

A.I. in health informatics  
**lecture 6** pattern classification  
(and regression)

kevin small &  
byron wallace

# today

- feature-spaces
- regression
  - least-squares
- ~~classification~~
  - decision trees

# what's in a feature-space?

- machine learning algorithms operate over vectors
- a *feature-space* is the vector-space used to represent things for ml algorithms
- suppose we want to predict housing prices; what attributes would you use and how would you encode them?

# what's in a feature-space?

- feature-engineering isn't considered 'sexy'
  - but in practice it is probably the most important contributor to performance
- ***representation is everything!***

# what's in a feature-space?

- want to classify people as 'diabetic' or not
- have the following information
  - their name
  - their favorite color
  - the month they were born



# what's in a feature-space?

suppose we are classifying objects

$$\{x_0 \ x_1 \ \dots \ x_n\}$$

a feature mapping **F** maps from the input into its vector representation (in the feature-space)

# text encoding

- many applications in health informatics involve text
- we will review 'bag-of-words' style text encoding (or, how to vectorize text)

# text encoding by example

suppose we want to encode the following sentences:

$S_1$  = "Boston drivers are frequently aggressive"

$S_2$  = "The Boston Red Sox frequently hit line drives"

# text encoding by example

## 1. remove stop-words

$S_1$  = "Boston drivers ~~are~~ frequently aggressive"

$S_2$  = "~~The~~ Boston Red Sox frequently hit line drives"

# text encoding by example

## 2. lower-case

$S_1 =$  "boston drivers ~~are~~ frequently aggressive"

$S_2 =$  "~~The~~ boston red sox frequently hit line drives"

# text encoding by example

## 3. stem

$S_1 =$  "boston drive~~s~~ ~~are~~ frequentl~~y~~ aggressiv~~e~~"

$S_2 =$  "The boston red sox frequentl~~y~~ hit line drive~~s~~"

# text encoding: indicator vectors

$S_1 =$  "*boston drivers* ~~are~~ ~~the~~ frequently ~~are~~ aggressive"

$S_2 =$  "~~The~~ *boston red sox* frequently ~~hit~~ line drive~~s~~"

$V =$  [hit, red, sox, line, boston, frequent, drive, aggressive]

# text encoding: indicator vectors

	hit	red	sox	line	boston	frequent	drive	aggressive
$S_1 =$	0	0	0	0	1	1	1	1
$S_2 =$	1	1	1	1	1	1	1	0

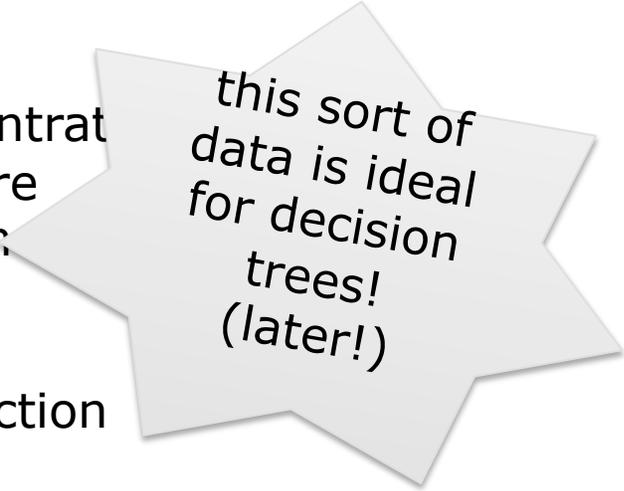
a new sentence,  $S_3$ , comes along it reads: "I hate the red sox". to which sentence is it most similar?

$S_3 =$	0	1	1	0	0	0	0	0
---------	---	---	---	---	---	---	---	---

# other sorts of spaces

- text tends to be *high-dimensional* and *sparse*
- pima dataset

0 # of times pregnant  
1 plasma glucose concentrat  
2 diastolic blood pressure  
3 triceps skin fold thickn  
4 2-Hour serum insulin  
5 body mass index  
6 diabetes pedigree function  
7 age  
**label** -1 or 1



this sort of  
data is ideal  
for decision  
trees!  
(later!)

