

Least Squares as Orthogonal Projection

Renato Assunção
ESRI Inc. and DCC-UFMG

Dois problemas clássicos em Machine Learning [😊]

- ⊕ Regressão
- ⊕ Classificação

⊕ Regressão: prever $Y =$ preço de ap^{to} a partir de características (features) tais como:

$X_1 =$ área do ap^{to}

$X_2 =$ idade do ap^{to}

$X_3 =$ nº de quartos

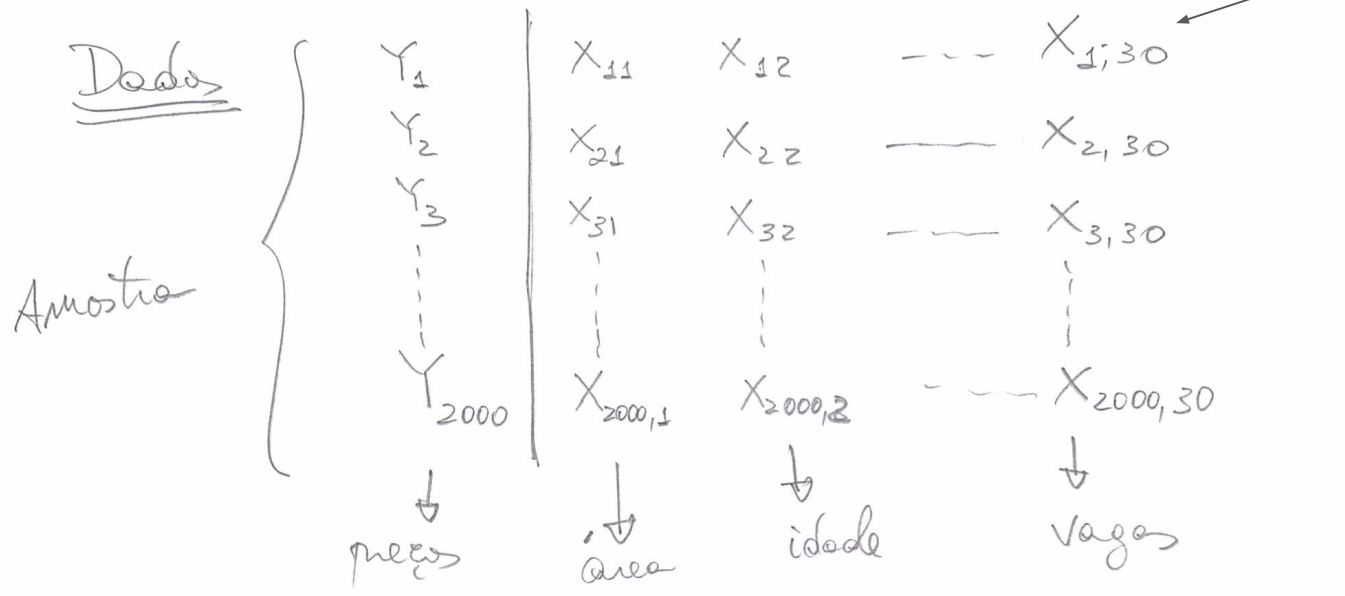
\vdots
 $X_{30} =$ nº vagas de garagem

$$Y = \text{preço do ap}^{\text{to}}$$

$$\tilde{X} = (X_1, X_2, \dots, X_{30}) \in \mathbb{R}^{30}$$

\downarrow area \downarrow idade \downarrow vagas

Numerical variables



Queremos encontrar uma fórmula (função) ⁽³⁾

$$g(\underline{x}) = g(\text{área}, \text{idade}, \dots, \text{vagas})$$

tal que $|Y - g(\underline{x})| \approx 0$

para todo elemento da amostra (todo ap^{to})

⊕ $Y - g(\underline{x}) =$ erro ao prever preço Y
usando $g(\underline{x})$ como preditor

⊕ O que eu realmente quero é (4)
prever preços de ap^{to} que eu
não observei ainda

⊕ Conceito:

⊕ selecione um novo ap^{to} ao
acaso $(Y, \underline{x}) = (Y, x_1, \dots, x_{30})$

⊕ Qual o valor esperado do
erro de previsão $E(|Y - g(\underline{x})|)$

- ⊕ Resultados práticos e teóricos para ⁽³⁾
- o erro ao quadrado:

$$MSE = E \left[\underbrace{(Y - g(x))^2}_{\text{erro de previsão ao quadrado}} \right]$$

↓
Mean
Squared
Error

↓
erro de previsão ao
quadrado

- ⊕ Vários algoritmos: {
- Regressão Linear,
 - SVM, Árvores de regressão
 - Redes Neurais,
-

⊕ Em todos, o objetivo é encontrar (6)
uma função $g(\underline{x})$ que minimize

$$\mathbb{E} \left[(Y - g(\underline{x}))^2 \right]$$

⊕ Minimizar num espaço de funções:

Cálculo



Agora

Qual a função
 $g(\underline{x})$ que minimize
a expressão $\mathbb{E} \left[(Y - g(\underline{x}))^2 \right]$

(7)

⊕ Temos a solução teórica e perfeita

Teorema Dentre as infinitas funções

$g(\underline{x})$, aquela que minimiza

$$\mathbb{E}[(Y - g(\underline{x}))^2]$$

é a esperança condicional

$$\mu(\underline{x}) = \mathbb{E}(Y | \underline{x})$$

(7)

⊕ Temos a solução teórica e perfeita

Teorema Dentre as infinitas funções

$g(\underline{x})$, aquela que minimiza

$$\mathbb{E}[(Y - g(\underline{x}))^2]$$

é a esperança condicional

$$\mu(\underline{x}) = \mathbb{E}(Y | \underline{x})$$

Provamos este resultado na aula passada. Ver notas no website

⊕ Esta solução teórica é pouco útil (8)
pois não sabemos calcular $\mu(\underline{x}) = E(Y|\underline{x})$

⊕ Em geral, não conhecemos a distribuição conjunta das variáveis aleatórias

$$(Y, \underline{x}) = (Y, X_1, X_2, \dots, X_{30})$$

⊕ O que fazemos então?

⊕ Todo modelo/algoritmo de ML estabelece ⁹
uma classe \mathcal{G} (um conjunto) de funções
possíveis $g(\underline{x})$

⊕ Dentro desse subconjunto \mathcal{G} , procura
achar aquela $g^*(\underline{x})$ que seja
próxima da função ótima
perfeita $\mu(\underline{x}) = E(Y | \underline{x})$

⊕ Exemplos :

(10)

Regressão linear

$$\mathcal{L} = \left\{ g(\underline{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_{30} x_{30} \right\}$$

combinação linear
das features

Cada $g(\underline{x}) \in \mathcal{L}$ é determinada pelos

coeficientes $(\theta_0, \theta_1, \dots, \theta_{30}) \in \mathbb{R}^{31}$

Queremos encontrar uma $g^*(\underline{x}) \in \mathcal{C}$ (11)

(isto é, escolher coeficientes $\theta_0, \theta_1, \dots, \theta_{30}$)

tal que

$$E \left[\left(Y - \underbrace{(\theta_0 + \theta_1 X_1 + \dots + \theta_{30} X_{30})}_{g(\underline{x})} \right)^2 \right]$$

seja mínimo

Queremos encontrar uma $g^*(\underline{x}) \in \mathcal{C}$ ⑪

(isto é, escolher coeficientes $\theta_0, \theta_1, \dots, \theta_{30}$)

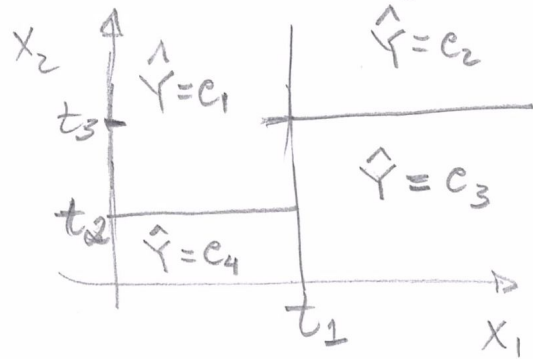
tal que

$$E \left[\left(Y - \underbrace{(\theta_0 + \theta_1 X_1 + \dots + \theta_{30} X_{30})}_{g(\underline{x})} \right)^2 \right]$$

seja mínimo

Justificativa para esta escolha linear para o conjunto \mathcal{C} : usar a melhor aprox polinomial para $E(Y | X) = \mu(x)$. Aqui, usamos aprox do polinomio de primeira ordem (vimos na aula passada).

Example: Regression tree with two features ⁽¹²⁾:



$$\hat{Y} = g(\underline{X}) = \text{prediction}$$

→ "função azulada"

$$\hat{Y} = g(X_1, X_2) = \begin{cases} c_1, & \text{if } X_1 \leq t_1 \quad \underline{\text{AND}} \quad X_2 > t_2 \\ c_2, & \text{if } X_1 > t_1 \quad \underline{\text{AND}} \quad X_2 > t_3 \\ c_3, & \text{if } X_1 > t_1 \quad \underline{\text{AND}} \quad X_2 \leq t_3 \\ c_4, & \text{if } X_1 \leq t_1 \quad \underline{\text{AND}} \quad X_2 \leq t_2 \end{cases}$$

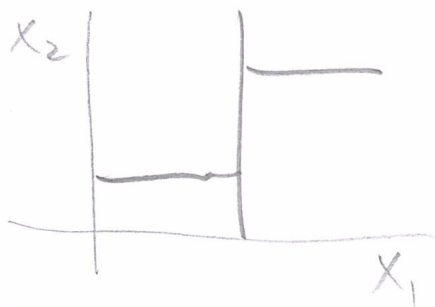
⊕ Achar a melhor $g(x_1, x_2)$ nesta

(13)

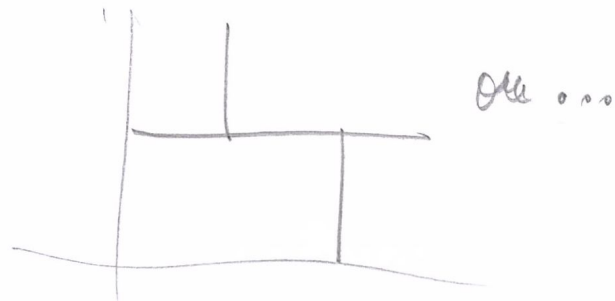
classe de "funções azulejas" \mathcal{F}

é encontrar os pontos de corte

t_1, t_2, t_3 e os valores c_1, c_2, c_3, c_4



ou



RESUMO Objetivo: achar a função de \underline{x} (14)

ótima (para minimizar o erro esperado)

$$\mu(\underline{x}) = E(Y | \underline{x})$$

Soluções Definir uma classe de funções \mathcal{F}

que seja:

- trade-off
- ⊕ flexível, maleável, que possa aproximar qualquer função
 - ⊕ "pequena" para ser trabalhada numericamente

Soluções: SVM, NN, Regressão linear...

Na classe \mathcal{C} escolhida pelo algoritmo,⁽¹⁵⁾
achar a funcp $g^*(x) \in \mathcal{C}$ que
melhor aproxime $\mu(x) = E(Y|X)$

{ Mas como encontrar esta $g^*(x)$ se não
conhecemos $\mu(x)$?

{ Na prática, buscamos minimizar a
soma dos erros de predição no conjunto
de treinamento

Para qualquer vetor aleatório

$$(Y, \underline{X}) = (Y, X_1, \dots, X_{30})$$

O valor esperado de uma função

$h(Y, \underline{X})$ é aprox. a média

aritmética de h calculada numa amostra de tamanho grande:

$$E[h(Y, \underline{X})] \approx \frac{1}{N} \sum_{i=1}^N h(Y_i, \underline{X}_i)$$

$$E \left(\underbrace{(Y - g(\underline{x}))^2}_{h(Y, \underline{x})} \right) \approx \underbrace{\frac{1}{N} \sum_{i=1}^N (Y_i - g(\underline{x}_i))^2}_{\text{Ache } g^* \in \mathcal{G} \text{ que minimiza esta soma}} \quad (17)$$

⊕ Como minimizar um objetivo num espaço de funções? Derivar e igualar a zero? Derivar em relação a quê?

