



האוניברסיטה
העברית
בירושלים
THE HEBREW
UNIVERSITY
OF JERUSALEM

Artificial Intelligence in Medicine

Regression

Nir Friedman and Tommy Kaplan

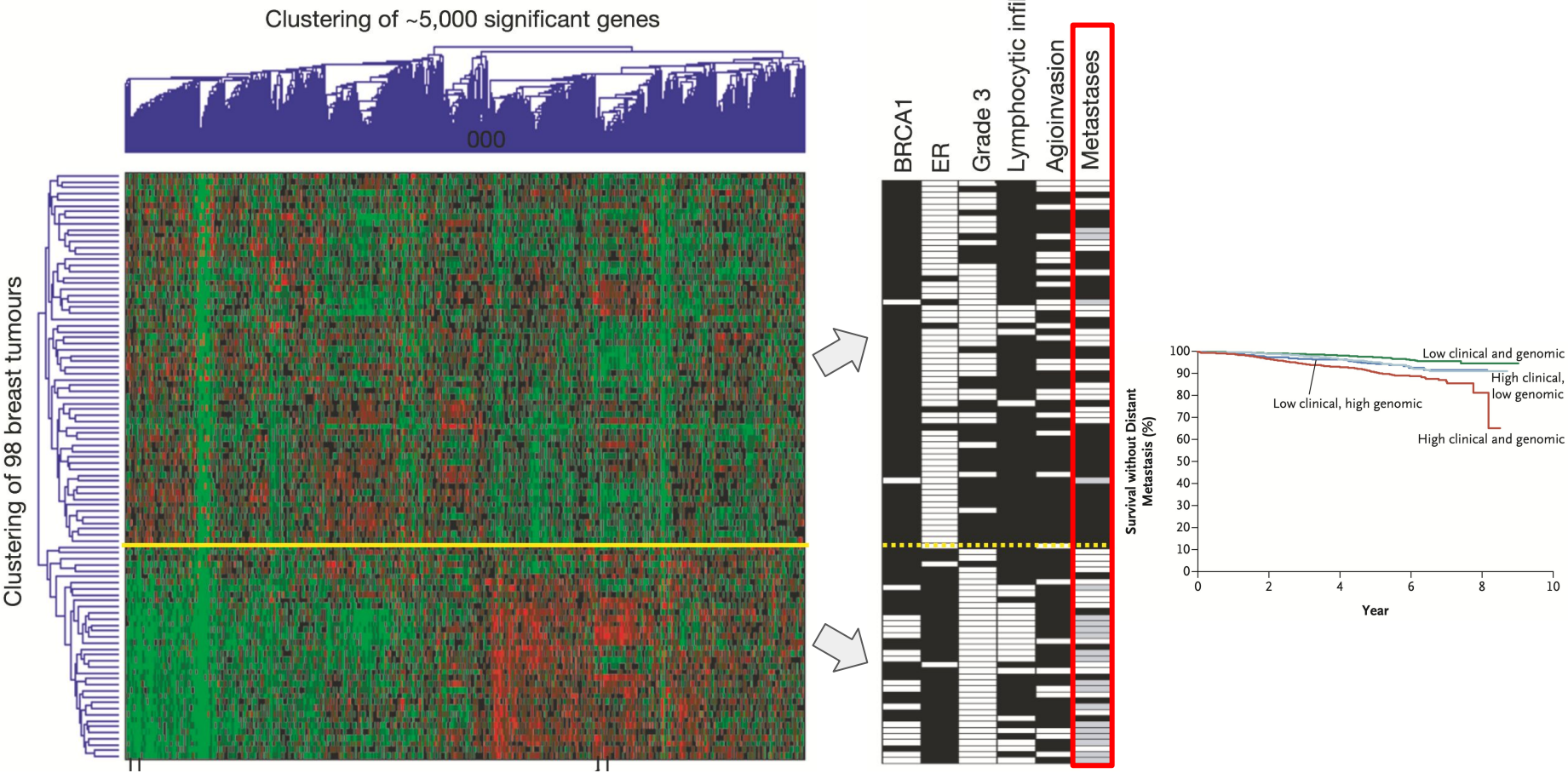
6/2/24

תוכן

- ניבוי רציף (כמותי) לעומת בינארי
- במקום קו, רוצים לנבא ערך, למשל תוחלת הערך הצפוי כפונקציה של פיצ'רים (ערכים אחרים)
- שערך טעות (בינארית) <- טעות ריבועית - בגלל שמניחים רעש עם התפלגות נורמלית)
- רגרסיה לינארית רגילה
- רגרסיה ריבועית, לוג, סינוס
- רגרסיה לינארית מרובה
- מורכבות ואובר-פיטינג
- רגולריזציה - רידג' $L1$, לאסו $L2$, אלסטיק-נט $L1+L2$
- עצי רגרסיה, יער
- מודל מבוסס קוהורט (שכנים)

Previously, we predicted metastases presence

what if we wanted to predict survival? dosage?



Regression models in medicine

- Pediatric dosage calculations (body weight)
- Chemotherapy (body surface area)

DuBois and DuBois¹ : 1.85 m²

$$\text{Equation: } \text{BSA (m}^2\text{)} = 0.007184 \times \text{Height(cm)}^{0.725} \times \text{Weight(kg)}^{0.425}$$

Gehan and George²: 1.85 m²

$$\text{Equation: } \text{BSA (m}^2\text{)} = 0.0235 \times \text{Height(cm)}^{0.42246} \times \text{Weight(kg)}^{0.51456}$$

Haycock³: 1.85 m²

$$\text{Equation: } \text{BSA (m}^2\text{)} = 0.024265 \times \text{Height(cm)}^{0.3964} \times \text{Weight(kg)}^{0.5378}$$

