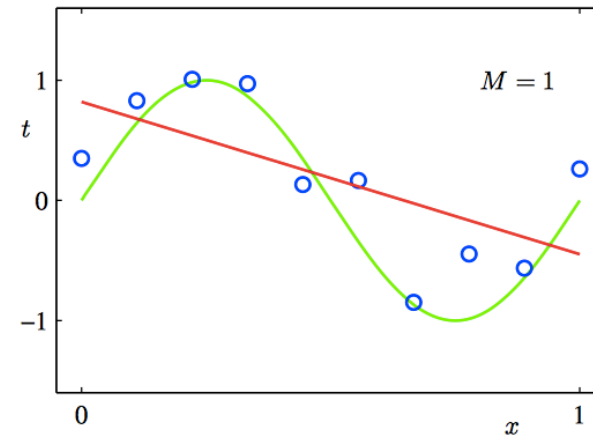
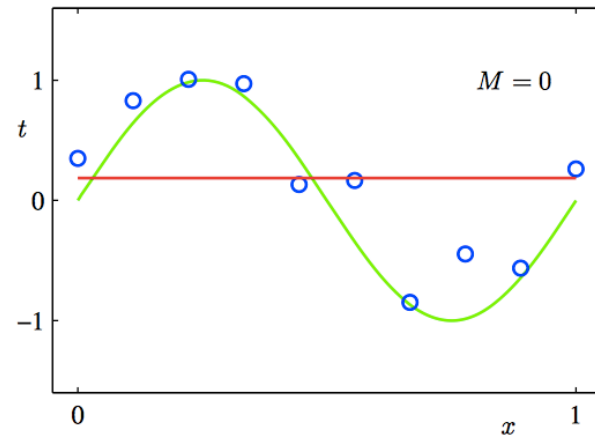
Bayes



$$p(y_* | \mathcal{D}) = \frac{\int p(y_* | \mathcal{D}, \theta) p(\mathcal{D} | \theta) p(\theta) d\theta}{\int p(\mathcal{D} | \theta') p(\theta') d\theta'}$$