Se \mathcal{C} é uma classe parametrizada por ⁽¹⁸⁾
um vetor de coeficientes $(\theta_0, \theta_1, \dots, \theta_{30})$,
derive um relacp a estes coeficientes.

EX: Regressão Linear

$$g(x) = \beta_0 + \beta_1 x$$

Achar $\boxed{\beta_0 \text{ e } \beta_1}$ que minimizam

$$\frac{1}{N} \sum_{i=1}^M (y_i - (\beta_0 + \beta_1 x_i))^2 = \frac{1}{N} \left(\begin{aligned} & (70.3 - (\beta_0 + \beta_1 150))^2 + \\ & + (90.7 - (\beta_0 + \beta_1 210))^2 + \dots \end{aligned} \right)$$

Ignorando a constante $1/N$, derive em relação a β_0 e β_1 e iguale a zero:

- Derivando e igualando a zero:

$$0 = \sum_i \frac{\partial}{\partial \beta_0} (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$0 = \sum_i \frac{\partial}{\partial \beta_1} (y_i - (\beta_0 + \beta_1 x_i))^2$$

- Ou seja:

$$0 = \sum_i 2(y_i - (\beta_0 + \beta_1 x_i)) (-1)$$

$$0 = \sum_i 2(y_i - (\beta_0 + \beta_1 x_i)) (-x_i)$$

- Temos

$$0 = - \sum_i y_i + \beta_0 n + \beta_1 \sum_i x_i$$

$$0 = - \sum_i (y_i x_i) + \beta_0 \sum_i x_i + \beta_1 \sum_i x_i^2$$

- Rearranjando:

$$\beta_0 n + \beta_1 \sum_i x_i = \sum_i y_i$$

$$\beta_0 \sum_i x_i + \beta_1 \sum_i x_i^2 = \sum_i (y_i x_i)$$

- Este é um sistema linear de duas equações com duas incógnitas, β_0 e β_1 .

- Sistema na forma matricial

$$\begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \sum_i y_i \\ \sum_i (x_i y_i) \end{bmatrix}$$

- Com solução:

$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix}^{-1} \begin{bmatrix} \sum_i y_i \\ \sum_i (x_i y_i) \end{bmatrix}$$

- Vamos usar uma notação para simplificar as expressões.

- Vamos denotar a média dos x e y 's por

- $\bar{x} = \frac{1}{n} \sum_i x_i$, média aritmética dos x_i 's

- $\bar{y} = \frac{1}{n} \sum_i y_i$

- $\overline{x^2} = \frac{1}{n} \sum_i x_i^2$, média aritmética dos x_i^2 's

- $\overline{xy} = \frac{1}{n} \sum_i (x_i y_i)$, média aritmética dos $x_i y_i$'s

- Sistema na forma matricial e com a notação introduzida

$$\begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \overline{x^2} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix}$$

- Com solução:

$$\begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \begin{bmatrix} 1 & \bar{x} \\ \bar{x} & \overline{x^2} \end{bmatrix}^{-1} \begin{bmatrix} \bar{y} \\ \overline{xy} \end{bmatrix}$$

- Como a inversa de uma matriz 2×2 é conhecida, podemos resolver de forma explícita a solução de mínimos quadrados.
- Após alguma manipulação algébrica, temos a solução como uma fórmula envolvendo os pontos:

$$\hat{\beta}_1 = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sum_i (x_i - \bar{x})^2}$$
$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

Regressão: caso geral

- Vetor \mathbf{Y} de dimensão n :
 - preço de n aptos
- $(k+1)$ features, vetores-coluna de dimensão n
 - Feature 0: vetor-coluna de 1's
 - Feature 1: vetor-coluna com área dos aptos n aptos
 -
 - Feature k : vetor-coluna com indicador binário “tem salão de festa?”
- Colete as features numa matriz \mathbf{X} de dimensão $n \times (k + 1)$
- Objetivo: minimizar

$$\sum_{i=1}^n (y_i - \beta_0 * 1 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})^2 = \sum_{i=1}^n \|\mathbf{Y} - \mathbf{X}\beta\|^2$$

Regressão: caso geral

- Derive
- $\sum_{i=1}^n (y_i - \beta_0 * 1 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})^2 = \sum_{i=1}^n \|\mathbf{Y} - \mathbf{X}\beta\|^2$
- em relação aos coeficientes β_j de cada feature
- Iguale cada derivada parcial a zero
- “Isole” os beta’s
- Teremos um sistema linear chamado equações normais:

$$(\mathbf{X}^t \mathbf{X}) \beta = \mathbf{X}^t \mathbf{Y}$$

Equações normais

$$\underbrace{(\mathbf{X}^t \mathbf{X})}_{(k+1) \times n \quad n \times (k+1)} \quad \underbrace{\beta}_{(k+1) \times 1} = \underbrace{\mathbf{X}^t \mathbf{Y}}_{(k+1) \times n \quad n \times 1}$$

$$\underbrace{(\mathbf{X}^t \mathbf{X})}_{(k+1) \times (k+1)} \quad \underbrace{\beta}_{(k+1) \times 1} = \underbrace{\mathbf{X}^t \mathbf{Y}}_{(k+1) \times 1}$$

Temos um sistema linear

$\mathbf{X}'\mathbf{X}$ é matriz quadrada $(k+1) \times (k+1)$

Se tiver inversa, podemos obter a solução

$$\hat{\beta} = (\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t \mathbf{Y}$$

Uma abordagem mais geral

Solução Least Squares foi vista como um problema numérico:

minimizar uma função objetivo

$$L(\beta_0, \beta_1, \dots, \beta_k) = \sum_i (\text{preço}_i - \beta_0 * 1 + \beta_1 \text{área}_i + \dots + \beta_k \text{banheiro?}_i)^2$$

Ao invés de olhar o problema como um problema numérico, vamos dar uma roupagem mais teórica.

Vai permitir uma solução mais elegante e ... generalizável para espaços de funções

⊕ Mas, e quando \mathcal{C} não é tão simples? (19)

⊕ O que é uma classe \mathcal{C} boa?

⊕ Podemos garantir que sempre vamos encontrar uma função $g^*(\underline{x})$ que seja uma boa aproximação para $\mu(\underline{x}) = E(Y|\underline{x})$?

⊕ Precisamos de três coisas:

- trabalhar no espaço de funções
 - medir distância entre funções
 - saber quando uma sequência de funções converge (aproxima) para outra
-

- ⊕ Isto nos leva a conceitos de análise funcional
- ⊕ Recordar espaço vetorial

Conjunto de objetos em que:

- podemos somar objetos e obter um 3º objeto
- podemos "espichar", "encolher" ou "reverter" objetos ao multiplicar por escalares

To have a vector space, the eight following **axioms** must be satisfied for every \mathbf{u} , \mathbf{v} and \mathbf{w} in \mathcal{V} , and a and b in F .^[3]

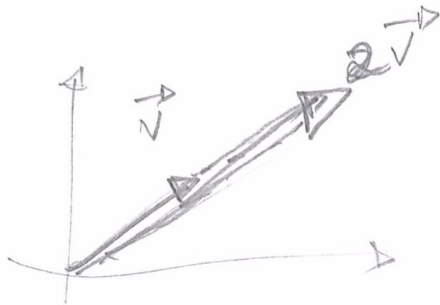
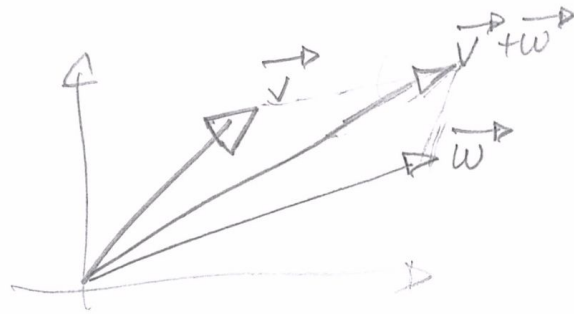
Axiom	Statement
Associativity of vector addition	$\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$
Commutativity of vector addition	$\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
Identity element of vector addition	There exists an element $\mathbf{0} \in \mathcal{V}$, called the <i>zero vector</i> , such that $\mathbf{v} + \mathbf{0} = \mathbf{v}$ for all $\mathbf{v} \in \mathcal{V}$.
Inverse elements of vector addition	For every $\mathbf{v} \in \mathcal{V}$, there exists an element $-\mathbf{v} \in \mathcal{V}$, called the <i>additive inverse</i> of \mathbf{v} , such that $\mathbf{v} + (-\mathbf{v}) = \mathbf{0}$.
Compatibility of scalar multiplication with field multiplication	$a(b\mathbf{v}) = (ab)\mathbf{v}$ ^[nb 3]
Identity element of scalar multiplication	$1\mathbf{v} = \mathbf{v}$, where 1 denotes the <i>multiplicative identity</i> in F .
Distributivity of scalar multiplication with respect to vector addition	$a(\mathbf{u} + \mathbf{v}) = a\mathbf{u} + a\mathbf{v}$
Distributivity of scalar multiplication with respect to field addition	$(a + b)\mathbf{v} = a\mathbf{v} + b\mathbf{v}$

Exemplos de espaços vetoriais:

(4)

$$\vec{v} \in \mathbb{R}^2 = \{ \vec{v} = (x, y), x \in \mathbb{R}, y \in \mathbb{R} \}$$

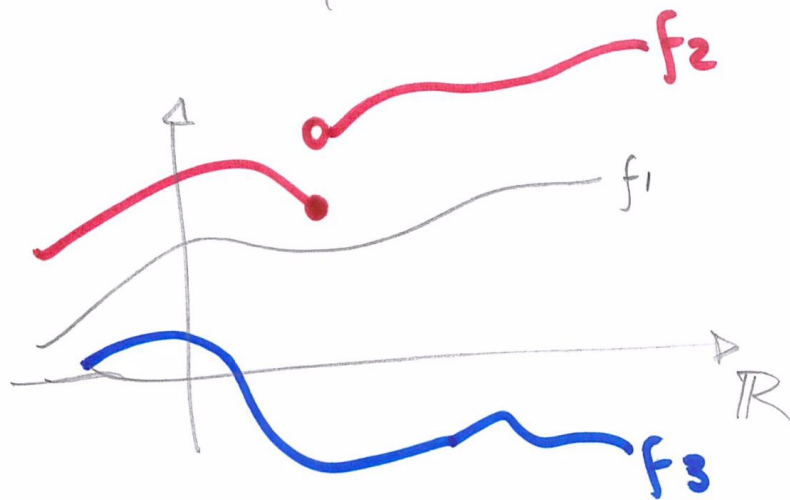
com escalares reais



$$\oplus \mathbb{R}^m = \{ (x_1, x_2, \dots, x_m), x_i \in \mathbb{R} \}$$

(22)

\oplus Espaço de funções \mathcal{F}
 $\mathcal{F} = \{ f: \mathbb{R} \rightarrow \mathbb{R} \}$

