

1.

a) O modelo é  $y = a + bt^2$ , o dataset fica:

$t$	$t^2$	$y$
5	25	722
10	100	1073
15	225	1178
20	400	1177

Using ridge regression:  $(X^T X + \lambda I) w = X^T y$

$$X = \begin{bmatrix} 25 \\ 100 \\ 225 \\ 400 \end{bmatrix} \quad w = \begin{bmatrix} a \\ b \end{bmatrix}, \quad y = \begin{bmatrix} 722 \\ 1073 \\ 1178 \\ 1177 \end{bmatrix}$$

$$\lambda = 0.5$$

Introduzindo o termo bias:

$$X = \begin{bmatrix} 1 & 25 \\ 1 & 100 \\ 1 & 225 \\ 1 & 400 \end{bmatrix}$$

Logo

$$\left( \begin{bmatrix} 1 & 1 & 1 & 1 \\ 25 & 100 & 225 & 400 \end{bmatrix} \begin{bmatrix} 1 & 25 \\ 1 & 100 \\ 1 & 225 \\ 1 & 400 \end{bmatrix} + 0.5 I \right) w = X^T y$$

$2 \times 4$                        $4 \times 2$

$$\begin{pmatrix} \begin{bmatrix} 4 & 750 \\ 750 & 221250 \end{bmatrix} + 0,5 \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{pmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 9 \\ 5 \end{bmatrix}$$

2x2

2x1

$$\begin{bmatrix} 4,5 & 750 \\ 750 & 221250,5 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 25 & 100 & 225 & 400 \end{bmatrix} \begin{bmatrix} 722 \\ 1073 \\ 1178 \\ 1177 \end{bmatrix}$$

2x2

2x1

2x4

4x1

$$\begin{bmatrix} 4,5a + 750b \\ 750a + 221250,5b \end{bmatrix} = \begin{bmatrix} 4150 \\ 861200 \end{bmatrix}$$

2x1

2x1

Convertendo em sistema de Eqs:

$$\begin{cases} 4,5a + 750b = 4150 \\ 750a + 221250,5b = 861200 \end{cases} \quad \Rightarrow \quad \begin{cases} a = \frac{4150 - 750b}{4,5} \\ 750 \left( \frac{4150 - 750b}{4,5} \right) + 221250,5b = 861200 \end{cases}$$

$$\quad \Rightarrow \quad \frac{3112500}{4,5} - \frac{562500b}{4,5} + 221250,5b = 861200$$

$$(2) \left\{ \begin{array}{l} \text{---} \\ (221250,5 - \frac{562500}{4,5})b = 861200 - \frac{3912500}{4,5} \end{array} \right.$$

$$(3) \left\{ \begin{array}{l} \text{---} \\ 96250,5b = \frac{508600}{3} \end{array} \right. \Rightarrow \left\{ \begin{array}{l} \text{---} \\ b = \frac{508600}{3} \times \frac{1}{96250,5} \end{array} \right. \Rightarrow$$

$$\left\{ \begin{array}{l} a = \frac{4150 - 250(1,76)}{4,5} \\ b = 1,76 \end{array} \right. \Rightarrow$$

$$\left\{ \begin{array}{l} a = 628,89 \\ b = 1,76 \end{array} \right.$$

$$\text{Modelo } \boxed{y = 628,89 + 1,76t^2}$$

$$b) y_{\text{related}} = 600 + 2t^2, \text{ logo } w = \begin{bmatrix} 600 \\ 2 \end{bmatrix}$$

Since we have initialized weights, that could help us now to define, the learning task should be defined not as Ridge Regression but a standard Regression that could be executed with Least-Square update or ~~least~~ steepest descent. so,

$$\text{for LR: } w^{(t+1)} = w^{(t)} + (y_i - w^T x_i) x_i$$

$$\text{for steepest descent: } w^{(t+1)} = w^{(t)} + \alpha \sum_n (y_n - w^T x_n) x_n$$

where  $w^{(t)} = \begin{bmatrix} 600 \\ 2 \end{bmatrix}$ , and  $\alpha$  is a learning rate