Custo exponencial com o número de parâmetros

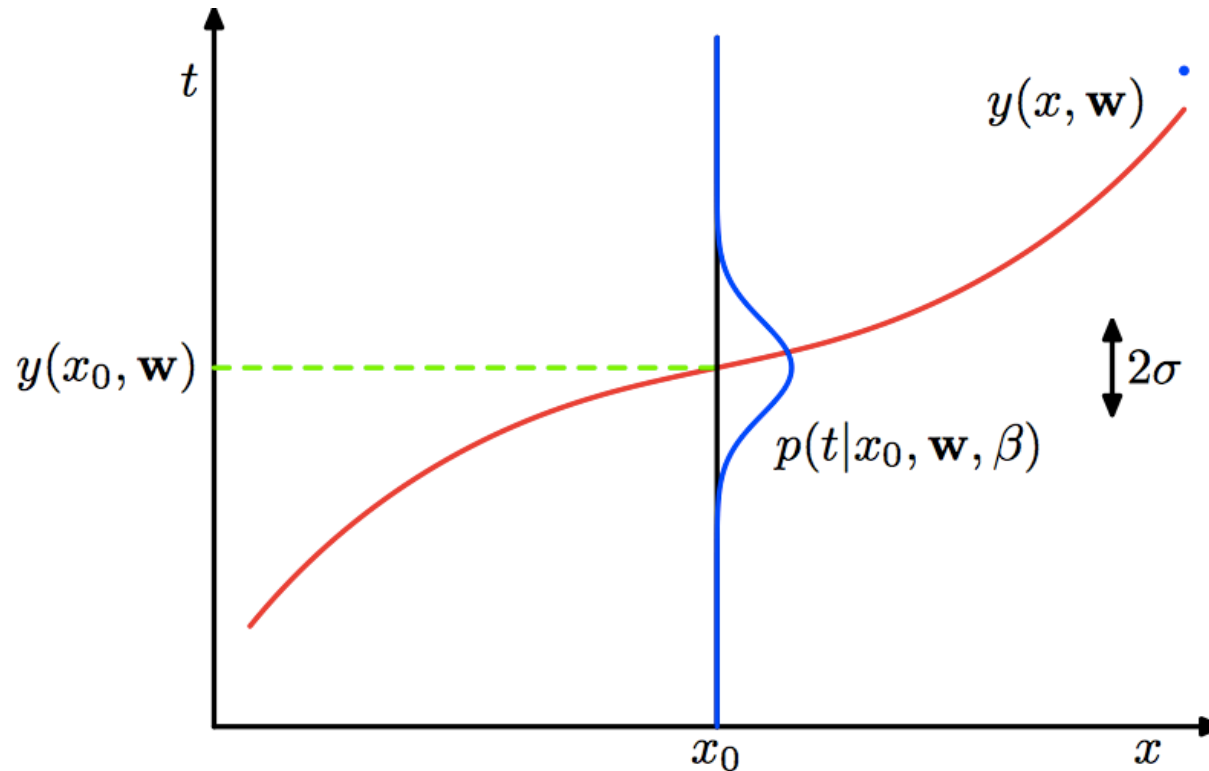
# Interpolação Polinomial – O Retorno



# A abordagem bayesiana

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$$

Modelo de Erro



$N$  input values  $\mathbf{x} = (x_1, \dots, x_N)^T$

$\mathbf{t}$  target values  $\mathbf{t} = (t_1, \dots, t_N)^T$

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y(x_n, \mathbf{w}), \beta^{-1}) \quad \text{Independência dos Erros}$$

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi).$$

Minimizando

$\mathbf{w}_{\text{ML}}$

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t | y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1}).$$



# A abordagem bayesiana completa

Priori sobre os coeficientes

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha).$$

MAP (Maximum a posteriori)

Equivale a  
Minimizar:

$$\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}. \quad \lambda = \frac{\alpha}{\beta}$$