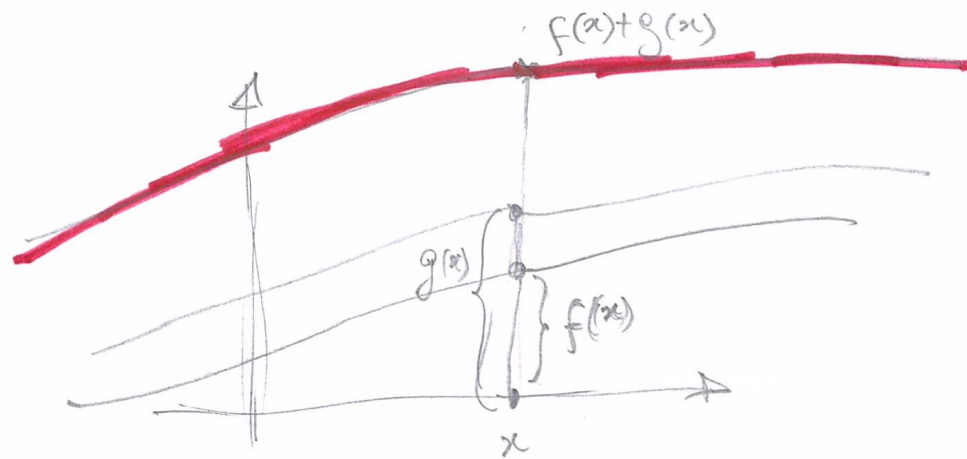
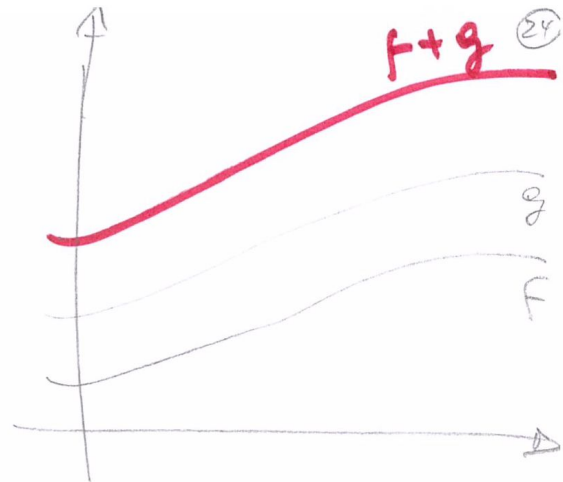
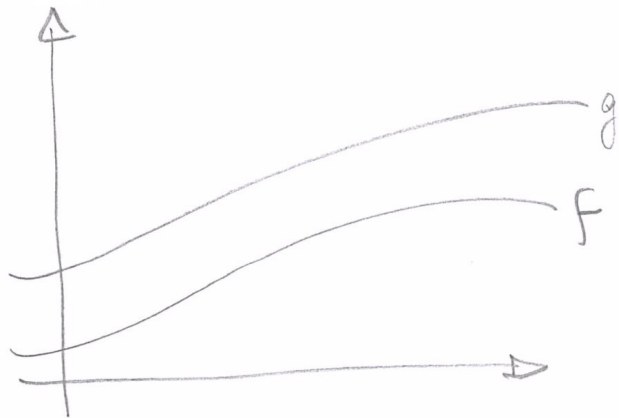


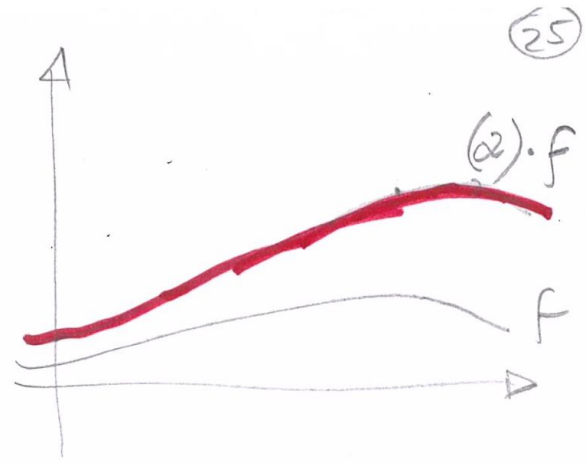
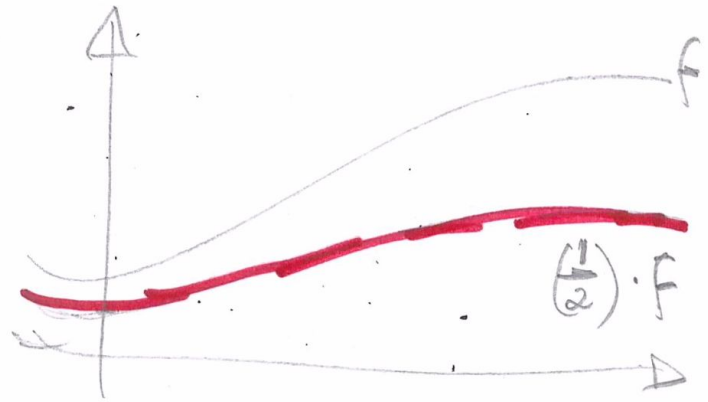
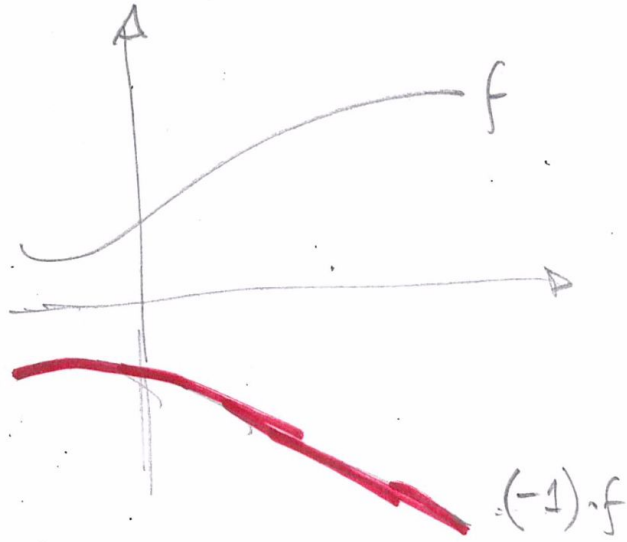
somar duas funções \rightarrow ainda é uma função (23)

Definir a soma $g + f = h$:

é a função tal que

$$h(x) = g(x) + f(x) \text{ para todo } x$$



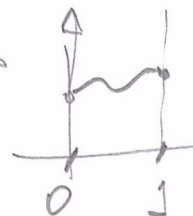


$\mathcal{F} = \{ f \text{ funções em qual domínio?} \}$ (26)

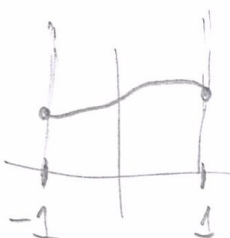
$\mathcal{F}(\mathbb{R}) = \{ f: \mathbb{R} \rightarrow \mathbb{R} \}$



$\mathcal{F}([0, 1]) = \{ f: [0, 1] \rightarrow \mathbb{R} \}$



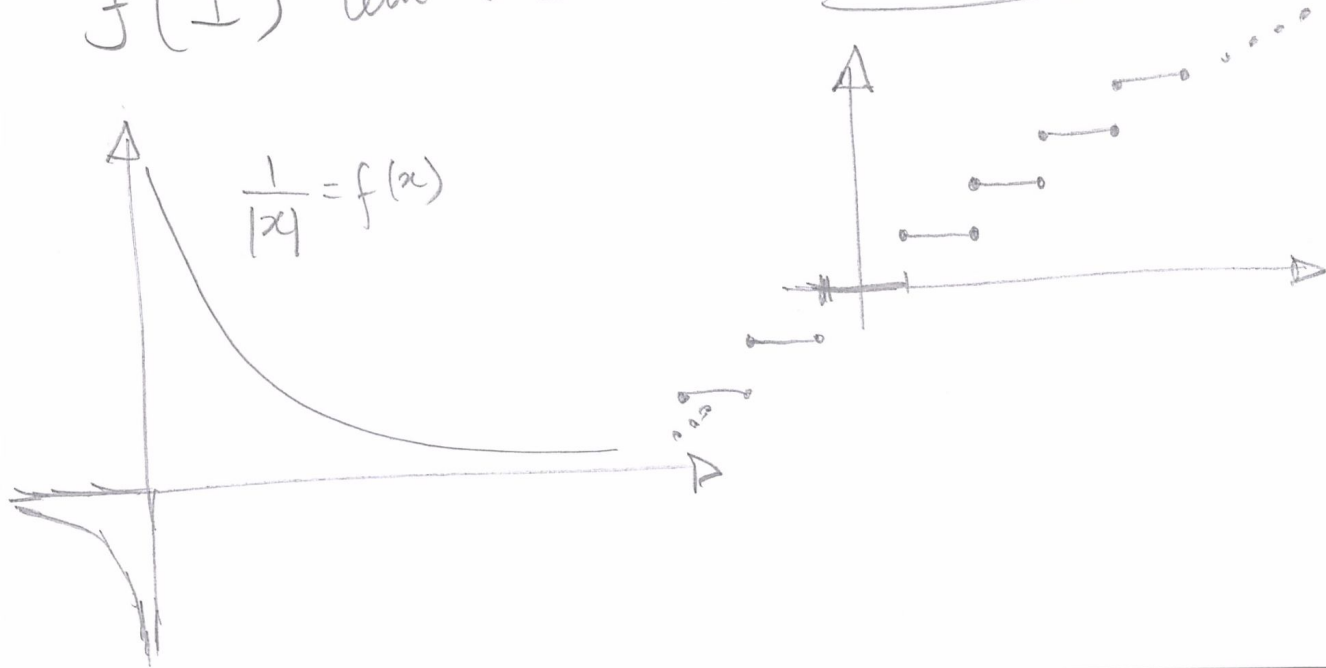
$\mathcal{F}([-1, 1]) = \{ f: [-1, 1] \rightarrow \mathbb{R} \}$



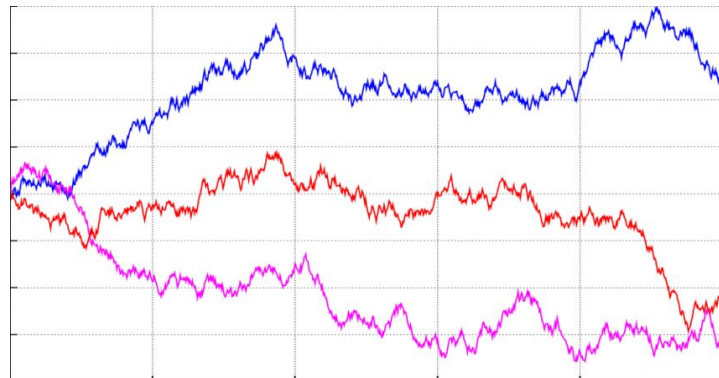
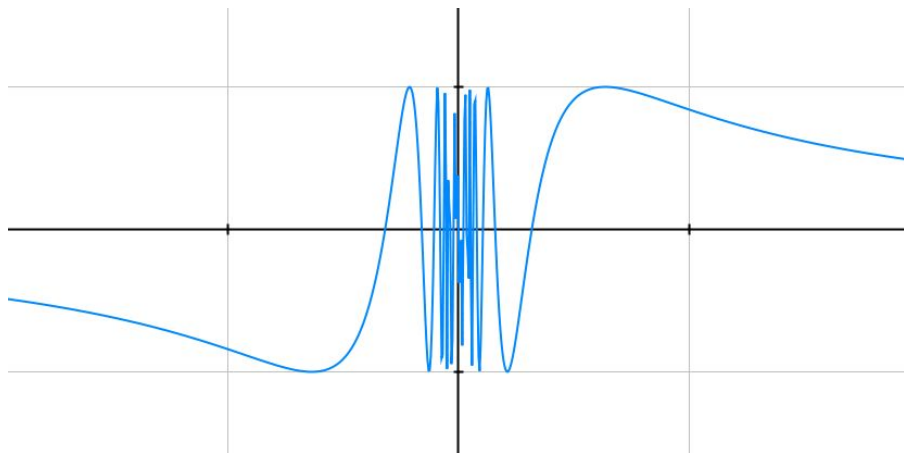
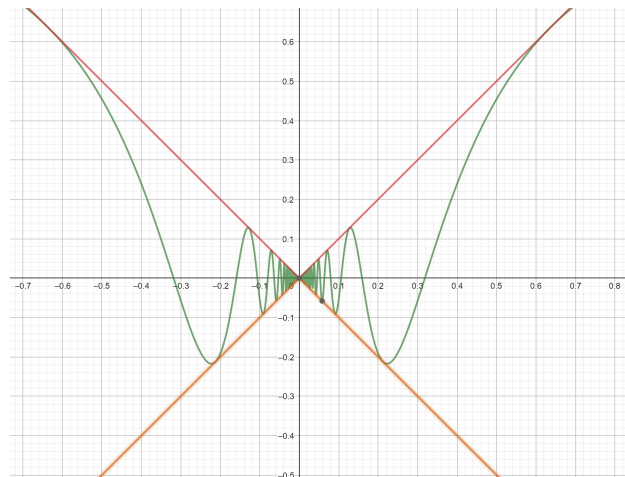
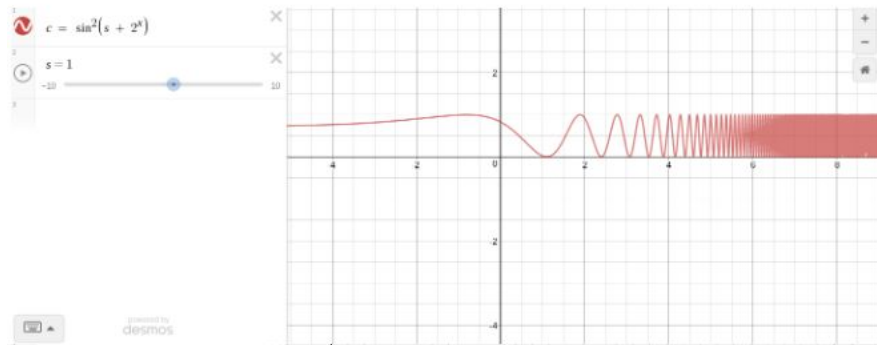
$$\mathcal{F}(I) = \{ f: \text{Intervalo } I \rightarrow \mathbb{R} \}$$