# other sorts of spaces

<b>feature</b>	<b>person P</b>
# of times pregnant	3
plasma glucose concentration	19
diastolic blood pressure	4.4
triceps skin fold thickness	5
2-hour serum insulin	34
BMI	2
diabetes pedigree function	32
age	2
	24


$$\mathbf{F}(\mathbf{P}) = [3, 19, 4.4, 5, 34, 2, 32, 2, 24]$$

# normalization, etc

$$\mathbf{F}(\mathbf{P}) = [3, \mathbf{19}, 4.4, 5, \mathbf{34}, 2, \mathbf{32}, 2, \mathbf{24}]$$

- **some** features are getting undue weight
- solution: normalize each column (eg., divide all entries in each column by the max entry in that column)

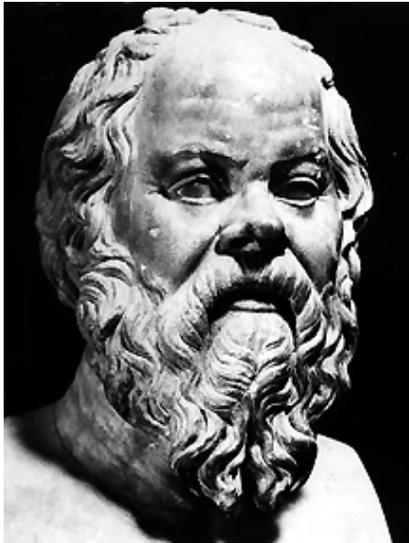
# standardizing

- in addition to column normalization, we sometimes want to row-standardize, ie., enforce a zero mean unit variance on each instance  $\mathbf{x}$

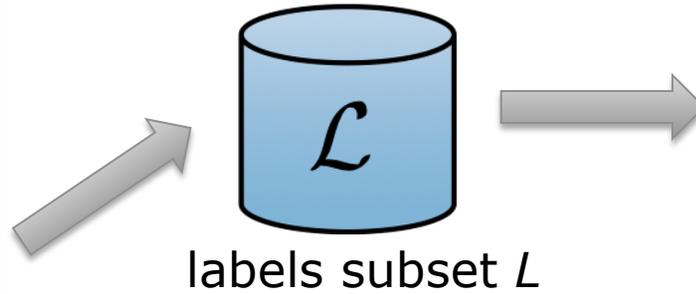
$$x'_i = \frac{x_i - \bar{x}_i}{sd(x)}$$

not necessary in general, but required for some algorithms, eg., PCA

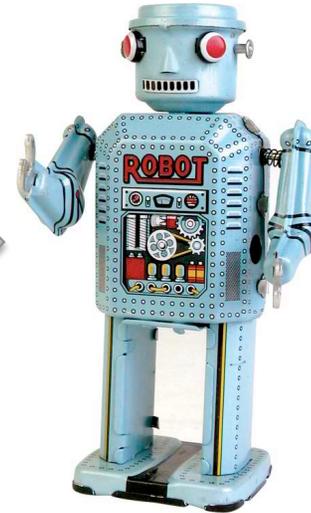
# supervised learning



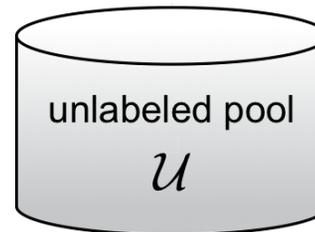
(human) expert / labeler



labels subset  $L$



classifier  $C$  is induced over  $L$



use  $C$  to classify unlabeled examples

# generalized supervised learning

$$\arg \min_{\theta} \underbrace{(h(x, \theta) - y)^2}_{\text{minimize error}} + \underbrace{\lambda R(h, \theta)}_{\text{penalize complexity}}$$

The diagram shows the equation  $\arg \min_{\theta} (h(x, \theta) - y)^2 + \lambda R(h, \theta)$ . A horizontal curly brace is placed under the term  $(h(x, \theta) - y)^2$ , with an arrow pointing from the text 'minimize error' below to the center of the brace. A second horizontal curly brace is placed under the term  $\lambda R(h, \theta)$ , with an arrow pointing from the text 'penalize complexity' below to the center of the brace.