# Predição Bayesiana

*(Tome médias levando em conta toda a incerteza)*

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t}) d\mathbf{w}.$$



$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

# Predição Bayesiana

$$p(t|x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t|m(x), s^2(x))$$

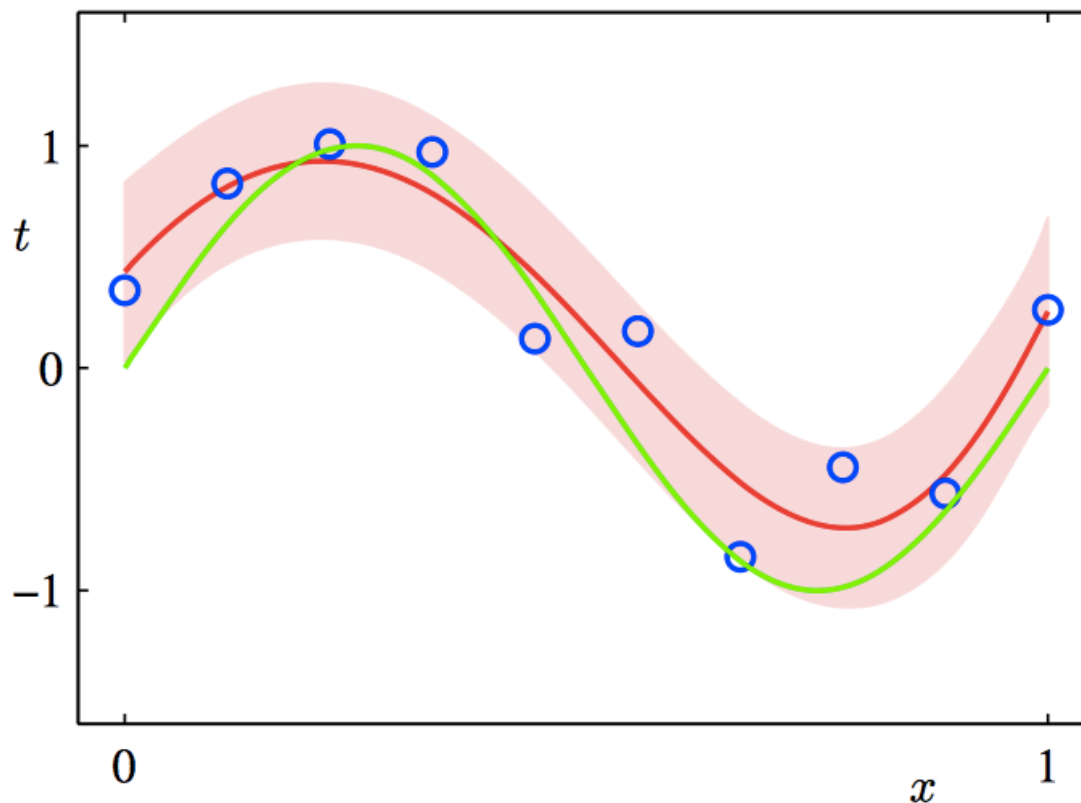
$$\alpha = 5 \times 10^{-3}$$

$$\beta = 11.1$$

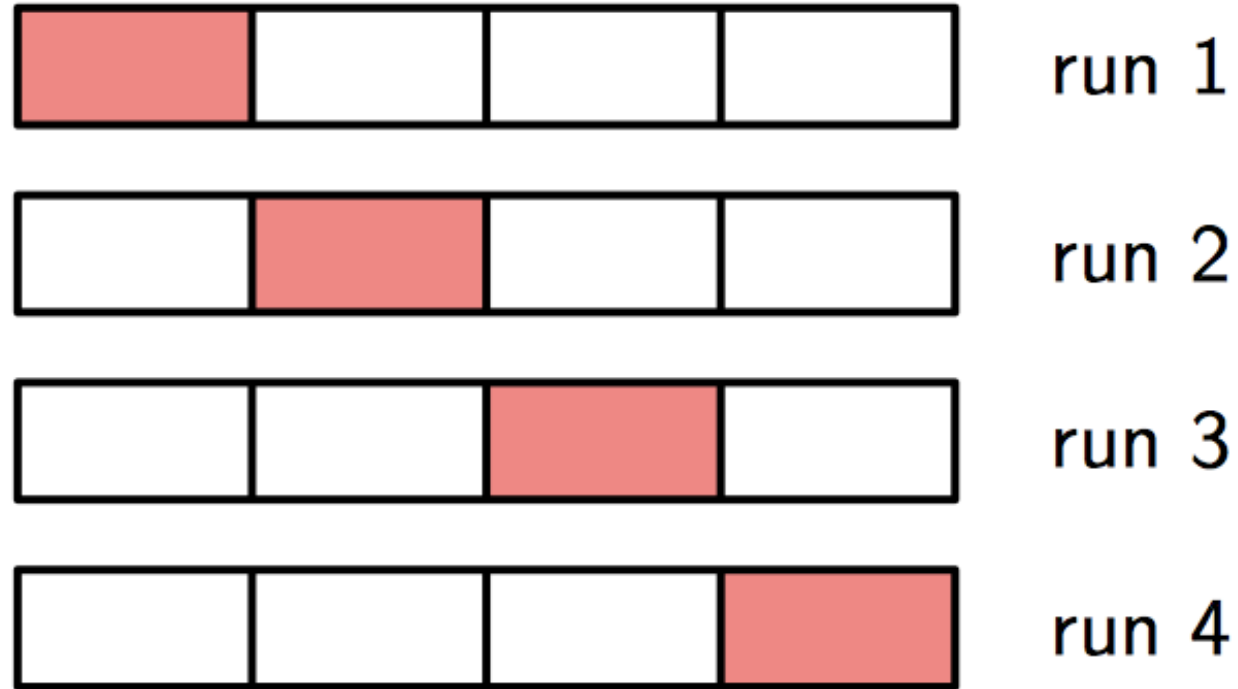
Como escolher ?



INCERTEZA DA OBSERVAÇÃO  
+  
INCERTEZA NOS PARÂMETROS  
DO MODELO



# Cross-validation and Model Choice



# Posterior, likelihood, prior, evidence

$$P(\boldsymbol{\theta} | D, \mathcal{H}) = \frac{P(D | \boldsymbol{\theta}, \mathcal{H})P(\boldsymbol{\theta} | \mathcal{H})}{P(D | \mathcal{H})}$$

$$\text{posterior} = \frac{\text{likelihood} \times \text{prior}}{\text{evidence}}.$$

## A Navalha de Occam

Quais os dois próximos números ?

-1, 3, 7, 11, ?, ?

15, 19 ? (somar 4)

Ou  
-19.9, 1043.8 ?

$-x^3/11 + 9/11x^2 + 23/11$  ?

"A theory with mathematical beauty is more likely to be correct than an ugly one that fits some experimental data" – Paul Dirac

# A Navalha de Occam

$\mathcal{H}_a$  – the sequence is an *arithmetic* progression, ‘add  $n$ ’, where  $n$  is an integer.

~~$\mathcal{H}_c$  – the sequence is generated by a *cubic* function of the form  $x \rightarrow cx^3 + dx^2 + e$ , where  $c$ ,  $d$  and  $e$  are fractions.~~

**OCCAM'S  
RAZOR**

*A Parsimonious  
Shave Every  
Time!*



**William of Occam**

1287 – 1347

*Entia non sunt multiplicanda sine necessitate*

# Lançamento de moedas – Seleção de Modelos

Observation of a sequence of the letters  $a$  and  $b$

The probability, given  $p_a$ , that  $F$  tosses result in a sequence  $s$  that contains  $\{F_a, F_b\}$  counts of the two outcomes

Probability of string  $s$

$$P(\mathbf{s} | p_a, F, \mathcal{H}_1) = p_a^{F_a} (1 - p_a)^{F_b}.$$

[For example,  $P(\mathbf{s} = \text{aaba} | p_a, F = 4, \mathcal{H}_1) = p_a p_a (1 - p_a) p_a.$ ]