(27)

$\mathcal{F}(I)$ tem elementos muito estranhos



$$c = \sin^2(s + 2^x) \quad s = 1$$



Ignorar espaços de funções (por enquanto)
Foco em \mathbb{R}^n e regressão linear

Exemplo de preço de apto

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{1499} \\ y_{1500} \end{pmatrix} \approx b_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix} + b_1 \begin{pmatrix} \text{área}_1 \\ \text{área}_2 \\ \vdots \\ \text{área}_{1499} \\ \text{área}_{1500} \end{pmatrix} + b_2 \begin{pmatrix} \text{idade}_1 \\ \text{idade}_2 \\ \vdots \\ \text{idade}_{1499} \\ \text{idade}_{1500} \end{pmatrix} + \dots + b_{30} \begin{pmatrix} \text{salão}_1 \\ \text{salão}_2 \\ \vdots \\ \text{salão}_{1499} \\ \text{salão}_{1500} \end{pmatrix}$$

- Y é um vetor de dimensão 1500 escrito como combinação linear de 31 vetores, cada um deles de dimensão 1500.
- Problema: encontrar os coeficientes b_0, b_1, \dots, b_{30} que tornem a aproximação acima a melhor possível.

A matriz de desenho X

- Seja X a matriz 1500×31 abaixo (note que ela tem uma coluna composta apenas de 1's):

$$X = \begin{pmatrix} 1 & \text{renda}_1 & \text{área}_1 & \cdots & \text{salão}_1 \\ 1 & \text{renda}_2 & \text{área}_2 & \cdots & \text{salão}_2 \\ \vdots & \vdots & \vdots & \vdots & \\ 1 & \text{renda}_{1499} & \text{área}_{1499} & \cdots & \text{salão}_{1499} \\ 1 & \text{renda}_{1500} & \text{área}_{1500} & \cdots & \text{salão}_{1500} \end{pmatrix}$$

Combinações lineares e a matriz X

- A combinação linear que buscamos

$$b_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix} + b_1 \begin{pmatrix} \text{área}_1 \\ \text{área}_2 \\ \vdots \\ \text{área}_{1499} \\ \text{área}_{1500} \end{pmatrix} + b_2 \begin{pmatrix} \text{idade}_1 \\ \text{idade}_2 \\ \vdots \\ \text{idade}_{1499} \\ \text{idade}_{1500} \end{pmatrix} + \dots + b_{30} \begin{pmatrix} \text{salão}_1 \\ \text{salão}_2 \\ \vdots \\ \text{salão}_{1499} \\ \text{salão}_{1500} \end{pmatrix}$$

- pode ser escrita como

$$X b = \begin{bmatrix} 1 & \text{renda}_1 & \text{área}_1 & \dots & \text{salão}_1 \\ 1 & \text{renda}_2 & \text{área}_2 & \dots & \text{salão}_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \text{renda}_{1499} & \text{área}_{1499} & \dots & \text{salão}_{1499} \\ 1 & \text{renda}_{1500} & \text{área}_{1500} & \dots & \text{salão}_{1500} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{30} \end{bmatrix}$$

Vetores próximos

Nosso problema é encontrar os coeficientes b_0, b_1, \dots, b_{30} tais que

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_{1499} \\ y_{1500} \end{pmatrix} \approx b_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix} + b_1 \begin{pmatrix} \text{área}_1 \\ \text{área}_2 \\ \vdots \\ \text{área}_{1499} \\ \text{área}_{1500} \end{pmatrix} + b_2 \begin{pmatrix} \text{idade}_1 \\ \text{idade}_2 \\ \vdots \\ \text{idade}_{1499} \\ \text{idade}_{1500} \end{pmatrix} + \dots + b_{30} \begin{pmatrix} \text{salão}_1 \\ \text{salão}_2 \\ \vdots \\ \text{salão}_{1499} \\ \text{salão}_{1500} \end{pmatrix}$$

Ou seja, encontrar b_0, b_1, \dots, b_{30} tais que $Y \approx Xb$ onde

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_{1498} \\ y_{1499} \\ y_{1500} \end{pmatrix} \approx \begin{pmatrix} 1 & \text{renda}_1 & \text{área}_1 & \cdots & \text{salão}_1 \\ 1 & \text{renda}_2 & \text{área}_2 & \cdots & \text{salão}_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \text{renda}_{1499} & \text{área}_{1499} & \cdots & \text{salão}_{1499} \\ 1 & \text{renda}_{1500} & \text{área}_{1500} & \cdots & \text{salão}_{1500} \end{pmatrix} \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_{30} \end{pmatrix} = Xb$$

onde $b = (b_0, \dots, b_{30})^t$.

Solução: minimizar norma

- X é uma matriz 1500×31 .
- Y e Xb são vetores 1500-dim.
- Além disso, Xb é uma combinação linear das colunas da matriz X .
- Queremos encontrar b tal que o vetor Xb seja o mais próximo possível do vetor Y .
- Queremos $Y - Xb$ aproximadamente igual AO VETOR ZERO.
- Queremos $\|Y - Xb\| \approx 0$ (o comprimento-norma é um número, não um vetor)

Solução melhor: minimizar norma ao quadrado

- Queremos $\|Y - Xb\| \approx 0$
- Queremos \hat{b} que minimize $\|Y - Xb\|$
- Mas norma euclidiana envolve a raiz quadrada da soma dos quadrados ...
- Mas se \hat{b} minimiza $\|Y - Xb\|$ então \hat{b} minimiza $\|Y - Xb\|^2$
- Esta segunda função é mais fácil de derivar.

Solução melhor: minimizar norma ao quadrado

- Então procuramos vetor b tal que $\|Y - Xb\|^2 \approx 0$.
- Queremos \hat{b} que minimize $\|Y - Xb\|^2$
- Matematicamente: queremos $\hat{b} = \arg \min_b \|Y - Xb\|^2$.
- Como encontrar este \hat{b} ?

Vetores e combinações lineares

- X é matriz 1500×31 . b é vetor 31×1
- Para qualquer vetor $b \in \mathbb{R}^{31}$, temos Xb em \mathbb{R}^{1500} .
- Varie b varrendo todos os vetores b possíveis. O que obtemos?
- Isto é, o que é o conjunto

$$\mathfrak{M}(X) = \{v \in \mathbb{R}^{1500} \text{ tais que } v = Xb \text{ para algum } b\} \quad ?$$

O que é $\mathfrak{M}(X)$?

- Colunas da matriz X estão fixadas, são vetores 1500×1 de números constantes, conhecidos.

$$\mathfrak{M}(X) = \left\{ b_0 \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{pmatrix} + b_1 \begin{pmatrix} \text{área}_1 \\ \text{área}_2 \\ \vdots \\ \text{área}_{1499} \\ \text{área}_{1500} \end{pmatrix} + b_2 \begin{pmatrix} \text{idade}_1 \\ \text{idade}_2 \\ \vdots \\ \text{idade}_{1499} \\ \text{idade}_{1500} \end{pmatrix} + \dots + b_{30} \begin{pmatrix} \text{salão}_1 \\ \text{salão}_2 \\ \vdots \\ \text{salão}_{1499} \\ \text{salão}_{1500} \end{pmatrix} \right\}$$

- $\mathfrak{M}(X)$ é um subconjunto de vetores do espaço vetorial \mathbb{R}^{1500} .
- Vetor zero pertence a $\mathfrak{M}(X)$.
- Somando duas combinações lineares de $\mathfrak{M}(X)$ ainda permanecemos em $\mathfrak{M}(X)$.
- Multiplicando um elemento de $\mathfrak{M}(X)$ por uma constante ainda permanecemos em $\mathfrak{M}(X)$.

Sub-espços vetoriais

- Informal: Sub-espço vetorial W de um espço vetorial V é um subconjunto de V tal que:
 - a soma de dois vetores de W permanece em W
 - multiplicar um vetor de W por um escalar permanece em W
 - O vetor 0 (nulo) pertence a W
- Só isto: W é um espço vetorial e não saímos de dentro dele ao manusear seus vetores com adição ou multiplicação por escalar.

Espaço $\mathfrak{M}(X)$ das combinações lineares

- $\mathfrak{M}(X)$ é um sub-espço vetorial de \mathbb{R}^{1500} .
- $\mathfrak{M}(X)$ é o sub-espço vetorial formado pelas combinações lineares dos 31 vetores-colunas de X .
- Se as colunas de X são linearmente independentes, então $\mathfrak{M}(X)$ é um sub-espço vetorial de dimensão igual ao número de colunas de X (que é 31, no nosso exemplo).
- Nosso problema então é: encontrar os coeficientes b da combinação linear $Xb \in \mathfrak{M}(X)$ tal que Xb seja o mais próximo possível do vetor Y .

Geometria dos Mínimos quadrados

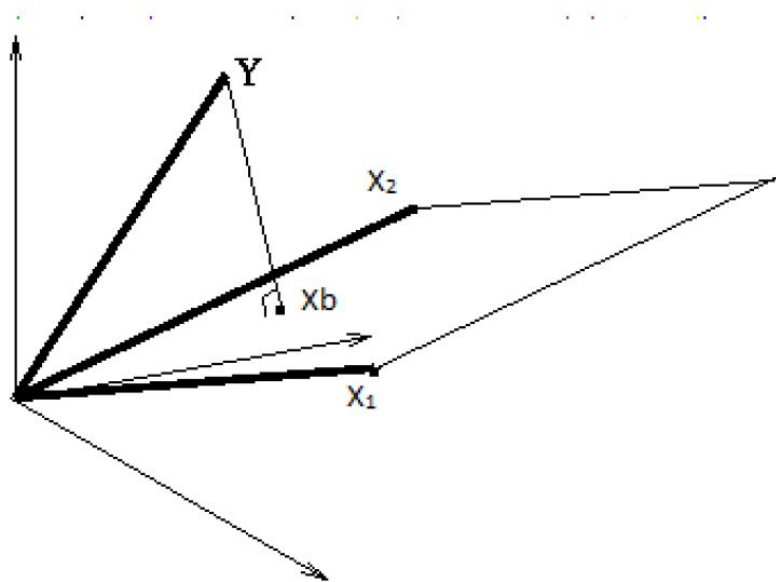


Figura: Representação do vetor $Y \in \mathbb{R}^{1500}$. O plano inclinado representa o sub-espço vetorial $\mathfrak{M}(X)$ gerado por uma matriz X com apenas duas colunas, os vetores X_1 e X_2 , ambos do \mathbb{R}^{1500} . O sub-espço vetorial $\mathfrak{M}(X)$ é de dimensão 2. Identifique visualmente o ponto-vetor em $\mathfrak{M}(X)$ que minimiza $\|Y - Xb\|^2$.