Mosteller⁴: 1.84 m²

$$\text{Equation: } \text{BSA (m}^2\text{)} = \text{SQRT} ([\text{Height(cm)} \times \text{Weight(kg)}] / 3600)$$

1) DuBois D, DuBois DF. A formula to estimate the approximate surface area if height and weight be known. Arch Int Med 1916;17:863-71.

2) Gehan EA, George SL. Estimation of human body surface area from height and weight. Cancer Chemother Rep 1970;54:225-35.

3) Haycock GB, Schwartz GJ, Wisotsky DH. Geometric method for measuring body surface area: A height-weight formula validated in infants, children and adults. J Pediatr 1978;93:62-6.

4) Mosteller RD. Simplified calculation of body-surface area. N Engl J Med 1987;317:1098.

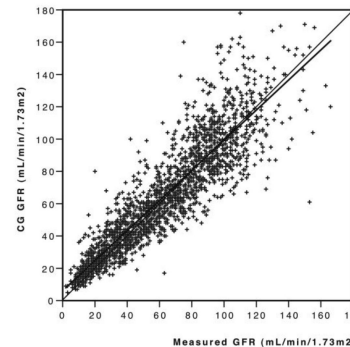
- GFR (Creatine, Creatinine, and kidney function)
(estimated Glomerular Filtration Rate)

Prediction of Creatinine Clearance from Serum Creatinine¹

DONALD W. COCKCROFT and M. HENRY GAULT

Departments of Medicine, Queen Mary Veterans' Hospital, Montreal, Quebec, and Memorial University, St. John's, Newfoundland

$$C_{Cr} = \frac{(140 - \text{age}) (\text{wt kg})}{72 \times S_{Cr}(\text{mg}/100 \text{ ml})}$$



החישוב נעשה כדלהלן:

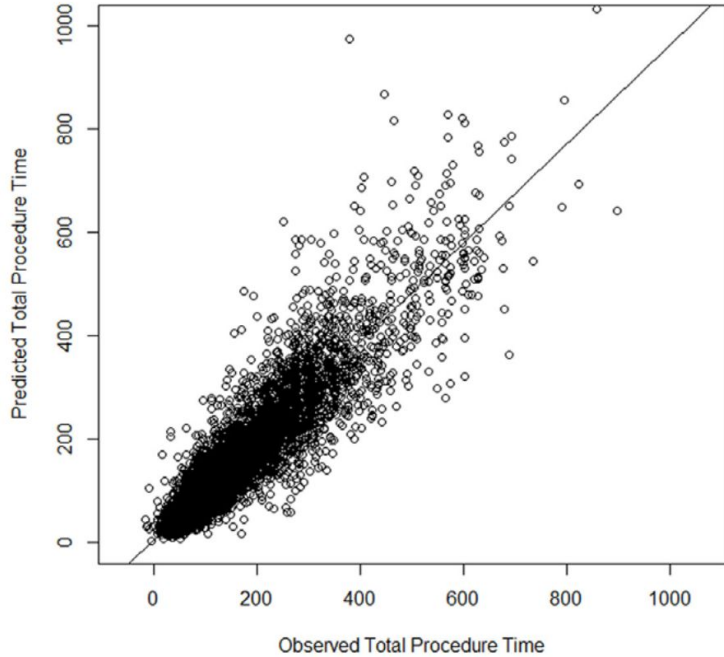
$$\text{Val} = \frac{[(140 - \text{גיל}) \times \text{משקל}] / (\text{serum creatinine} \times 72)}{\text{אם הנבדק היא נבדקת יש להכפיל את הערך ב-0.85}}$$

משמעות התוצאות

משמעות	ערכ
תפקוד הכליות תקין	מעל 50
אין צורך להתחשב בו למתן תרופות	
ליקוי בינוני בתפקוד הכליות	בין 10 ל 50
יש להתחשב בו במתן תרופות מסוימות	
ליקוי חמור בתפקוד הכליות	מתחת ל 10

Improving the Prediction of Total Surgical Procedure Time Using Linear Regression Modeling

Eric R. Edelman^{1*}, Sander M. J. van Kuijk², Ankie E. W. Hamaekers³,
Marcel J. M. de Korte³, Godefridus G. van Merode⁴ and Wolfgang F. F. A. Buhre³



Multiple regression model to analyze the total LOS for patients undergoing laparoscopic appendectomy

Teresa Angela Trunfio¹, Arianna Scala^{2*}, Cristiana Giglio³, Giovanni Rossi⁴, Anna Borrelli⁴, Maria Romano⁵ and Giovanni Improta^{2,6}

Noise is all around us

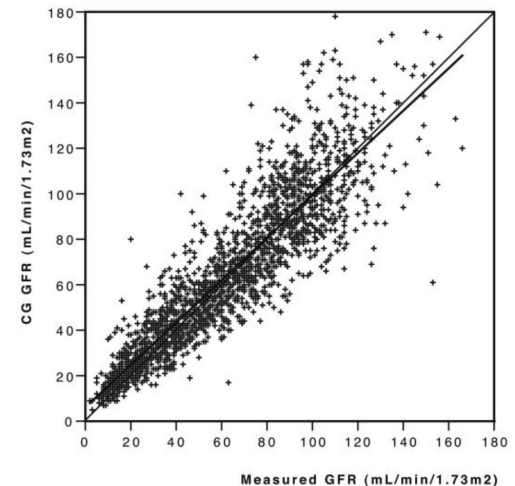
Sum of multiple effects \Rightarrow Central limit theorem