Priori

$$P(p_a | \mathcal{H}_1) = 1, \quad p_a \in [0, 1]$$

and  $p_b \equiv 1 - p_a.$



# Probabilidade posteriori de $p_a$

$$P(p_a | \mathbf{s}, F, \mathcal{H}_1) = \frac{P(\mathbf{s} | p_a, F, \mathcal{H}_1)P(p_a | \mathcal{H}_1)}{P(\mathbf{s} | F, \mathcal{H}_1)} \quad \text{BAYES}$$



$$P(p_a | \mathbf{s}, F, \mathcal{H}_1) = \frac{p_a^{F_a}(1 - p_a)^{F_b}}{P(\mathbf{s} | F, \mathcal{H}_1)}$$

Constante  
Normalizadora  
= Evidência

$$P(\mathbf{s} | F, \mathcal{H}_1) = \int_0^1 dp_a p_a^{F_a}(1 - p_a)^{F_b} = \frac{\Gamma(F_a + 1)\Gamma(F_b + 1)}{\Gamma(F_a + F_b + 2)} = \frac{F_a!F_b!}{(F_a + F_b + 1)!}$$

# Do ajuste para a previsão

$$P(\mathbf{a} | \mathbf{s}, F) = \int dp_{\mathbf{a}} P(\mathbf{a} | p_{\mathbf{a}}) P(p_{\mathbf{a}} | \mathbf{s}, F)$$

$$\begin{aligned} P(\mathbf{a} | \mathbf{s}, F) &= \int dp_{\mathbf{a}} p_{\mathbf{a}} \frac{p_{\mathbf{a}}^{F_{\mathbf{a}}} (1 - p_{\mathbf{a}})^{F_{\mathbf{b}}}}{P(\mathbf{s} | F)} \\ &= \int dp_{\mathbf{a}} \frac{p_{\mathbf{a}}^{F_{\mathbf{a}}+1} (1 - p_{\mathbf{a}})^{F_{\mathbf{b}}}}{P(\mathbf{s} | F)} \\ &= \left[ \frac{(F_{\mathbf{a}} + 1)! F_{\mathbf{b}}!}{(F_{\mathbf{a}} + F_{\mathbf{b}} + 2)!} \right] / \left[ \frac{F_{\mathbf{a}}! F_{\mathbf{b}}!}{(F_{\mathbf{a}} + F_{\mathbf{b}} + 1)!} \right] = \frac{F_{\mathbf{a}} + 1}{F_{\mathbf{a}} + F_{\mathbf{b}} + 2} \end{aligned}$$

REGRA DE  
LAPLACE

# Seleção de modelos

Outro cientista introduz um outro modelo:

Hypothesis  $H_0$ : probability of  $a$  is  $p_0 = \frac{1}{6}$



$$P(\mathcal{H}_1 | \mathbf{s}, F) = \frac{P(\mathbf{s} | F, \mathcal{H}_1)P(\mathcal{H}_1)}{P(\mathbf{s} | F)}.$$

Como selecionar o modelo ?

vs.

$$P(\mathcal{H}_0 | \mathbf{s}, F) = \frac{P(\mathbf{s} | F, \mathcal{H}_0)P(\mathcal{H}_0)}{P(\mathbf{s} | F)}$$

$$P(\mathbf{s} | F) = P(\mathbf{s} | F, \mathcal{H}_1)P(\mathcal{H}_1) + P(\mathbf{s} | F, \mathcal{H}_0)P(\mathcal{H}_0)$$

# Seleção de modelos

Defining  $p_0$  to be  $1/6$ , we have

$$P(\mathbf{s} | F, \mathcal{H}_0) = p_0^{F_a} (1 - p_0)^{F_b} \quad \underline{= \text{Evidência}}$$

Thus the posterior probability ratio of model  $\mathcal{H}_1$  to model  $\mathcal{H}_0$  is

$$\begin{aligned} \frac{P(\mathcal{H}_1 | \mathbf{s}, F)}{P(\mathcal{H}_0 | \mathbf{s}, F)} &= \frac{P(\mathbf{s} | F, \mathcal{H}_1)P(\mathcal{H}_1)}{P(\mathbf{s} | F, \mathcal{H}_0)P(\mathcal{H}_0)} \\ &= \frac{F_a! F_b!}{(F_a + F_b + 1)!} \bigg/ p_0^{F_a} (1 - p_0)^{F_b}. \end{aligned}$$

**How model comparison works:** The evidence for a model is usually the normalizing constant of an earlier Bayesian inference.