2.

$$a) \quad p(x|y_0) = \frac{e^{2\cos(x-1)}}{2\pi(2.2796)}, \quad \text{if } x=0 \quad \left\{ \begin{array}{l} p(x|y_0) = 0,206 \\ p(x|y_1) = 0,210 \end{array} \right.$$

$$p(x|y_1) = \frac{e^{3\cos(x+0,9)}}{2\pi(4.8808)}$$

como  $p(x|y_1) > p(x|y_0)$ , a prediction sera  $y_1$ .

b). Para a prediction ser  $y_0$ :

$$p(x|y_0) > p(x|y_1), \text{ então:}$$

$$\frac{\frac{e^{2\cos(x-1)}}{2\pi(2.2796)}}{2\pi(2.2796)} > \frac{\frac{e^{3\cos(x-1)}}{2\pi(4.8808)}}{2\pi(4.8808)}$$

$$\ln\left(\frac{e^{2\cos(x-1)}}{2\pi(2.2796)}\right) > \ln\left(\frac{e^{3\cos(x-1)}}{2\pi(4.8808)}\right)$$

$$\ln(e^{2\cos(x-1)}) - \ln(2\pi(2.2796)) > \ln(e^{3\cos(x-1)}) - \ln(2\pi(4.8808))$$

$$2\cos(x-1) - 3\cos(x-1) > \ln(2\pi(2.2796)) - \ln(2\pi(4.8808))$$

$$-\cos(x-1) < -(\ln(2\pi(2.2796)) - \ln(2\pi(4.8808)))$$

$$x-1 < \cos^{-1}(0,761)$$

$$x < \cos^{-1}(0,761) + 1 \quad \text{ou} \quad x < 1,76 \text{ rad}$$

hence, the predictor is equal to  $y_0 = [0, 176]$  and.

c) we can use the Bayes Formula to convert the prior probability to the posterior probability:

$$P(C_j | x) = \frac{P(x | C_j) P(C_j)}{P(x)}$$

$$\text{where } P(x) = \sum_{j=1}^2 P(x | C_j) P(C_j)$$

we want:

$$P(y_0 | x) = \frac{P(x | y_0) P(y_0)}{P(x)}$$

let's assume

$$P(y_0) = P(y_1)$$

$$P(y_0) = P(y_1) = 0.5$$

$$P(x) = P(x | y_0) P(y_0) + P(x | y_1) P(y_1)$$

$$= 0.5 (P(x | y_0) + P(x | y_1))$$

so

$$P(y_0 | x) = \frac{P(x | y_0) P(y_0)}{0.5 (P(x | y_0) + P(x | y_1))} = \frac{P(x | y_0) \cancel{0.5}}{\cancel{0.5} (P(x | y_0) + P(x | y_1))}$$

(A)

$$= \frac{P(x | y_0)}{P(x | y_0) + P(x | y_1)} = \frac{1}{\frac{P(x | y_0) + P(x | y_1)}{P(x | y_0)}} = \frac{1}{1 + P(x | y_1) / P(x | y_0)}$$

$$p(x|y_1) = \frac{e^{k_1 \cos(x - \mu_1)}}{2\pi I(k_1)} = \frac{e^{k_1 \cos(x - \mu_1)}}{e^{\ln(2\pi I(k_1))}}$$

$$= \frac{k_1 \cos(x - \mu_1) - \ln(2\pi I(k_1))}{1}$$

$$(2) = \boxed{\frac{-\ln(2\pi I(k_1)) + k_1 \cos(x - \mu_1)}{1}}$$

untuk:

$$\cos(x) = \sin\left(x + \frac{\pi}{2}\right)$$

$$\text{Jika } x = x+1$$

$$\cos(x+1) = \sin\left(x+1 + \frac{\pi}{2}\right)$$

$$\text{Jika } x = x - \mu_1$$

$$\cos(x - \mu_1) = \sin\left(x - \mu_1 + \frac{\pi}{2}\right)$$

Logo:

$$(1) = \frac{-\ln(2\pi I(k_1)) + k_1 \sin\left(x - \mu_1 + \frac{\pi}{2}\right)}{1}$$

$$\text{Turunan: } -\ln(2\pi I(k_1)) = w_0$$

$$k_1 = w_1$$

$$\theta = \mu_1 + \frac{\pi}{2}$$

Logo, fungsi (1) & (2)

$$p(y|x) = \frac{1}{1 + e^{w_0 + w_1 \sin(x - \theta)}}$$

Gg



3. see midterm\_sol.pdf

a)

(a) and (b) are appropriate to use in classification. In (c), there is very little penalty for extremely misclassified examples, which correspond to very negative  $y F(x)$ . In (d) and (e), correctly classified examples are penalized, whereas misclassified examples are not. In general,  $L$  should approximate the 0-1 loss, and it should be a non-increasing function of  $y F(x)$ .

b) Function (b) is more robust to outliers. For outliers,  $y F(x)$  is often very negative. In (a), outliers are heavily penalized. So, the resulting classifier is largely affected by the outliers. On the other hand, in (b), the loss of outliers is bounded. So, the resulting classifier is less affected by the outliers, and thus more robust.

c) To obtain the parameters in  $F(x)$ , we need to ~~maximize~~ minimize  $\sum_i L(y^i F(x^i)) = \sum_i \frac{1}{1 + \exp(y^i F(x^i))}$

$$= \sum_i \frac{1}{1 + \exp(y^i (w_0 + \sum_{j=1}^d w_j x_j^i))}$$

gradient descent, logit:

$$\frac{\partial}{\partial w_0} \sum_i L(y^i F(x^i)) = \sum_i \frac{\partial}{\partial w_0} \left( \frac{1}{1 + \exp(y^i (w_0 + \sum_{j=1}^d w_j x_j^i))} \right)$$

$$= - \sum_i \frac{y^i \exp(y^i F(x^i))}{(1 + \exp(y^i F(x^i)))^2}$$

Again,  $\forall k = 1, \dots, d$

$$\frac{\partial}{\partial w_k} \sum_i L(y^i F(x^i)) = \sum_i \frac{\partial}{\partial w_k} \left( \frac{1}{1 + \exp(y^i (w_0 + \sum_{j=1}^d w_j x_j^i))} \right)$$

$$= - \sum_i \frac{y^i x_k^i \exp(y^i F(x^i))}{(1 + \exp(y^i F(x^i)))^2}$$

Therefore, the update rules are as follows:

$$w_0^{(t+1)} = w_0^{(t)} - \mu \frac{\partial}{\partial w_0} \sum_i L(y^i F(x^i))$$

$$w_k^{(t+1)} = w_k^{(t)} + \mu \sum_i \frac{y^i x_k^i \exp(y^i F(x^i))}{(1 + \exp(y^i F(x^i)))^2}$$



Para  $k = 1, \dots, d$

$$w_k^{(t+1)} = w_k^{(t)} - \mu \frac{\partial}{\partial w_k} \sum_i L(y^i F(x^i))$$

$$w_k^{(t+1)} = w_k^{(t)} + \mu \sum_i \frac{y^i x_k^i \exp(y^i F(x^i))}{(1 + \exp(y^i F(x^i)))^2}$$

4.

$$P(A) = \frac{1}{2} ; P(B) = \mu ; P(C) = 3\mu ; P(D) = \frac{1}{2} - 4\mu$$

c students get a "(C)"

d students get a "(D)"

h students get a "(A)" or "(B)"