Teorema da Projeção Ortogonal

- Seja \mathbb{R}^n um espaço vetorial real de dimensão n .
- Seja \mathcal{W} um sub-espaço vetorial de \mathcal{V} com dimensão m .
- Seja $Y \in \mathbb{R}^n$ um vetor qualquer.
- **Teorema:** Existe um único vetor $\hat{w} \in \mathcal{W}$ que minimiza $\|Y - w\|$ com $w \in \mathcal{W}$.
- Além disso, este $\hat{w} \in \mathcal{W}$ é o único vetor tal que $Y - \hat{w}$ é ortogonal a \hat{w} . Isto é, \hat{w} é o único vetor tal que $(Y - \hat{w}) \perp \hat{w}$.

PROVA: a seguir

Teorema da Projeção Ortogonal

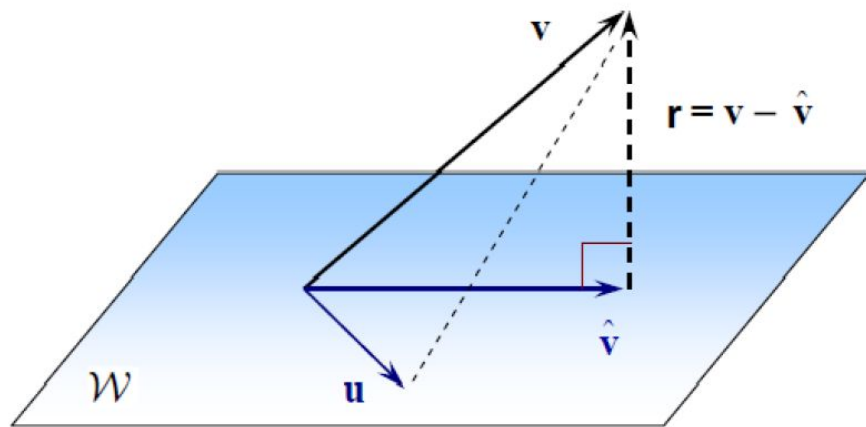


Figura: \mathbf{v} é um vetor do \mathbb{R}^3 . O plano \mathcal{W} é um sub-espço vetorial de dimensão 2. dado um vetor \mathbf{u} do sub-espço, $\|\mathbf{v} - \mathbf{u}\|$ (linha tracejada fina) é o comprimento do vetor $\mathbf{v} - \mathbf{u}$.

De todos os vetores \mathbf{u} do sub-espço \mathcal{W} , aquele que minimiza o comprimento $\|\mathbf{v} - \mathbf{u}\|$ é a projeção ortogonal $\hat{\mathbf{v}}$. O vetor $\hat{\mathbf{v}}$ é a aproximação de mínimos quadrados em \mathcal{W} para \mathbf{v} . O vetor $r = \mathbf{v} - \hat{\mathbf{v}}$ é o vetor de resíduos.

Geometria dos Mínimos quadrados

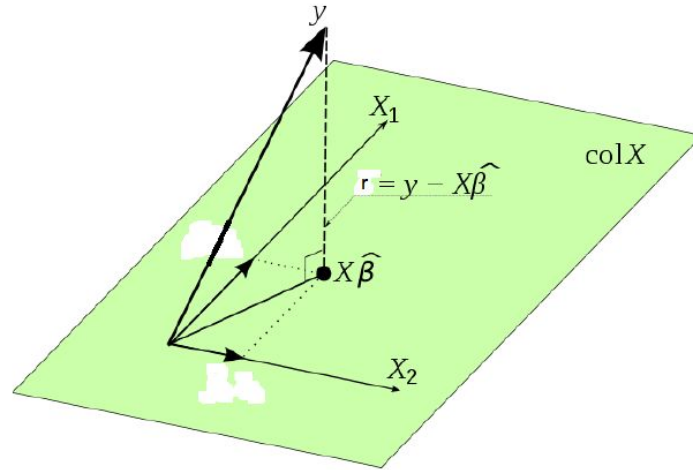


Figura: Projeção ortogonal de $Y \in \mathbb{R}^{1500}$ no sub-espaço vetorial $\mathfrak{M}(X)$ minimiza $\|Y - X\beta\|^2$. Esta projeção é o vetor $X\hat{\beta}$. Vetor de resíduos é $r = Y - X\hat{\beta}$ e é \perp a $X\hat{\beta}$. Imagem retirada de

<https://commons.wikimedia.org/w/index.php?curid=7309159>

- Demonstração do Teorema da Projeção.
- Não veremos a demonstração geral para espaços vetoriais arbitrários.
- Vamos fazer apenas o caso especial da regressão linear.

Produto interno com vetores-coluna

- Produto interno de dois vetores, \mathbf{v} e \mathbf{w} , no mesmo espaço vetorial de dimensão n :

$$\langle \mathbf{v}, \mathbf{w} \rangle = \sum_{i=1}^n v_i w_i$$

- Como \mathbf{v} e \mathbf{w} são vetores-coluna, temos

$$\langle \mathbf{v}, \mathbf{w} \rangle = \sum_i v_i w_i = \begin{bmatrix} v_1 & v_2 & \dots & v_n \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} = \mathbf{v}^t \mathbf{w}$$

- Temos $\mathbf{v} \perp \mathbf{w}$ se, e só se, $\langle \mathbf{v}, \mathbf{w} \rangle = \mathbf{v}^t \mathbf{w} = 0$.
- Imitando o Scilab, vamos denotar \mathbf{v}^t por \mathbf{v}' .

- Primeiro, queremos achar $\hat{Y} = X\hat{\beta}$ tal que $(Y - \hat{Y}) \perp \hat{Y}$
- Depois, queremos mostrar que este $\hat{Y} = X\hat{\beta}$ é o vetor que minimiza a distância $\|Y - X\beta\|^2$

$$\hat{Y} = X\beta \perp (Y - \hat{Y}), \forall Y$$

$$\langle X\beta, Y - \hat{Y} \rangle = 0$$

$$\begin{aligned} 0 &= (X\beta)^t (Y - X\beta) \\ &= \beta^t X^t (Y - X\beta) \\ &= \beta^t (X^t Y - X^t X\beta) \end{aligned}$$

- Ou $\beta^t = 0$ (o que implica que $\beta = 0$ e que $\hat{Y} = 0$), solução sem sentido
- Ou $X^t Y - X^t X\beta = 0 \implies$ com solução $\hat{\beta} = (X^t X)^{-1} X^t Y$

- Seja $\hat{Y} = X\hat{\beta} = X(X^tX)^{-1}X^tY$
- Vamos calcular agora $\|Y - X\beta\|^2$ para um β arbitrário:
- Some e subtraia: $Y - X\beta = Y - X\hat{\beta} + X\hat{\beta} - X\beta$
- Calculando:

$$\begin{aligned}
 \|Y - X\beta\|^2 &= (Y - X\beta)^t (Y - X\beta) = \\
 &= \left((Y - X\hat{\beta}) + (X\hat{\beta} - X\beta) \right)^t \left((Y - X\hat{\beta}) + (X\hat{\beta} - X\beta) \right) \\
 &= (Y - X\hat{\beta})^t (Y - X\hat{\beta}) + (Y - X\hat{\beta})^t (X\hat{\beta} - X\beta) + \\
 &+ (X\hat{\beta} - X\beta)^t (Y - X\hat{\beta}) + (X\hat{\beta} - X\beta)^t (X\hat{\beta} - X\beta) \\
 &= \left\| Y - X\hat{\beta} \right\|^2 + \underbrace{2 (Y - X\hat{\beta})^t (X\hat{\beta} - X\beta)}_A + \left\| X\hat{\beta} - X\beta \right\|^2
 \end{aligned}$$

- Vamos mostrar agora que $A = 0$

$$\begin{aligned}
 (Y - X\hat{\beta})^t (X\hat{\beta} - X\beta) &= (Y - X(X^tX)^{-1}X^tY)^t (X\hat{\beta} - X\beta) \\
 &= \left((I - X(X^tX)^{-1}X^t) Y \right)^t (X\hat{\beta} - X\beta) \\
 &= \left(Y^t (I - X(X^tX)^{-1}X^t)^t \right) (X(\hat{\beta} - \beta)) = (*)
 \end{aligned}$$

- Temos $(I - X(X^tX)^{-1}X^t)^t = I^t - (X^t)^t \left((X^tX)^{-1} \right)^t X^t$

$$= I - X \left((X^tX)^t \right)^{-1} X^t$$

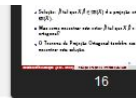
$$= I - X (X^tX)^{-1} X^t$$

- $$\begin{aligned}
\bullet \text{ Assim } (*) &= Y^t \left(I - X (X^t X)^{-1} X^t \right) X \left(\hat{\beta} - \beta \right) \\
&= Y^t \left[\left(I - X (X^t X)^{-1} X^t \right) X \right] \left(\hat{\beta} - \beta \right) \\
&= Y^t \left[X - X \underbrace{(X^t X)^{-1} X^t X}_I \right] \left(\hat{\beta} - \beta \right) \\
&= Y^t [X - X] \left(\hat{\beta} - \beta \right) \\
&= Y^t [0] \left(\hat{\beta} - \beta \right) = 0
\end{aligned}$$

- Portanto:

$$\|Y - X\beta\|^2 = \|Y - X\hat{\beta}\|^2 + 0 + \|X\hat{\beta} - X\beta\|^2 \geq \|Y - X\hat{\beta}\|^2$$

- pois, sendo uma distância, $\|X\hat{\beta} - X\beta\|^2$ é sempre maior ou igual a zero.
- Isto é, $\hat{\beta} = (X^t X)^{-1} X^t Y$ é o vetor de coeficientes β que minimiza $\|Y - X\beta\|^2$ para todo vetor Y



- Nosso problema: encontrar $\hat{\beta}$ tal que o vetor $X\hat{\beta}$ do subespaço $\mathfrak{M}(X)$ seja o mais próximo possível do vetor Y .
- O Teorema da Projeção Ortogonal garante que existe uma solução. Além disso,...
- Solução: $\hat{\beta}$ tal que $X\hat{\beta} \in \mathfrak{M}(X)$ é a projeção ortogonal de Y em $\mathfrak{M}(X)$.
- Mas como encontrar este vetor $\hat{\beta}$ tal que $X\hat{\beta}$ e seja esta projeção ortogonal?
- O Teorema da Projeção Ortogonal também nos dá a dica de como encontrar esta solução.

Encontrando a solução de mínimos quadrados

- Solução é $\hat{\beta}$, um vetor tal que $X\hat{\beta} \in \mathfrak{M}(X)$ é a projeção ortogonal de Y em $\mathfrak{M}(X)$.
- O Teorema da Projeção Ortogonal diz que a projeção $X\hat{\beta}$ é única e é o vetor tal que $X\hat{\beta} \perp (Y - X\hat{\beta})$.
- Em resumo, devemos ter o produto interno zerado:
 $\langle X\hat{\beta}, (Y - X\hat{\beta}) \rangle = 0$.
- Isto implica que

$$\hat{\beta} = (X^t X)^{-1} X^t Y$$

Example


Concrete Compressive Strength

- A Resistência à Compressão do Concreto é uma medida da capacidade do concreto de suportar cargas que tendem a comprimir ou reduzir seu volume.
- É uma das propriedades mais críticas do concreto, indicando sua capacidade de suportar cargas estruturais sem romper.
- É definida como a tensão compressiva máxima que o concreto pode resistir antes da falha.
- É determinada pela aplicação de uma carga compressiva a uma amostra de concreto (geralmente um cubo ou um cilindro) sob condições controladas até que ela se quebre.
- <https://www.youtube.com/shorts/yGpwrL0zwHI>

Kaggle Dataset

- Aim: To predict the compressive strength of concrete based on material composition.