# supervised learning

- two main types of learning tasks:
  - regression
  - classification
- today start with regression, then do classification

# regression

- the regression task:

given a **training set**  $\{(\mathbf{x}_0, y_0), (\mathbf{x}_1, y_1) \dots (\mathbf{x}_m, y_m)\}$   
where  $y_i$ s are scalars (in  $\mathbf{R}$ ), induce a model  $\mathbf{h}$  such  
that given a new instance  $\mathbf{x}_{m+1}$ ,  $\mathbf{h}(\mathbf{x}_{m+1})$  is  
minimized

- eg., predict the price of a house given its attributes; predict someone's BMI given some of their physical characteristics...

# regression

- many possible ways to do this
- assume that we're interested in the linear case  $\mathbf{w} \cdot \mathbf{x}_i = y_i$ 
  - doesn't always make sense (y could be, e.g., polynomial in  $\mathbf{w}$ )

# very brief linear algebra review

- vector  $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]$        $\mathbf{x}^t = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$

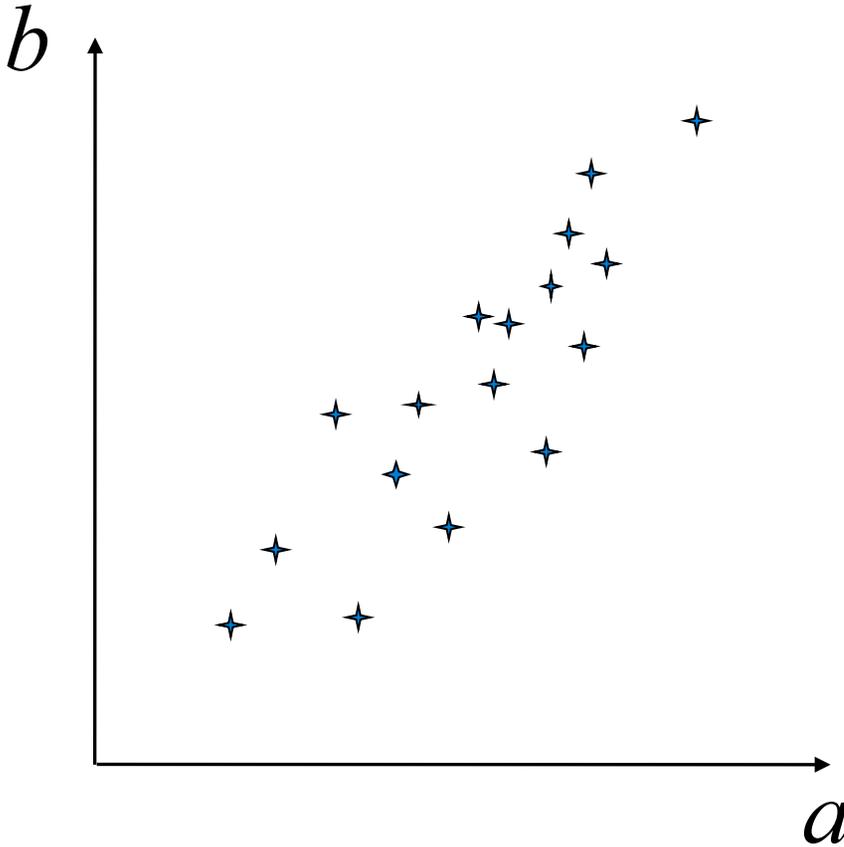
- dot product  $\mathbf{a} \cdot \mathbf{b} = \sum_i \mathbf{a}_i \cdot \mathbf{b}_i$

- m x n matrix  $\mathbf{X} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1n} \\ X_{21} & X_{22} & \dots & X_{2n} \\ \dots & & & \\ X_{m1} & X_{m2} & \dots & X_{mn} \end{bmatrix}$

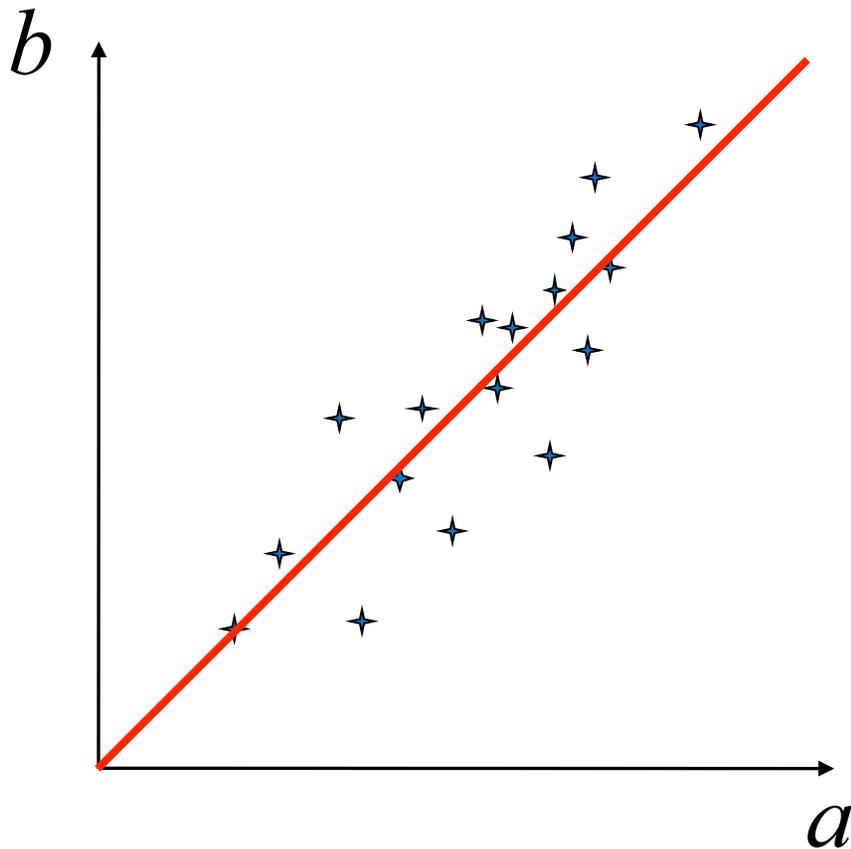
***acknowledgement***: many of the slides from here on (those regarding least-squares linear regression) are borrowed heavily from

*[graphics.stanford.edu/~jplewis/lscourse/ols\\_slides.ppt](http://graphics.stanford.edu/~jplewis/lscourse/ols_slides.ppt)*

# one-dimensional regression



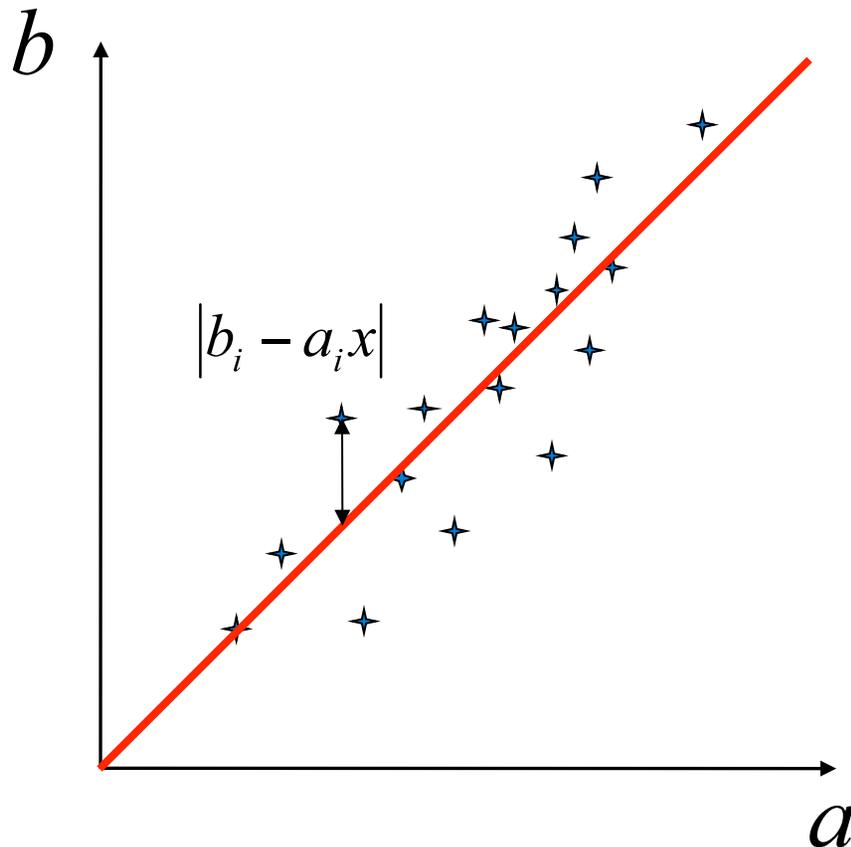
# one-dimensional regression



aim: find a line that  
represent the "best" linear  
relationship:

$$b = ax$$

# one-dimensional regression

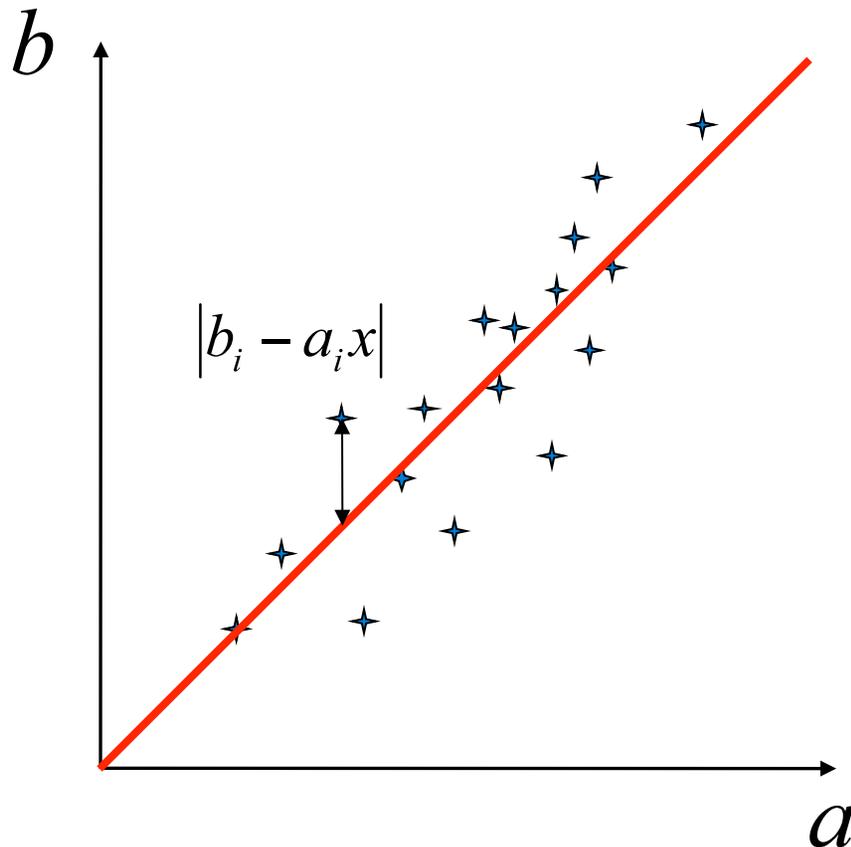


**problem** the data does not go through a line

the error is:

$$e_i = b_i - a_i x$$

# one-dimensional regression



**problem** the data does not go through a line

$$e_i = b_i - a_i x$$

**objective** find the line that minimizes the sum:

$$\sum_i (b_i - a_i x)^2$$

... in terms of our general aim

$$\arg \min_{\theta} \underbrace{(h(x, \theta) - y)^2}_{\text{minimize error}} + \lambda \underbrace{R(h, \theta)}_{\text{penalize complexity}}$$

in the case of (one-dimensional) LS linear regression:

$$\arg \min_x \sum_i (b_i - a_i x)^2$$

(just sort of ignoring model complexity)

re-expressed in matrix notation

$$\begin{aligned} & \sum_i (b_i - a_i x)^2 \\ &= (\mathbf{b} - x\mathbf{a})^T (\mathbf{b} - x\mathbf{a}) \\ &= \underbrace{\|\mathbf{b} - x\mathbf{a}\|}^2 \end{aligned}$$

this is what we want to minimize!

# more generally (multidimensional linear regression)

model with  $m$  parameters

$$b = a_1 x_1 + \dots + a_m x_m = \sum_j a_j x_j$$

and  $n$  measurements (examples)

$$e(\mathbf{x}) = \sum_{i=1}^n (b_i - \sum_{j=1}^m a_{i,j} x_j)^2$$

$$= \left\| \mathbf{b} - \begin{bmatrix} \sum_{j=1}^m a_{i,j} x_j \end{bmatrix} \right\|^2$$

find  $\mathbf{x}$  that  
minimizes this!

$$e(\mathbf{x}) = \|\mathbf{b} - \mathbf{Ax}\|^2$$

# minimizing the error

$$\mathbf{b} - \mathbf{Ax} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} - \begin{bmatrix} a_{1,1} & \dots & a_{1,m} \\ \vdots & & \vdots \\ a_{n,1} & \dots & a_{n,m} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$$

$$= \begin{bmatrix} b_1 - (a_{1,1}x_1 + \dots + a_{1,m}x_m) \\ \vdots \\ b_n - (a_{n,1}x_1 + \dots + a_{n,m}x_m) \end{bmatrix}$$

$$\mathbf{b} - \mathbf{Ax} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} - \begin{bmatrix} a_{1,1} & \dots & a_{1,m} \\ \vdots & & \vdots \\ a_{n,1} & \dots & a_{n,m} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

parameter 1

instance  $n$

$$= \begin{bmatrix} b_1 - (a_{1,1}x_1 + \dots + a_{1,m}x_m) \\ \vdots \\ b_n - (a_{n,1}x_1 + \dots + a_{n,m}x_m) \end{bmatrix}$$

solving for  $\mathbf{x}$

$$\begin{aligned} e &= (\mathbf{b} - \mathbf{x}\mathbf{A})^T (\mathbf{b} - \mathbf{x}\mathbf{A}) \\ &= \mathbf{b}^T \mathbf{b} - 2\mathbf{b}^T \mathbf{A}\mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A}\mathbf{x} \end{aligned}$$

# solving for $\mathbf{x}$

as usual, we take the derivative, set it to 0 and solve for the variable of interest ( $\mathbf{x}$ )

$$\frac{\partial e}{\partial \mathbf{x}} = -2\mathbf{A}^t \mathbf{b} + 2\mathbf{A}^t \mathbf{A} \mathbf{x}$$

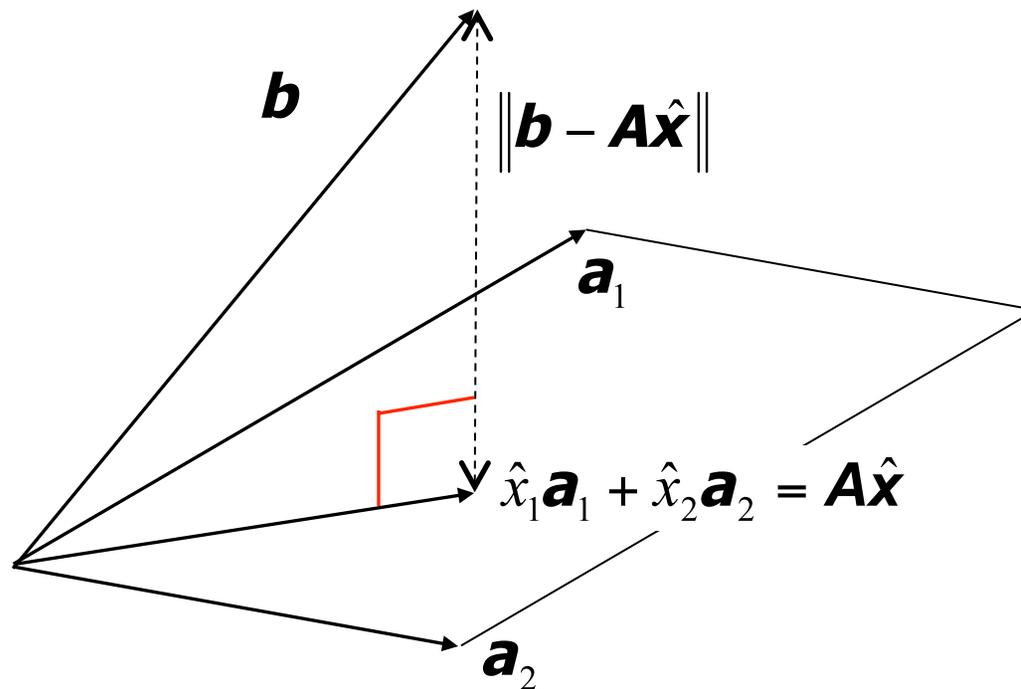
$$0 = -2\mathbf{A}^t \mathbf{b} + 2\mathbf{A}^t \mathbf{A} \mathbf{x}$$

$$\rightarrow \mathbf{A}^t \mathbf{A} \mathbf{x} = \mathbf{A}^t \mathbf{b}$$

$$\rightarrow \hat{\mathbf{x}} = (\mathbf{A}^t \mathbf{A})^{-1} \mathbf{A}^t \mathbf{b}$$

# the geometric interpretation

- $\mathbf{A}\hat{\mathbf{x}}$  is the orthogonal projection of  $\mathbf{b}$  onto  $\text{range}(\mathbf{A})$   
 $\Leftrightarrow \mathbf{A}^T(\mathbf{b} - \mathbf{A}\hat{\mathbf{x}}) = \mathbf{0} \Leftrightarrow \mathbf{A}^T \mathbf{A}\hat{\mathbf{x}} = \mathbf{A}^T \mathbf{b}$



awesome!

wait, what were we doing again?

given a new instance,  $\mathbf{a}$  our prediction is:

$$\hat{\mathbf{x}} \cdot \mathbf{a}$$

# what about penalizing for complexity?

- here we sort of ignored model complexity
- how might we penalize for this?
- **Lasso** constrain the norm of  $\mathbf{x}$ , our parameter vector, to be  $\leq r$

# that's it for regression

- often we want to *classify* things into categories, as opposed to predicting a continuous value
- this is known as *classification* (next class!)