המשפט החזק של המספרים החלשים

In the absence of a better model or knowledge

Non-additive noise \Rightarrow multiplicative cases

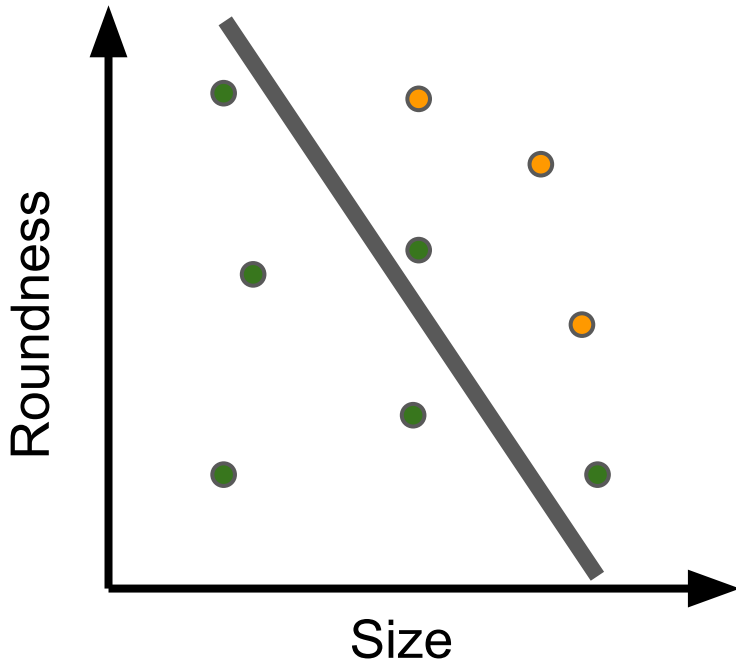
Creatinine = noise depends on body mass



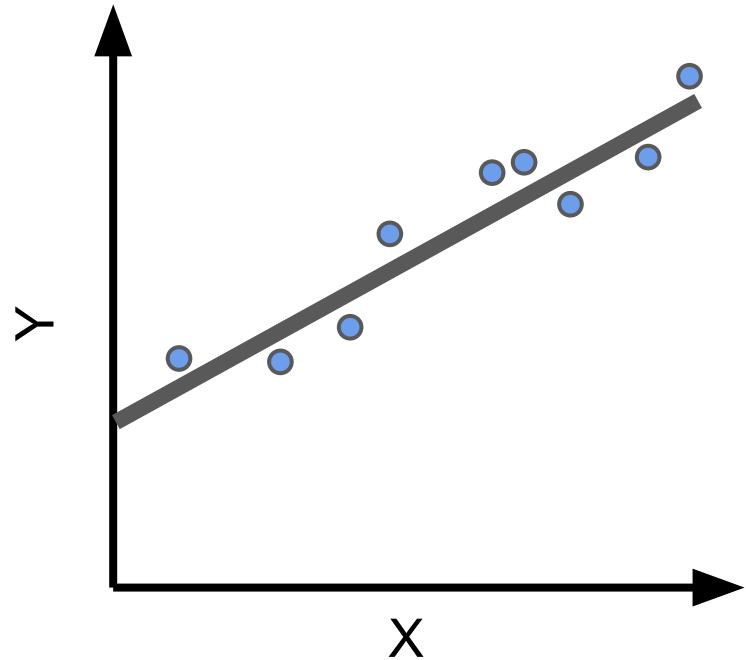
Error estimation

“No more counting dollars, we'll be counting stars”

Minimize classification errors

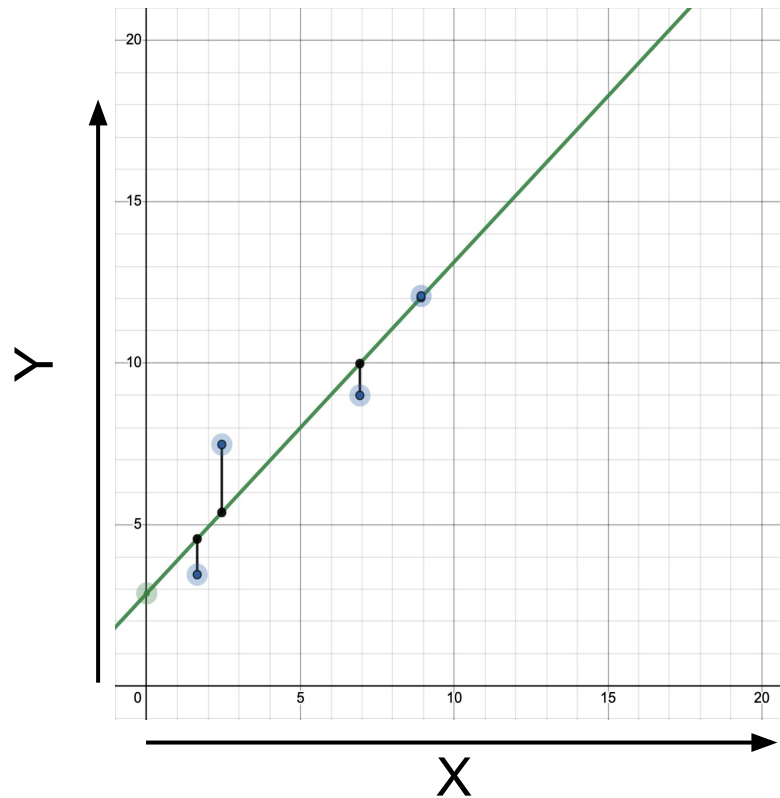


Accurately predict Y from X



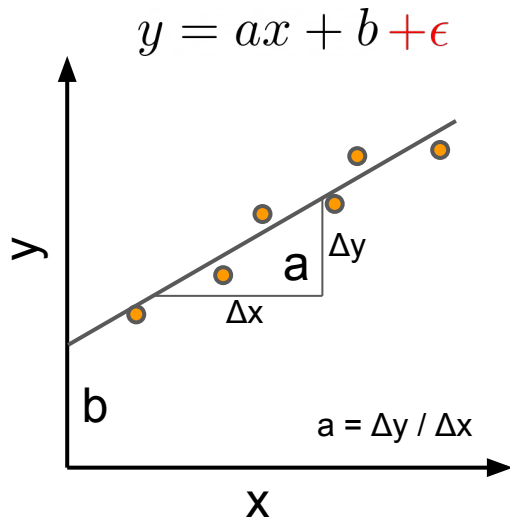
Error estimation

Absolute distances? Squared distances? ...



Linear regression (univariate)

Least squares method



$$\arg \min_{a,b} \sum (y_i - (ax_i + b))^2 =$$

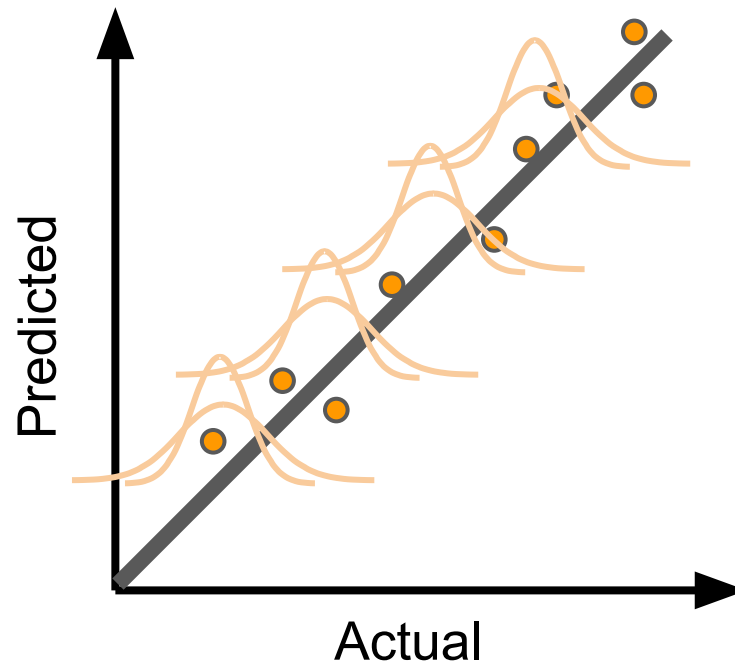
$$\arg \min_{a,b} \sum (y_i - \hat{y}_i)^2 =$$

$$\arg \min_{a,b} \sum \epsilon_i^2$$

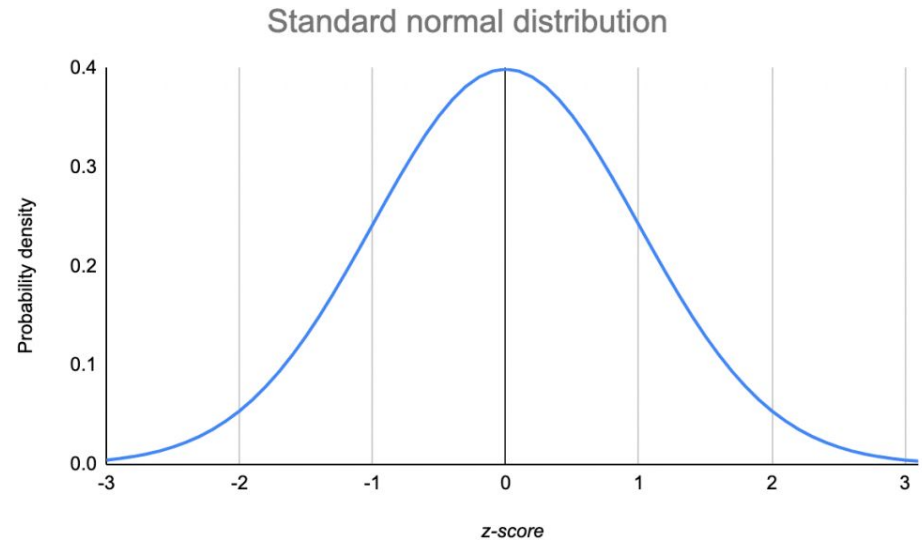
Why squared errors?