Target Variable (Response Variable)

Feature Name	Description	Units	Typical Range
Compressive Strength	The maximum compressive stress the concrete can withstand. 	MPa (MegaPascals)	2.33 - 82.6

- Number of Samples: 1,030 observations
- Number of Features: 8 predictors

Features (Predictor Variables)

Feature Name	Description	Units	Typical Range
Cement	The amount of cement used in the mix.	kg/m ³	102 - 540
Blast Furnace Slag	By-product of steel production, often used as a cement substitute.	kg/m ³	0 - 359.4
Fly Ash	A by-product of coal combustion, used as a partial cement replacement.	kg/m ³	0 - 200.1
Water	The amount of water used in the mix.	kg/m ³	121.8 - 247
Superplasticizer	Chemical additive to enhance workability and strength.	kg/m ³	0 - 32.2
Coarse Aggregate	Gravel or crushed stone used as a filler material.	kg/m ³	801 - 1145
Fine Aggregate	Sand used as a filler material.	kg/m ³	594 - 992.6
Age	Age of the concrete sample when tested.	days	1 - 365

X'X matrix (9x9) scaled by 10^7 :

```
[[ 0.    0.03  0.01  0.01  0.02  0.    0.1   0.08  0.   ]
 [ 0.03  9.27  1.88  1.3   5.24  0.19 28.08 22.21  1.38]
 [ 0.01  1.88  1.33  0.23  1.4   0.05  7.21  5.69  0.32]
 [ 0.01  1.3   0.23  0.72  0.98  0.05  5.43  4.36  0.19]
 [ 0.02  5.24  1.4   0.98  3.44  0.11 18.16 14.39  0.89]
 [ 0.    0.19  0.05  0.05  0.11  0.01  0.61  0.51  0.02]
 [ 0.1   28.08  7.21  5.43 18.16  0.61 98.12 77.41  4.57]
 [ 0.08 22.21  5.69  4.36 14.39  0.51 77.41 62.3   3.56]
 [ 0.    1.38  0.32  0.19  0.89  0.02  4.57  3.56  0.63]]
```

X'Y vector (9x1) scaled by 10^7 : [0. 1.13 0.29 0.19 0.66 0.03 3.57 2.83 0.2]

The Normal Equations are: $X'X * B = X'Y$

Where B is the vector of regression coefficients (intercept + slopes).


```
#generate OLS regression results for all features
```

```
import statsmodels.api as sm
```

```
X_sm = sm.add_constant(X)
```

```
model = sm.OLS(y,X_sm)
```

```
print(model.fit().summary())
```

OLS Regression Results

```
=====
```

Dep. Variable:	csMPa	R-squared:	0.616
Model:	OLS	Adj. R-squared:	0.613
Method:	Least Squares	F-statistic:	204.3
Date:	Fri, 15 Oct 2021	Prob (F-statistic):	6.29e-206
Time:	16:43:15	Log-Likelihood:	-3869.0
No. Observations:	1030	AIC:	7756.
Df Residuals:	1021	BIC:	7800.
Df Model:	8		
Covariance Type:	nonrobust		

```
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	-23.3312	26.586	-0.878	0.380	-75.500	28.837
cement	0.1198	0.008	14.113	0.000	0.103	0.136
slag	0.1039	0.010	10.247	0.000	0.084	0.124
flyash	0.0879	0.013	6.988	0.000	0.063	0.113
water	-0.1499	0.040	-3.731	0.000	-0.229	-0.071
superplasticizer	0.2922	0.093	3.128	0.002	0.109	0.476
coarseaggregate	0.0181	0.009	1.926	0.054	-0.000	0.037
fineaggregate	0.0202	0.011	1.887	0.059	-0.001	0.041
age	0.1142	0.005	21.046	0.000	0.104	0.125