# Comparando modelos

$F$	Data ( $F_a, F_b$ )	$\frac{P(\mathcal{H}_1   \mathbf{s}, F)}{P(\mathcal{H}_0   \mathbf{s}, F)}$	
6	(5, 1)	222.2	
6	(3, 3)	2.67	
6	(2, 4)	0.71	= 1/1.4
6	(1, 5)	0.356	= 1/2.8
6	(0, 6)	0.427	= 1/2.3
20	(10, 10)	96.5	
20	(3, 17)	0.2	= 1/5
20	(0, 20)	1.83	

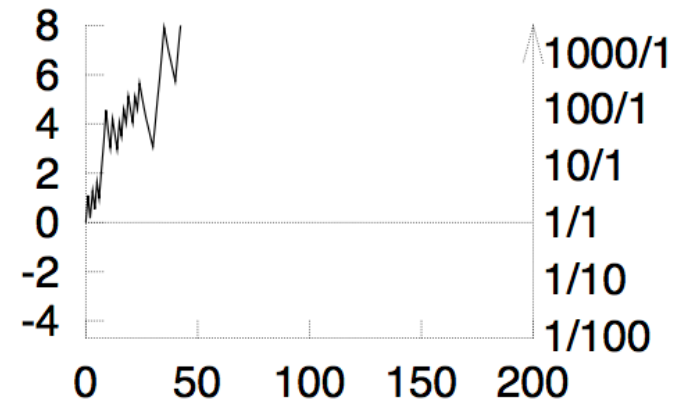