Assuming normally distributed noise

$$f(\epsilon) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\epsilon}{\sigma}\right)^2}$$



Normal Distribution



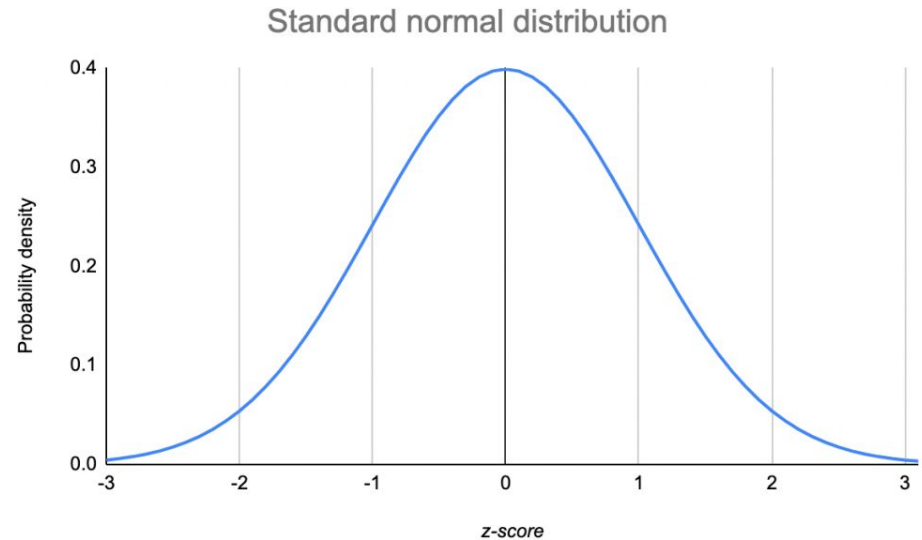
$$f(\epsilon) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\epsilon}{\sigma}\right)^2}$$

Function of the square distance from the mean

Normal Distribution

Properties:

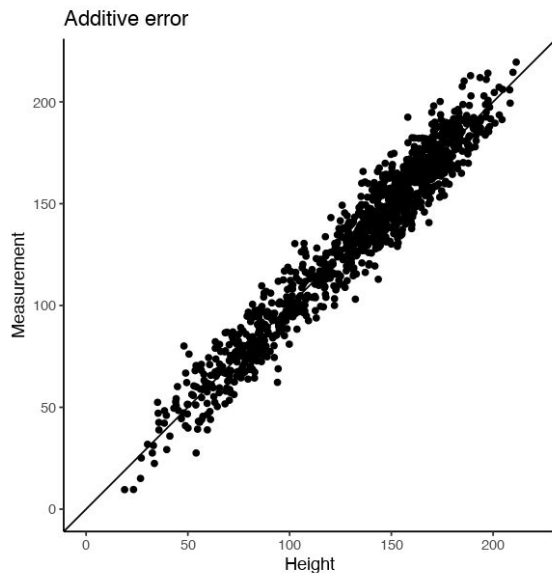
- Symmetrical
- Central limit theorem:
sum of small effects results in a normal distribution
- Formally, the simplest distribution with a given mean+variance
- Mathematical convenience
- Maintained under addition, shifting and rescaling



Additive vs Multiplicative Noise

Additive measurement noise:

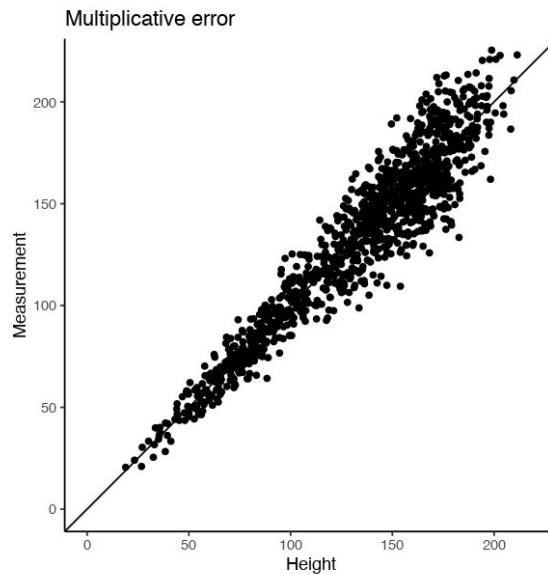
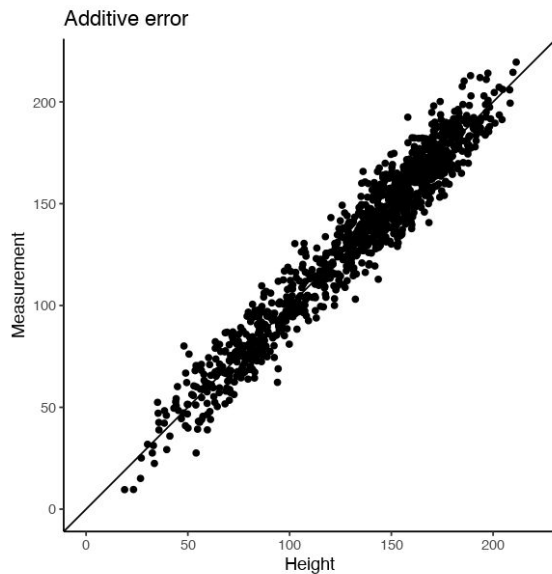
- “Height +/- 10cm”
- Error size does not depend on the actual value



Additive vs Multiplicative Noise

Multiplicative measurement noise:

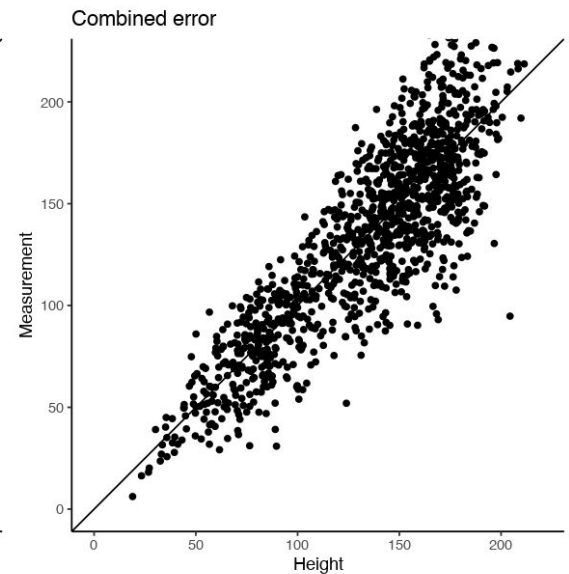
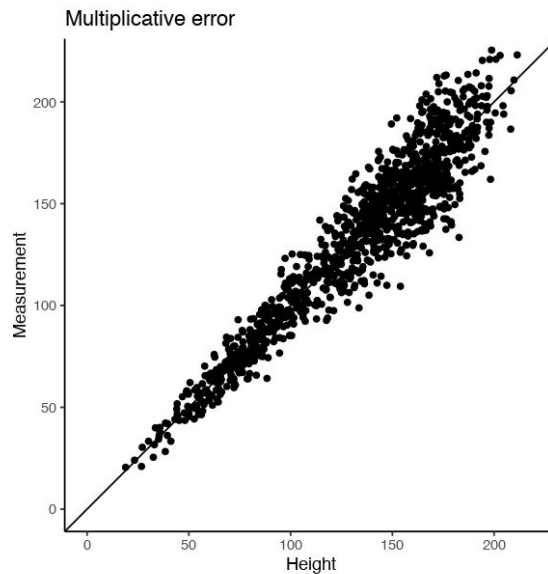
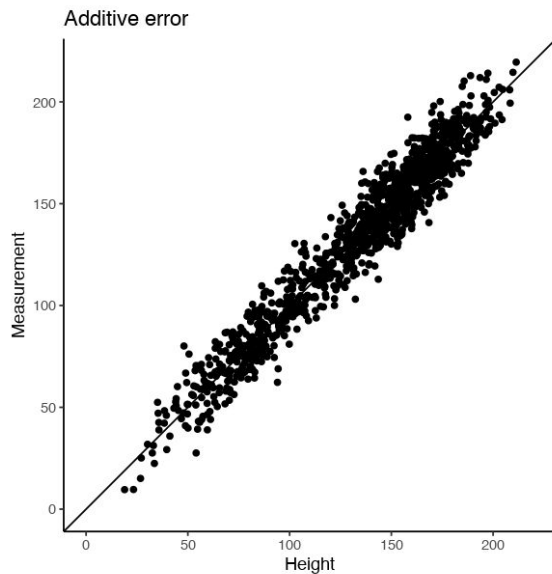
- “Height +/- 10%”
- Error size increase with actual value