```
=====
```

Projeção e predição

- Se as variáveis (colunas) da matriz X realmente servirem para prever o valor de Y e
- se o modelo de regressão linear for uma boa aproximação para o relacionamento das variáveis,

- então esperamos que

$$Y \approx \hat{Y} = X\hat{\beta}$$

- Como medir o grau de aproximação?
- É possível obter uma decomposição do vetor Y em componentes ortogonais. A partir daí extraímos uma medida de qualidade do ajuste.

Decomposição em soma de quadrados

- Seja $\bar{y} = \sum_i y_i / 1500$, o preço médio dos 1500 apartamentos.
- Defina o vetor 1500×1 dado por $\bar{Y} = (\bar{y}, \bar{y}, \dots, \bar{y})' = \bar{y}(1, 1, \dots, 1)'$
- O vetor Y pode ser escrito como

$$Y = \hat{Y} + (Y - \hat{Y}) = \bar{Y} + \hat{Y} - \bar{Y} + (Y - \hat{Y})$$

- Isto é,

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{1499} \\ y_{1500} \end{bmatrix} = \begin{bmatrix} \bar{y} \\ \bar{y} \\ \vdots \\ \bar{y} \\ \bar{y} \end{bmatrix} + \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_{1499} \\ \hat{y}_{1500} \end{bmatrix} - \begin{bmatrix} \bar{y} \\ \bar{y} \\ \vdots \\ \bar{y} \\ \bar{y} \end{bmatrix} + \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_{1499} - \hat{y}_{1499} \\ y_{1500} - \hat{y}_{1500} \end{bmatrix}$$

Decomposição em soma de quadrados

- Isto é,

$$\begin{bmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_{1499} - \bar{y} \\ y_{1500} - \bar{y} \end{bmatrix} = \begin{bmatrix} \hat{y}_1 - \bar{y} \\ \hat{y}_2 - \bar{y} \\ \vdots \\ \hat{y}_{1499} - \bar{y} \\ \hat{y}_{1500} - \bar{y} \end{bmatrix} + \begin{bmatrix} y_1 - \hat{y}_1 \\ y_2 - \hat{y}_2 \\ \vdots \\ y_{1499} - \hat{y}_{1499} \\ y_{1500} - \hat{y}_{1500} \end{bmatrix}$$

- $\mathbf{Y} - \bar{y}\mathbf{1} = (\hat{\mathbf{Y}} - \bar{y}\mathbf{1}) + (\mathbf{Y} - \hat{\mathbf{Y}})$
- Os vetores do lado direito são ortogonais um ao outro. Em consequência,

$$\|\mathbf{Y} - \bar{y}\mathbf{1}\|^2 = \|\hat{\mathbf{Y}} - \bar{y}\mathbf{1}\|^2 + \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2$$

Remember

$$\hat{\beta}_{(p+1) \times 1} = (X'X)^{-1} X'Y$$

@

$$\hat{Y}_{n \times 1} = X \hat{\beta} = \underbrace{X (X'X)^{-1} X'}_H Y = HY$$

H = hat matrix (put the hat on top of Y)

H is $n \times n$ ~~Projection matrix~~

Properties of H:

Properties of $H = X(X'X)^{-1}X'$

(5)

i) H is $n \times n$

ii) H is the \perp projection matrix

\perp Projects where?

\rightarrow Im $\underbrace{M(X)}$
linear combination
of columns of X

iii) $H^2 = H$ (show this)

iv) H is symmetric (show this)

Fact If $\vec{v} \perp \vec{w}$ then

©

$$\|\vec{v} + \vec{w}\|^2 = \|\vec{v}\|^2 + \|\vec{w}\|^2$$

Proof: $\|\vec{v} + \vec{w}\|^2 = \langle \vec{v} + \vec{w}, \vec{v} + \vec{w} \rangle$

$$= (\vec{v} + \vec{w})^T (\vec{v} + \vec{w})$$

$$= (v^T + w^T) \cdot (v + w)$$

$$= v^T v + \underbrace{v^T w}_{\overset{0}{\parallel}} + \underbrace{w^T v}_{\overset{0}{\parallel}} + w^T w$$

$$= \|v\|^2 + \|w\|^2 \quad \overset{0}{\parallel} \text{ pois são } \perp \text{'s}$$

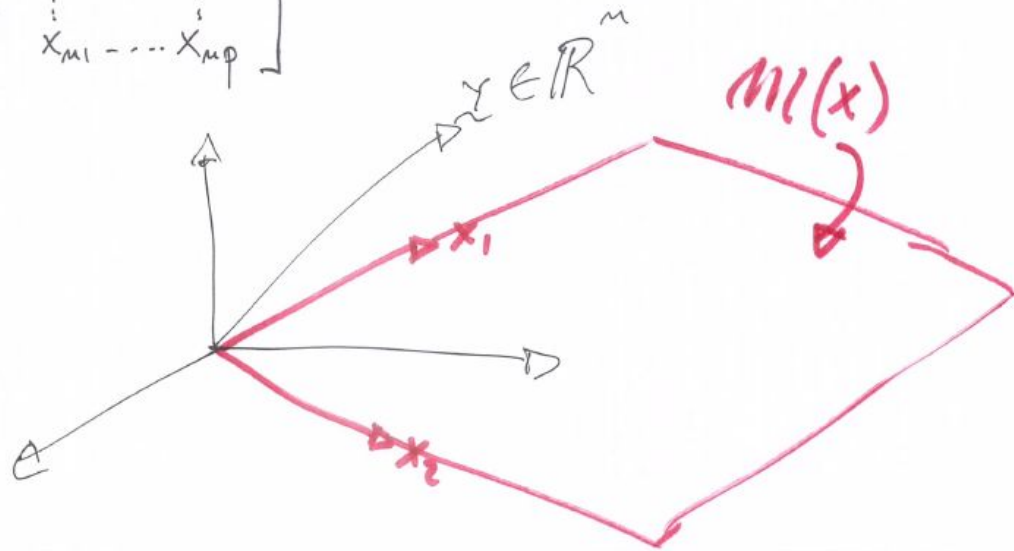
Assumption: Column $\underline{1}$ is in X ①

$$X = \begin{bmatrix} \underline{1} & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ \underline{1} & x_{m1} & \dots & x_{mp} \end{bmatrix}$$

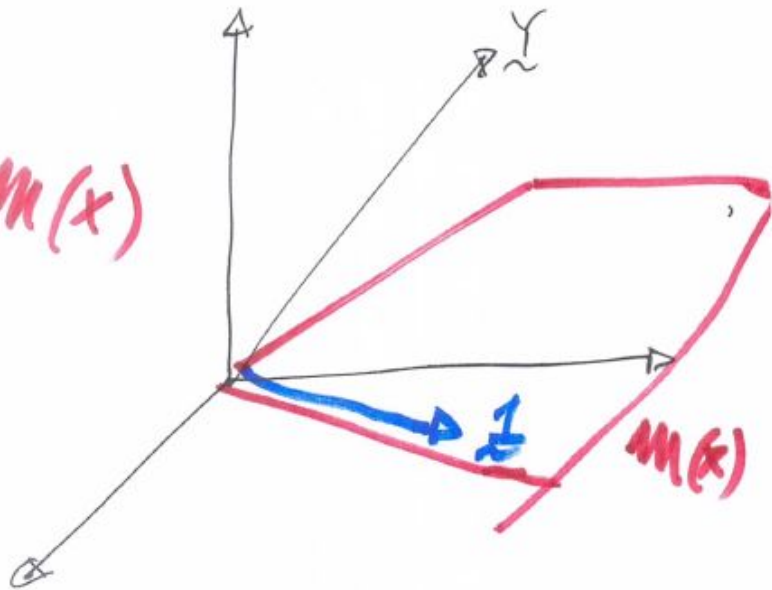
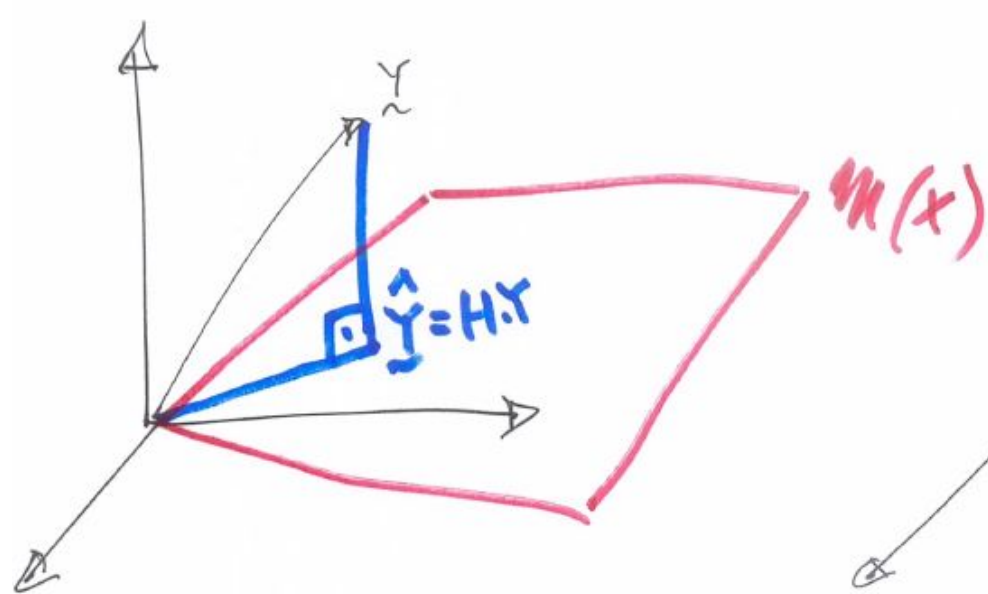
$m \times (p+1)$

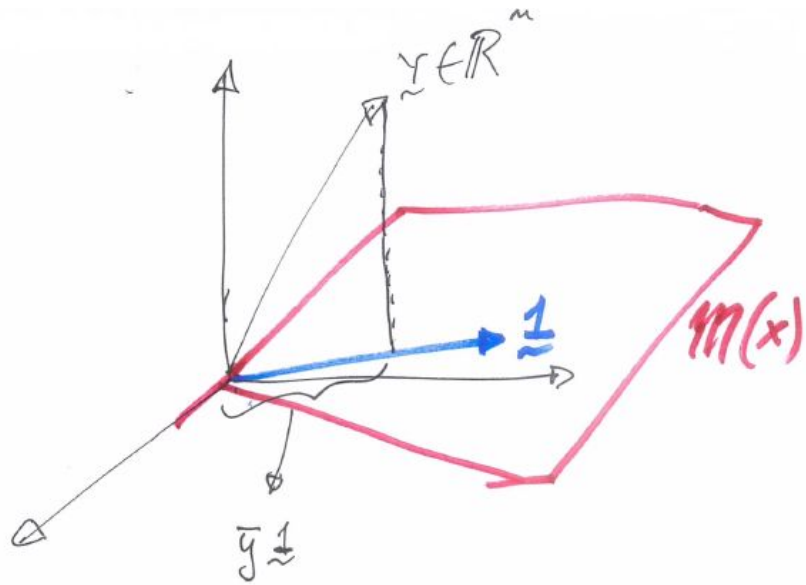
$$\underline{1} \in \mathcal{M}(X)$$

$$\underline{1} = 1 \cdot \underline{1} + 0 \cdot \underline{x}_1 + \dots + 0 \cdot \underline{x}_p$$



(e)





$$\begin{aligned}
 \bar{y} \cdot \underline{1} &= \left(\frac{1}{n} \sum y_i \right) \cdot \underline{1} = \frac{1}{n} \underbrace{\left(\underline{1}' \cdot \underline{y} \right)}_{\substack{\parallel \\ (1, \dots, 1) \cdot \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}}} = \frac{1}{n} \underline{1} \cdot \left(\underline{1}' \cdot \underline{y} \right) \\
 &= \frac{1}{n} \left(\underline{1} \cdot \underline{1}' \right) \cdot \underline{y} \\
 &= U \cdot \underline{y} \\
 &\quad \hookrightarrow \text{matriz associada} \\
 &\quad \text{com } \underline{1}
 \end{aligned}$$

é um escalar

→ U é matriz de projecção \perp no vetor $\underline{1}$ (9)

Prova: basta mostrar que, para todo $Y \in \mathbb{R}^n$, temos

$$(UY) \perp (Y - UY).$$

$$(UY)' \cdot (Y - UY) = (Y' U') \cdot ((I - U) \cdot Y) = Y' (U' (I - U)) \cdot Y$$

$$= Y' (U' - U'U) Y.$$

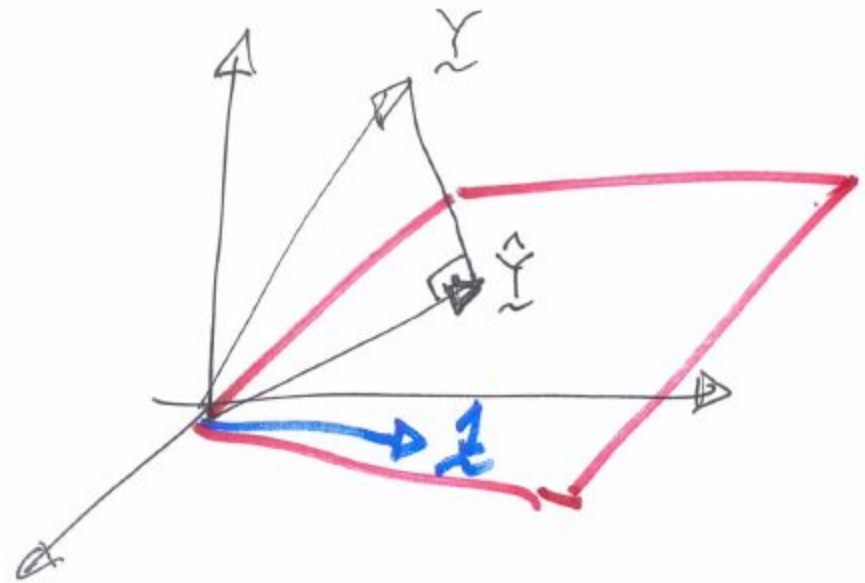
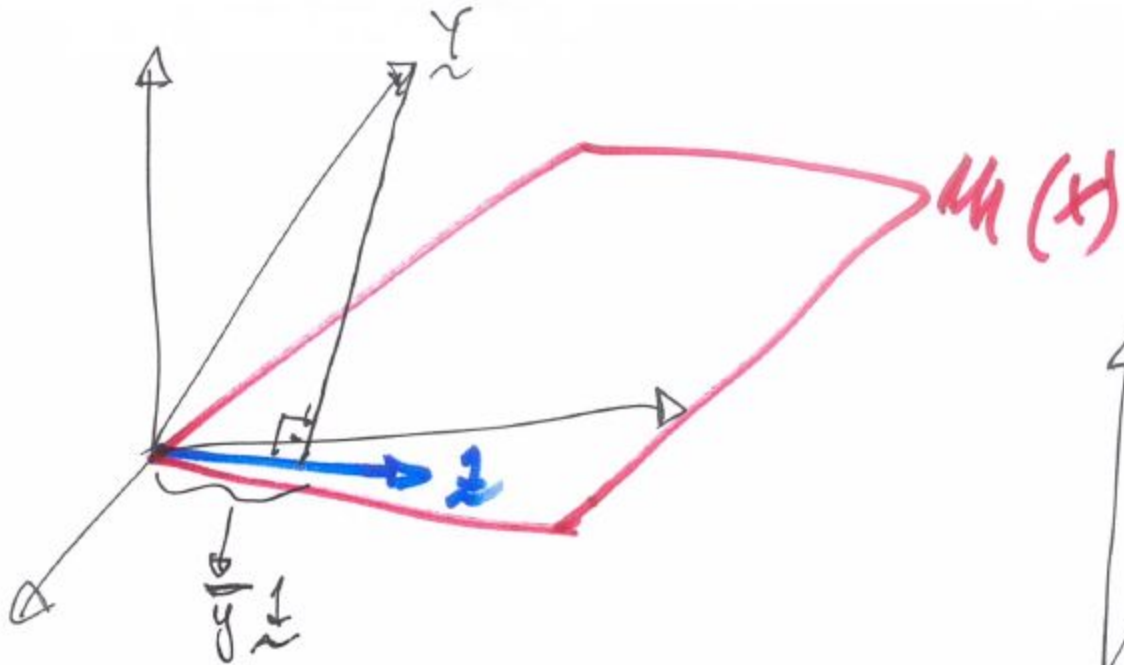
Mas $U' = U$ (show this)

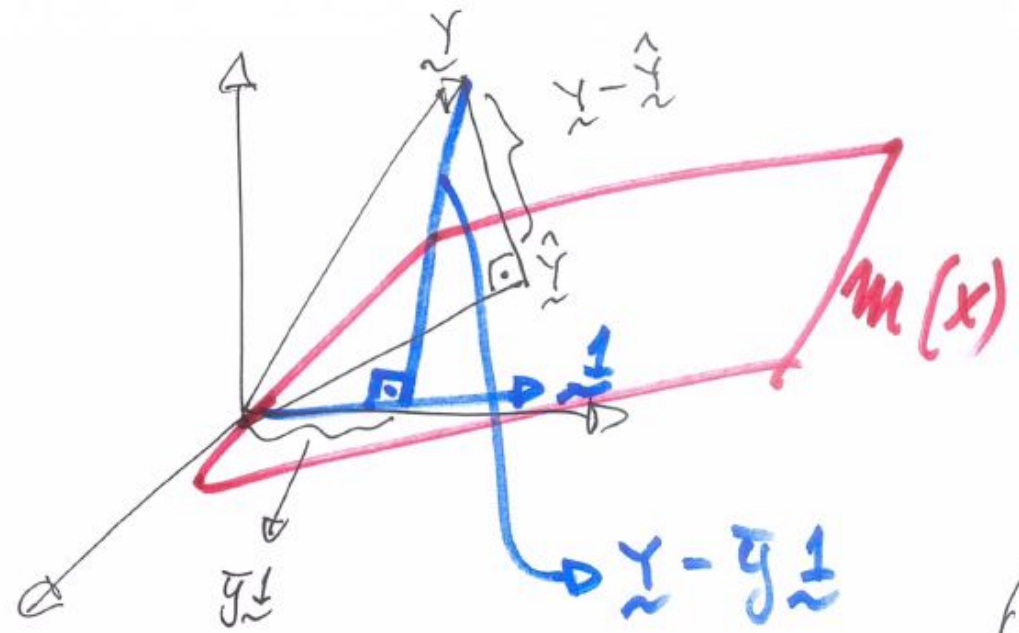
$U^2 = U$ (show this)

$$\Rightarrow U' - U'U = U - U^2 = 0$$

~~q.e.d.~~
q.e.d.

(h)





~~Y - \bar{y}~~

$$Y - \bar{y} = \underbrace{\hat{Y} - \bar{y}} + \underbrace{Y - \hat{Y}}$$

~~Y - \bar{y}~~

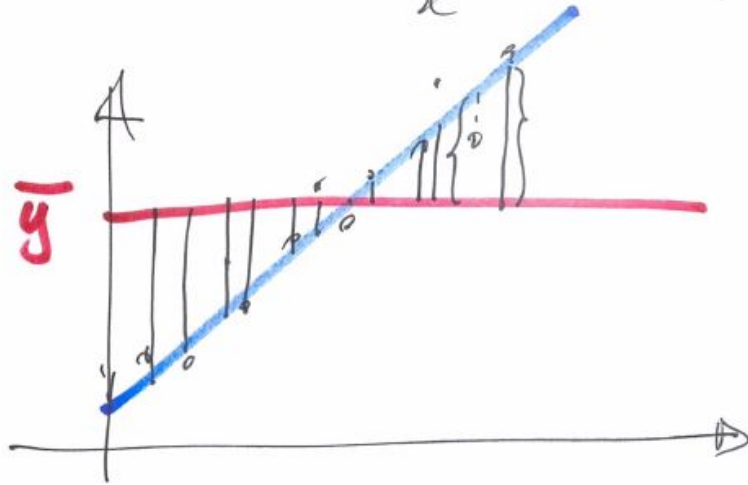
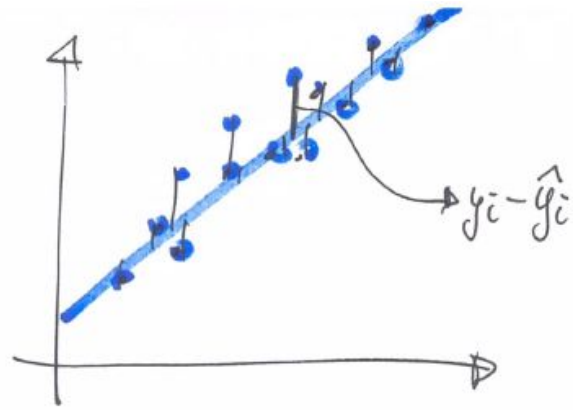
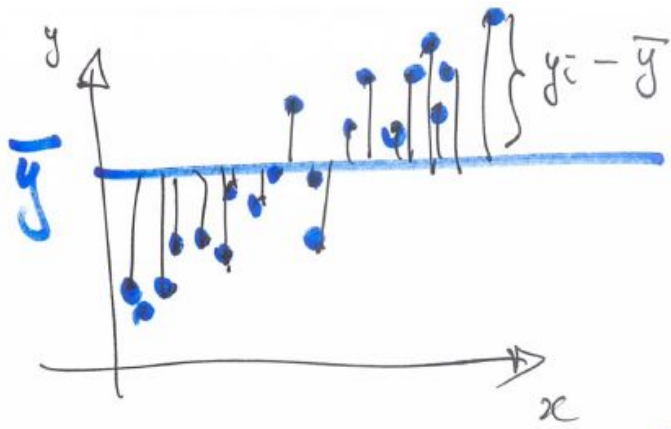
(identifique $\hat{Y} - \bar{y}$ na figura)

Assum, $(H-U)Y \perp (I-H)Y$ e

(K)

$$\| \tilde{Y} - \bar{y} \mathbf{1} \|^2 = \| \hat{Y} - \bar{y} \mathbf{1} \|^2 + \| Y - \hat{Y} \|^2$$

$$\underbrace{\sum_i (y_i - \bar{y})^2}_{SS_{Total}} = \underbrace{\sum_i (\hat{y}_i - \bar{y})^2}_{SS_{Reg}} + \underbrace{\sum (y_i - \hat{y}_i)^2}_{SS_{Res}}$$



$$\hat{y}_i - \bar{y}$$

$$(\text{beta} - \bar{y})$$

The sum of squares

- When the residual vector

$$\|\mathbf{r}\|^2 = \|\mathbf{Y} - \hat{\mathbf{Y}}\|^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

is small, we have a good fit.

- The idea is to compare this remaining variability with the original variability in \mathbf{Y} BEFORE any regressors were considered.
- The variation of \mathbf{Y} around \bar{y} , the mean of \mathbf{Y} , is equal to:

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \|\mathbf{Y} - \bar{y}\mathbf{1}\|^2$$

Finally, the R^2

- That is, we consider the ratio

$$\frac{\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2}{\|\mathbf{Y} - \bar{y}\mathbf{1}\|^2}$$

- If we have a good fit, we should have this ratio close to zero.
- We can prove that this ratio is always smaller than 1.
- Hence, it is more common to use R^2 :

$$R^2 = 1 - \frac{\|\mathbf{Y} - \hat{\mathbf{Y}}\|^2}{\|\mathbf{Y} - \bar{y}\mathbf{1}\|^2}$$

- A good fit should have $R^2 \approx 1$.

```
#generate OLS regression results for all features
```

```
import statsmodels.api as sm
```

```
X_sm = sm.add_constant(X)
```

```
model = sm.OLS(y,X_sm)
```

```
print(model.fit().summary())
```

OLS Regression Results

```
=====
Dep. Variable:          csMPa   R-squared:          0.616
Model:                 OLS     Adj. R-squared:     0.613
Method:                Least Squares   F-statistic:        204.3
Date:                  Fri, 15 Oct 2021   Prob (F-statistic): 6.29e-206
Time:                  16:43:15          Log-Likelihood:     -3869.0
No. Observations:     1030             AIC:                7756.
Df Residuals:         1021             BIC:                7800.
Df Model:              8
Covariance Type:      nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	-23.3312	26.586	-0.878	0.380	-75.500	28.837
cement	0.1198	0.008	14.113	0.000	0.103	0.136
slag	0.1039	0.010	10.247	0.000	0.084	0.124
flyash	0.0879	0.013	6.988	0.000	0.063	0.113
water	-0.1499	0.040	-3.731	0.000	-0.229	-0.071
superplasticizer	0.2922	0.093	3.128	0.002	0.109	0.476
coarseaggregate	0.0181	0.009	1.926	0.054	-0.000	0.037
fineaggregate	0.0202	0.011	1.887	0.059	-0.001	0.041
age	0.1142	0.005	21.046	0.000	0.104	0.125

Matriz de projeção ortogonal

Seja X uma matriz de números reais de dimensão $n \times (p + 1)$, seja $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ um vetor $(p + 1) \times 1$ e \mathbf{y} um vetor $n \times 1$.

Sejam v_1, \dots, v_k vetores do \mathbb{R}^n . Verifique que o conjunto das combinações lineares desses vetores forma um sub-espço vetorial do \mathbb{R}^n .

- Verifique que $X\beta = \beta_0 X_0 + \beta_1 X_1 + \dots + \beta_p X_p$ onde X_0, X_1, \dots, X_p são os vetores colunas de X . Assim, o conjunto $\mathfrak{M}(X)$ das combinações lineares das colunas de X é igual a $\mathfrak{M}(X) = \{X\beta \mid \beta \in \mathbb{R}^{p+1}\}$.

Seja X uma matriz de números reais de dimensão $n \times (p + 1)$, seja $\beta = (\beta_0, \beta_1, \dots, \beta_p)$ um vetor $(p + 1) \times 1$ e \mathbf{y} um vetor $n \times 1$.

- Verifique que $X\beta = \beta_0 X_0 + \beta_1 X_1 + \dots + \beta_p X_p$ onde X_0, X_1, \dots, X_p são os vetores colunas de X . Assim, o conjunto $\mathfrak{M}(X)$ das combinações lineares das colunas de X é igual a $\mathfrak{M}(X) = \{X\beta \mid \beta \in \mathbb{R}^{p+1}\}$.
- Seja W um subespaço do espaço vetorial V . Definimos o espaço ortogonal de W como sendo

$$W^\perp = \{u \in V \mid \langle u, w \rangle = 0 \quad \forall w \in W\}$$

Mostre que W^\perp é um subespaço vetorial de V .

Solução: $\vec{0} \in W^\perp$ pois $\langle \vec{0}, w \rangle = 0$ para todo $w \in W$. E também $\langle a_1 u_1 + a_2 u_2, w \rangle = 0$ se $\langle u_1, w \rangle = 0$ e $\langle u_2, w \rangle = 0$

A matriz de projeção ortogonal

- Seja $H = X(X'X)^{-1}X'$ de dimensão $n \times n$. Verifique que H é idempotente ($H^2 = H$) e simétrica ($H' = H$).

A matriz de projeção ortogonal

- Seja $H = X(X'X)^{-1}X'$ de dimensão $n \times n$. Verifique que H é idempotente ($H^2 = H$) e simétrica ($H' = H$).
- Seja $y \in \mathbb{R}^n$. A matriz P de dimensão $n \times n$ é dita de projeção ortogonal num certo subespaço vetorial se $y - Py \perp Py$ para todo $y \in \mathbb{R}^n$. Mostre que $H = X(X'X)^{-1}X'$ é uma matriz de projeção ortogonal usando que H é idempotente e simétrica.

A matriz de projeção ortogonal

- Seja $H = X(X'X)^{-1}X'$ de dimensão $n \times n$. Verifique que H é idempotente ($H^2 = H$) e simétrica ($H' = H$).
- Seja $y \in \mathbb{R}^n$. A matriz P de dimensão $n \times n$ é dita de projeção ortogonal num certo subespaço vetorial se $y - Py \perp Py$ para todo $y \in \mathbb{R}^n$. Mostre que $H = X(X'X)^{-1}X'$ é uma matriz de projeção ortogonal usando que H é idempotente e simétrica.
- Como $H = X(X'X)^{-1}X'$ é uma matriz de projeção ortogonal, resta saber em que sub-espaço vetorial W a matriz H projeta os vetores $y \in \mathbb{R}^n$. Mostre que H projeta ortogonalmente em $\mathfrak{M}(X)$ (isto é, mostre que $W = \mathfrak{M}(X)$.)

Solução: Para todo $y \in \mathbb{R}^n$, temos $Hy = X(X'X)^{-1}X'y = Xb$ onde $b = (X'X)^{-1}X'y$. Assim, $Hy \in \mathfrak{M}(X)$ para todo y e portanto $W \subset \mathfrak{M}(X)$. Por outro lado, tome um elemento Xb qualquer de $\mathfrak{M}(X)$. Por definição, $Hy \in W$ para todo y . Em particular, tomando $y = Xb$, temos então $HXb \in W$. Mas $HXb = X(X'X)^{-1}X'Xb = Xb$. Isto é, $Xb \in W$ e portanto $\mathfrak{M}(X) \subset W$. Concluimos então que $W = \mathfrak{M}(X)$.

LS = projeção ortogonal

- Seja $H = X(X'X)^{-1}X'$ a matriz de projeção ortogonal no espaço $\mathfrak{M}(X)$ das combinações lineares das colunas de X . Mostre que ao escolher β tal que $X\beta = Hy$ estamos minimizando a distância $\|y - X\beta\|^2$. DICA: Escreva some e subtraia Hy em $\|y - X\beta\|^2$ e use que $\|v\|^2 = \langle v, v \rangle$

Solução: $\|y - X\beta\|^2 = \langle y - X\beta, y - X\beta \rangle$. Somando e subtraindo Hy obtemos

$$\begin{aligned}\|y - X\beta\|^2 &= \langle y - Hy + Hy - X\beta, y - Hy + Hy - X\beta \rangle \\ &= \langle y - Hy, y - Hy \rangle + \langle Hy - X\beta, Hy - X\beta \rangle - 2 \langle y - Hy, Hy - X\beta \rangle \\ &= \|y - Hy\|^2 + \|Hy - X\beta\|^2 + 0.\end{aligned}$$

O último termo acima é zero pois $Hy - X\beta \in \mathfrak{M}(X)$ já que $Hy \in \mathfrak{M}(X)$ e $X\beta \in \mathfrak{M}(X)$ e o conjunto $\mathfrak{M}(X)$ é um sub-espaço vetorial (e portanto contém a diferença dos vetores). Além disso, $y - Hy \in \mathfrak{M}(X)^\perp$. Portanto, o produto interno $\langle y - Hy, Hy - X\beta \rangle$ é nulo.

Assim, $\|y - X\beta\|^2 = \|y - Hy\|^2 + \|Hy - X\beta\|^2$. O primeiro termo do lado direito não depende de β e o segundo é não-negativo. Ele será minimizado se for igual a zero, o que ocorre se tomamos $X\beta = Hy$.