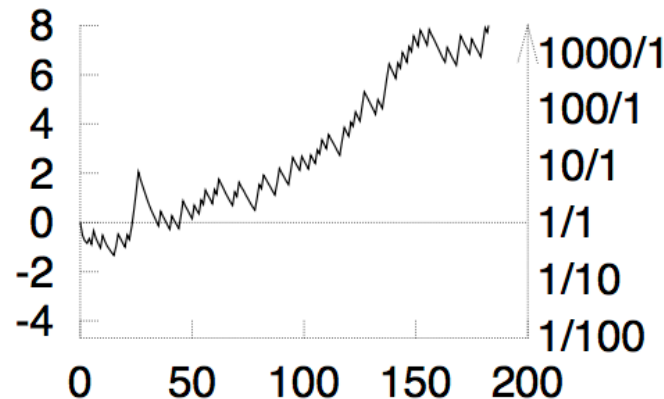
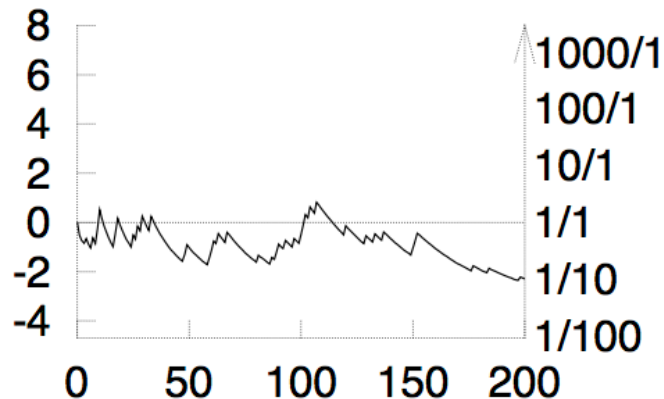
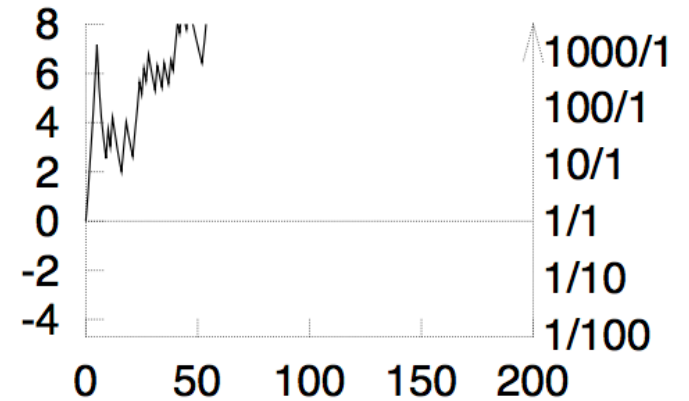
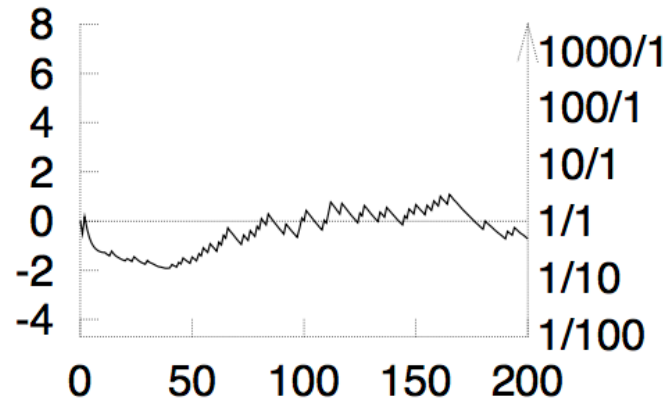
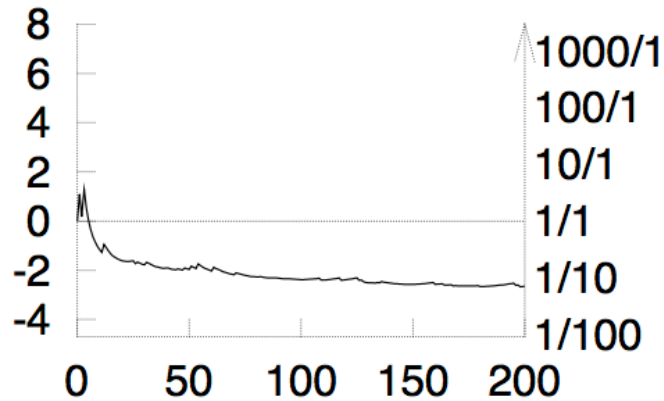
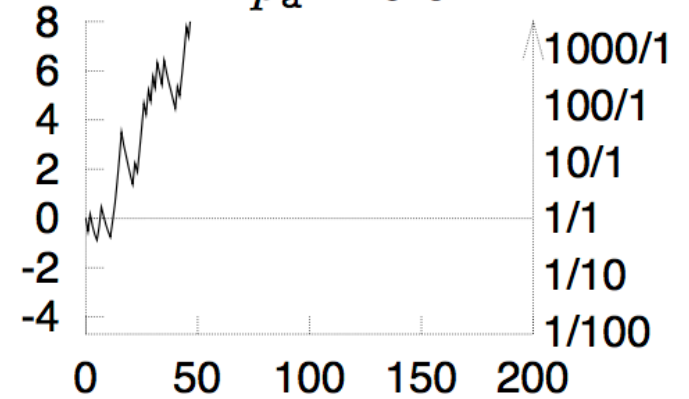
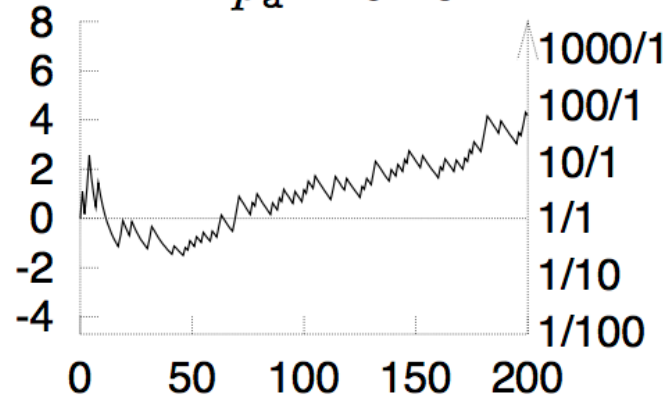
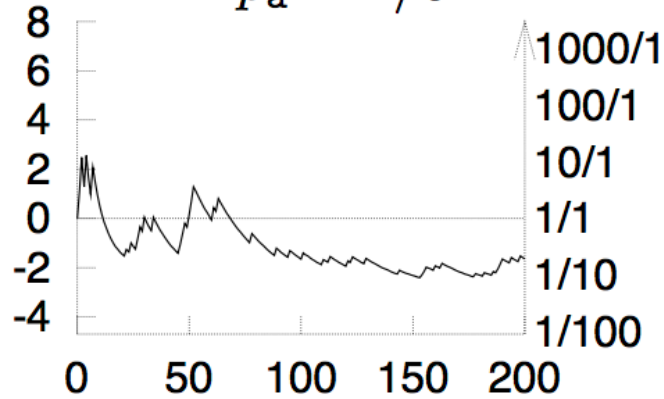
$\mathcal{H}_0$  is true