Total der

$$c + d + h$$

$$\boxed{c + d + a + b}$$

$$a + b = h$$

EM algorithm to obtain  $\mu$  : See final 2008 - solution.pdf

a)

$$\hat{a} = \frac{P(A)}{P(A) + P(B)} h ; \hat{b} = \frac{P(B)}{P(A) + P(B)} h$$

$$\hat{a} = \frac{\frac{1}{2}}{\frac{1}{2} + \mu} h$$

$$\hat{b} = \frac{\mu}{\frac{1}{2} + \mu} h$$

b)

$$\hat{u} = \frac{h-a+c}{6(h-a+c+d)} \quad \hat{u} = \frac{h-a+c}{(3 \times 4)(h-a+c+d)} \quad \boxed{\frac{h-a+c}{12(h-a+c+d)}}$$

See example on:

2008f - solution.pdf

5. HMM (see midten - solution.pdf)

a) hidden prob. matrices:

Transition matrix:

$$\begin{matrix} & S & V \\ S & \begin{bmatrix} 0.8 & 0.2 \end{bmatrix} \\ V & \begin{bmatrix} 0.4 & 0.6 \end{bmatrix} \end{matrix}$$

Emission Matrix:

$$\begin{matrix} & \text{Grin} & \text{Frown} \\ S & \begin{bmatrix} 0.5 \\ 0.8 \end{bmatrix} \\ V & \begin{bmatrix} 0.5 \\ 0.2 \end{bmatrix} \end{matrix}$$

$$\pi_0 = \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \end{bmatrix}$$

Viterbi algorithm: (slide 43 lecture 11)

Input: 6 6 F 6

initialization

$$V_i^k = P(x_i | y_i^k = 1) \pi_k$$

iteration

$$V_t^k = P(x_t | y_t = 1) \max_i a_{i,k} V_{t-1}^i$$

$$P_{tn}(k, t) = \arg \max_i a_{i,k} V_{t-1}^i$$

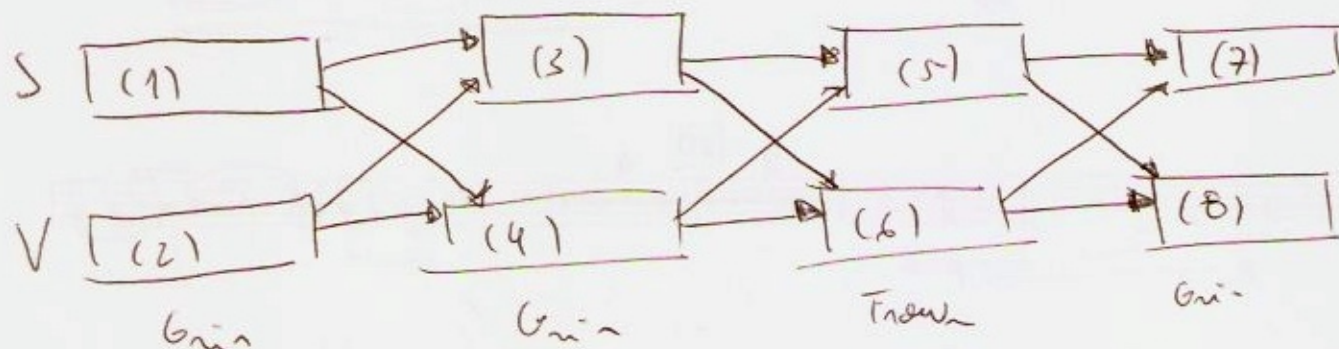
$$\text{True } P(x, y^k) = \max_k \frac{V_T^k}{T}$$



Traceback

$$y_t^* = \operatorname{argmax}_k v_t^k \quad ; \quad y_{t+1}^* = P_{t+1}(y_t^*, t)$$

Redup:



Calculus

$$1. p_{i0}^s \times p(\text{Grin} | s) = \frac{1}{2} \times 0.5 = \frac{1}{2} \times \frac{1}{2} = \boxed{\frac{1}{4}}$$

$$2. p_{i0}^v \times p(\text{Grin} | v) = \frac{1}{2} \times 0.8 = \frac{1}{2} \times \frac{4}{5} = \frac{4}{10} = \boxed{\frac{2}{5}}$$

$$3. p(\text{Grin} | s) \times \max \begin{cases} (1) \times p(s | s) = \frac{1}{4} \times 0.8 = \boxed{0.2} \Rightarrow \text{the max!} \\ (2) \times 0.4 = \frac{2}{5} \times p(v | s) = 0.16 \end{cases}$$

$$\hookrightarrow 0.5 \times 0.2 = \boxed{0.1}$$

$$4. p(\text{Grin} | v) \times \max \begin{cases} (1) \times p(s | v) = \frac{1}{4} \times 0.2 = 0.05 \\ (2) \times p(v | v) = \frac{2}{5} \times 0.6 = \boxed{0.24} \Rightarrow \text{the max!} \end{cases}$$

$$\hookrightarrow 0.8 \times 0.24 = \boxed{0.192}$$

5.

$$P(\text{Frown} | S) \times \max \begin{cases} (3) \times P(S|S) = 0.1 \times 0.8 = \underline{0.08} \text{ max!} \\ (4) \times P(S|V) = 0.192 \times 0.4 = 0.0768 \end{cases}$$

$$0.5 \times 0.08 = \underline{0.04}$$

6.

$$P(\text{Frown} | V) \times \max \begin{cases} (3) \times P(V|S) = 0.1 \times 0.2 = 0.02 \\ (4) \times P(V|V) = 0.192 \times 0.6 = \underline{0.1152} \text{ max!} \end{cases}$$

$$0.2 \times 0.1152 = \underline{0.02304}$$

7.

$$P(\text{Grin} | S) \times \max \begin{cases} (5) \times P(S|S) = 0.04 \times 0.8 = \underline{0.032} \text{ max!} \\ (6) \times P(S|V) = 0.02304 \times 0.4 = \underline{0.009216} \end{cases}$$

$$0.5 \times 0.032 = \underline{0.016}$$

8.

$$P(\text{Grin} | V) \times \max \begin{cases} (5) \times P(V|S) = 0.04 \times 0.2 = 8 \times 10^{-3} \\ (6) \times P(V|V) = 0.02304 \times 0.6 = \underline{0.013824} \text{ max} \end{cases}$$

$$0.8 \times 0.013824 = \underline{0.010592}$$

A probabilidade do sorriso mais provável ter-se em  $\varphi_t = k$ ; logo no último refugo. (7) é mais que (8), então o sorriso mais provável é "SSSS".



5) Para mudar, há 25% de chance:

Nº de ~~S~~ sabendo ~~S~~: 3

Nº de ~~F~~ sabendo S: 8

Nº de Transições S → S: 6

Nº de Transições S → V: 5

Nº de 6 sabendo V: 8

Nº de ~~S~~ sabendo V: 1

Nº de Transições V → S: 4

Nº de Transições V → V: 4

Transition Matrix

$$\begin{matrix} & S & V \\ \begin{matrix} S \\ V \end{matrix} & \begin{bmatrix} \frac{6}{11} & \frac{5}{11} \\ \frac{4}{8} & \frac{4}{8} \end{bmatrix} & = \begin{bmatrix} \frac{6}{11} & \frac{5}{11} \\ \frac{1}{2} & \frac{1}{2} \end{bmatrix} \end{matrix}$$

Apri o matrix

$$\lambda = 0$$

Se for  $\lambda = 0.1$

na 10ª edição

0.1 ao invés  
(de 0.25)

Emission Matrix:

$$\begin{matrix} & 6 & F \\ \begin{matrix} S \\ V \end{matrix} & \begin{bmatrix} \frac{3}{11} & \frac{8}{11} \\ \frac{0}{9} & \frac{1}{9} \end{bmatrix} \end{matrix}$$

tabela

E	$P(E   X = S)$
smm	3/11
tor	8/11

E	$P(E   X = V)$
smm	0/9
tor	1/9

HMM:

