

Lesson Plan 4/5

Thursday, April 5, 2018 11:31 AM

Admin

- Project 2 due on Friday
- HW4 due on Tuesday
- Grades for Project 1 posted
- HW3 grading almost done

7

HDSC:

- Example wrangling with R and Python: <https://www.superdatascience.com/wrangling-in-data-r-python/>

Project 2 Questions?

- Added note re. dates to project description

HW4

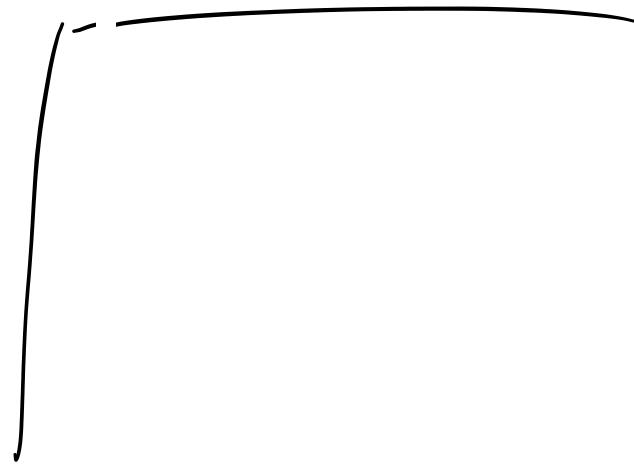
- Questions?

Data Analysis with Geometry

- Preliminaries
- Geometry and distances
- KNN classifier
- Some vector algebra

Linear Regression

- Formulation
- Interpretation



Notation:

y : outcome (continuous numeric)

X : predictors

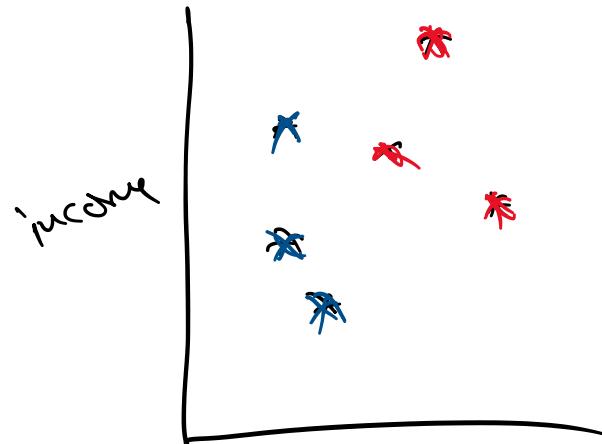
G : categorical
outcome

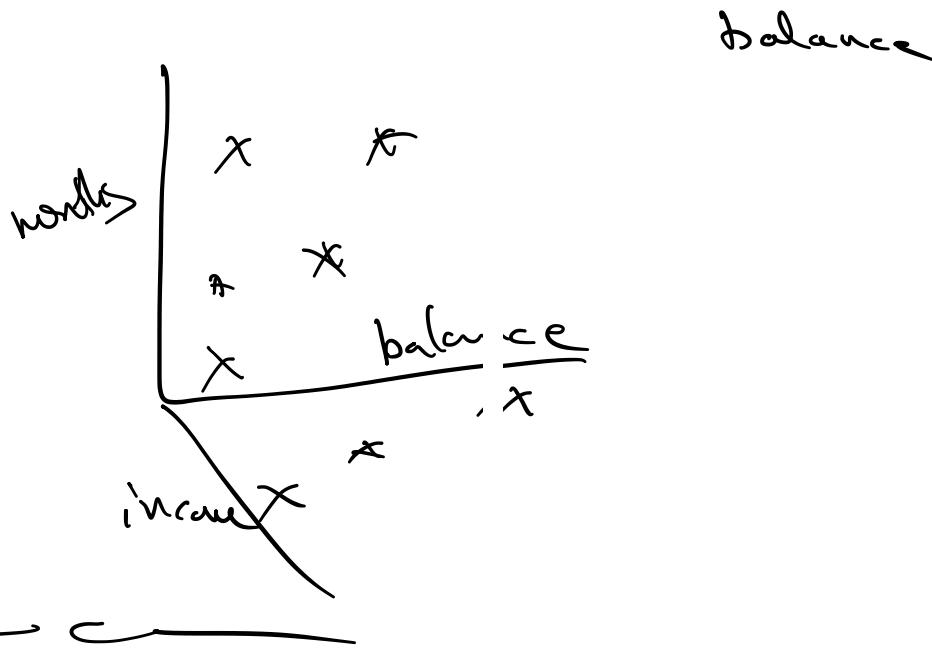
Represent entities as points
in Euclidean Space

Entities are represented

as $\langle x_1, x_2, \dots, x_p, y \rangle$

$\in \mathbb{R}^{p+1}$





$D \rightarrow$ distances

- Euclidean Distance

$$d(x_1, x_2) = \sqrt{\sum_{i=1}^p (x_{1i} - x_{2i})^2}$$

v'

Inductive Bias

→ Examples that are "close"
have similar outcomes

$$P(Y=y | X=x)$$

— — —

↑↑

points that
are "close"

— . . . —

n .. 1

(1) Categorical attributes
are transformed into
numerical

(2) Scaling & centering important

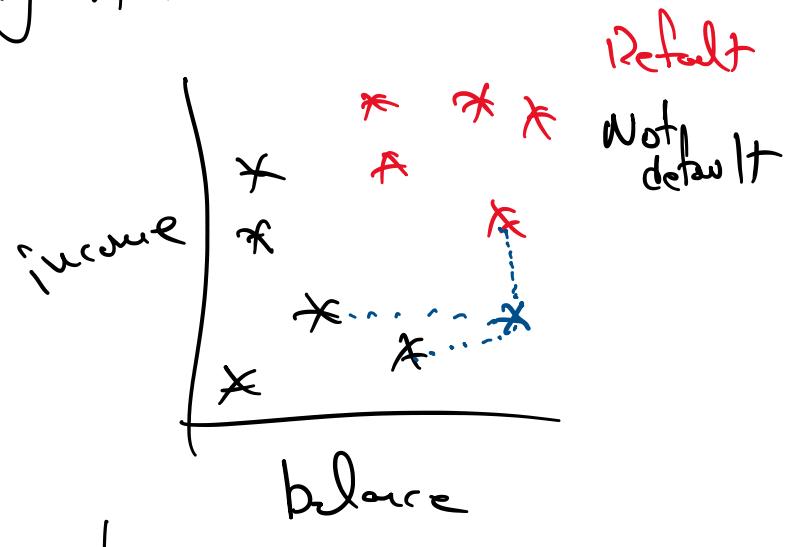
$$\overline{C} \leftarrow \underline{C}$$

K-nearest neighbor algorithm

Given training set

$$\{ \langle x_{11}, \dots, x_{1P}, g_1 \rangle, \\ \langle x_{21}, \dots, x_{2P}, g_2 \rangle, \\ \vdots \\ \vdots \}$$

$$\{ \langle x_{n1}, \dots, x_{nP}, g_n \rangle \}$$



Generate predictions for new

examples $\langle x_{t1}, \dots, x_{tp} \rangle = x_t$.

→ Compute distance
between x_t &
each point in
training set

→ Find K nearest
points

→ Compute majority
class

→ Return majority class

Why choose K ?

(Case 1) $K=n$

All predictions
are the same
(under fitting)

(Case 2) $K=1$

"Unstable"
(over fitting)

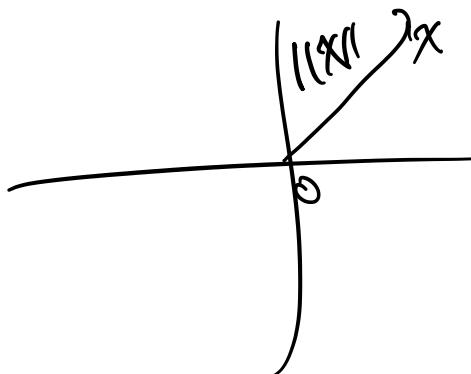
$\cup \cup \dots$

K : hyper-parameter

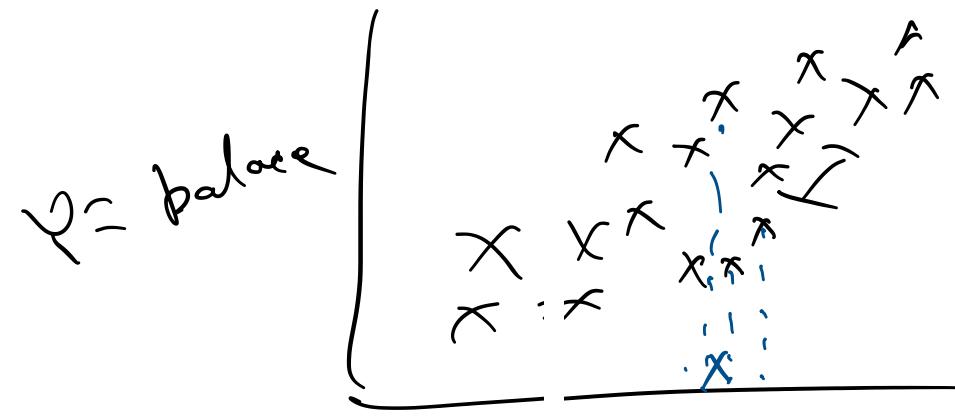
interested over vs. underfitting

generalization: ability to accurately
predict outcomes for unseen data

$$\|x\| = \sqrt{\sum_{j=1}^n (x_j - \bar{o})^2}$$



$\curvearrowleft \curvearrowright$

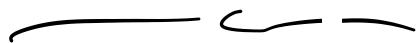


Now I have

x_i - income

$y_1, y_2, y_3, \dots, y_k$ for k nearest neighbors

\bar{y} minimizes $\min \sum_{i=1}^k (y_i - \mu)^2$



Assume

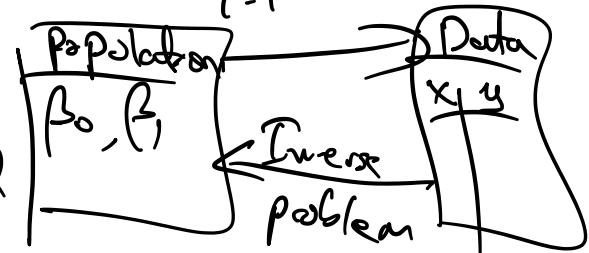
$$Y = \beta_0 + \beta_1 X$$

Linear regression model

Estimation problem:

- Given training data
- Learn parameters β_0, β_1

$$y_i = \beta_0 + \beta_1 x_i + \epsilon^{\text{variation}}$$



$$\min \sum (y_i - (\beta_0 + \beta_1 x_i))^2$$

$$f(\beta_i; \tilde{\beta}_i) = \left(\sum_i \frac{y_i - \hat{y}_i}{\hat{y}_i} \right)^p$$

Model

→ Closed-form solution for $p=1$

→ Algorithms to minimize this function