$\mathcal{H}_1$  is true

$p_a = 1/6$

$p_a = 0.25$

$p_a = 0.5$



# Navalha de Occam – Ajuste de modelo

## 1 – First level of inference

$$P(\mathbf{w}|D, \mathcal{H}_i) = \frac{P(D|\mathbf{w}, \mathcal{H}_i)P(\mathbf{w}|\mathcal{H}_i)}{P(D|\mathcal{H}_i)}$$

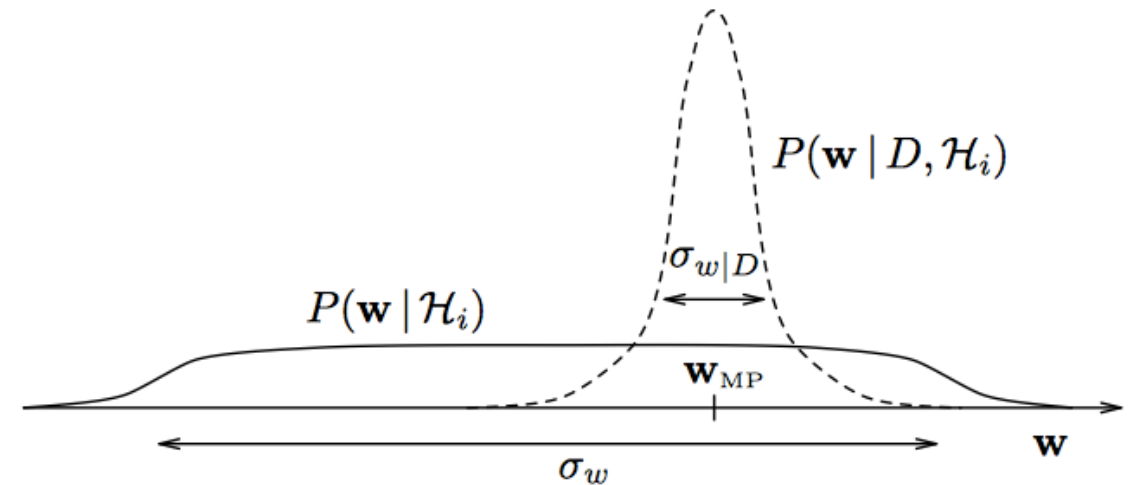
$$\text{Posterior} = \frac{\text{Likelihood} \times \text{Prior}}{\text{Evidence}}$$

# Navalha de Occam – Comparação de modelos

## 2 – Second level of inference

$$P(\mathcal{H}_i | D) \propto P(D | \mathcal{H}_i) P(\mathcal{H}_i).$$

$$P(D | \mathcal{H}_i) = \int P(D | \mathbf{w}, \mathcal{H}_i) P(\mathbf{w} | \mathcal{H}_i) d\mathbf{w}$$



For many problems the posterior  $P(\mathbf{w} | D, \mathcal{H}_i) \propto P(D | \mathbf{w}, \mathcal{H}_i) P(\mathbf{w} | \mathcal{H}_i)$  has a strong peak at the most probable parameters  $\mathbf{w}_{MP}$

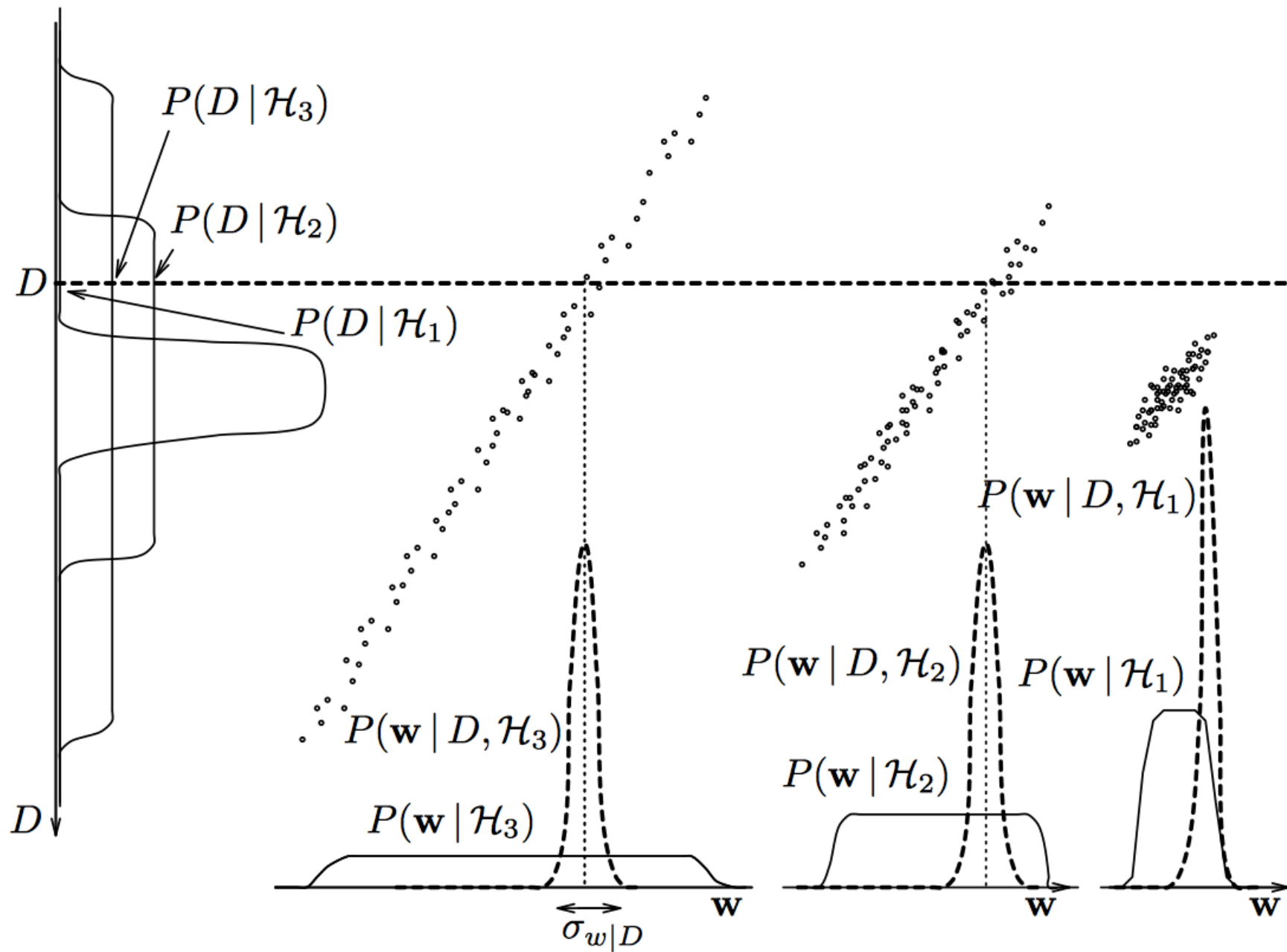
$$P(\mathbf{w}_{MP} | \mathcal{H}_i) = 1/\sigma_w.$$

$$P(D | \mathcal{H}_i) \simeq \underbrace{P(D | \mathbf{w}_{MP}, \mathcal{H}_i)}_{\text{Best fit likelihood}} \times \underbrace{P(\mathbf{w}_{MP} | \mathcal{H}_i) \sigma_{w|D}}_{\text{Occam factor}}$$

$$\text{Evidence} \simeq \text{Best fit likelihood} \times \text{Occam factor}$$

$$\text{Occam factor} = \frac{\sigma_{w|D}}{\sigma_w}$$





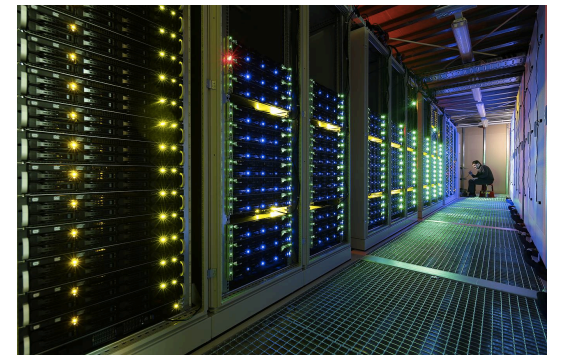
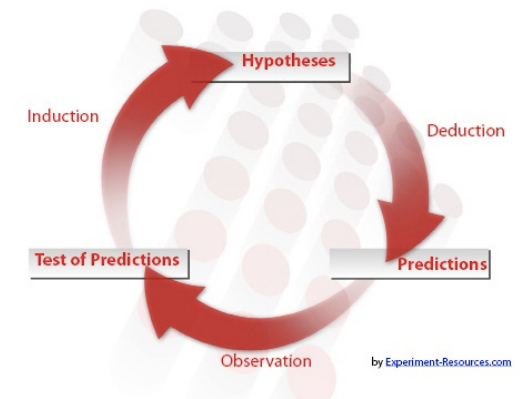
# Probability theory reaches parts that ad hoc methods cannot reach

$$P(\mathbf{t} | D, I) = \sum_{\mathcal{H}} P(\mathbf{t} | D, \mathcal{H}, I) P(\mathcal{H} | D, I). \quad \text{Previs\~{a}o Bayesiana}$$

$$P(\mathcal{H} | D, I) = \frac{P(D | \mathcal{H}, I) P(\mathcal{H} | I)}{P(D | I)}$$

$$P(\lambda | D, \mathcal{H}) = \frac{P(D | \lambda, \mathcal{H}) P(\lambda | \mathcal{H})}{P(D | \mathcal{H})}$$

EVIDENCE



# Machine Learning: Treino, Teste, Previsão e Ação