Additive vs Multiplicative Noise

Combined noise

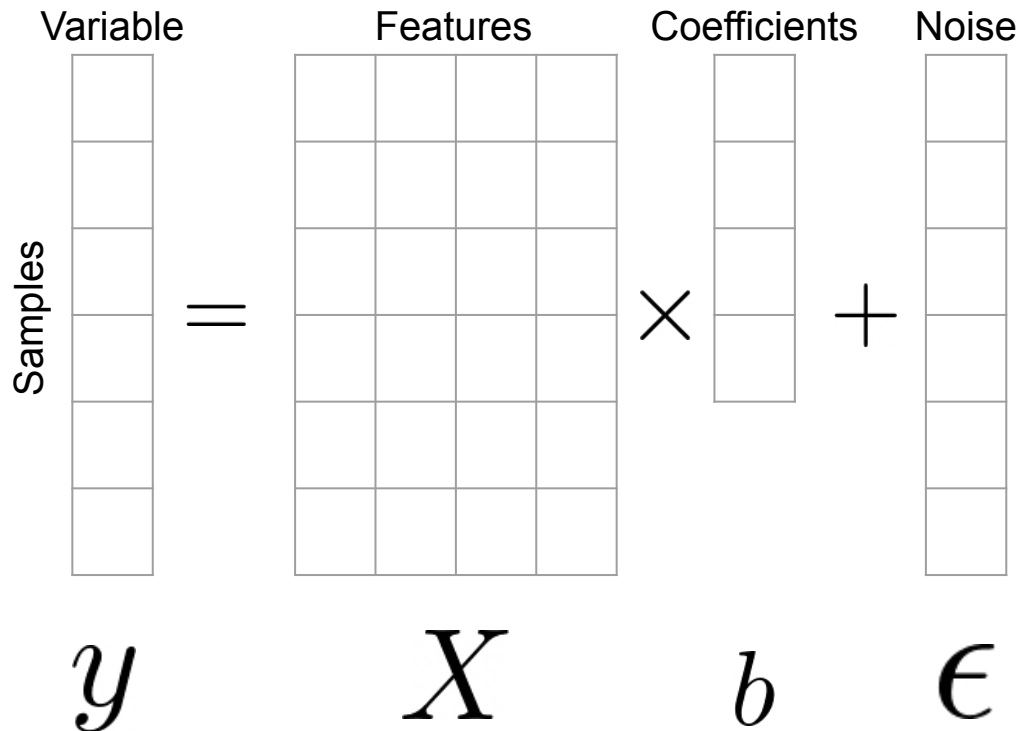
- Additive + multiplicative components



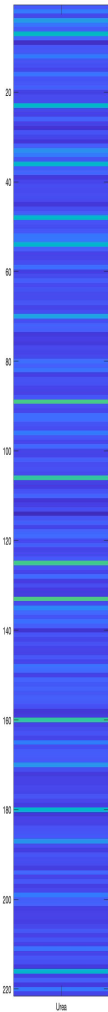
Linear regression (multivariate)

Definition, formalities, High-dimensional interpretation

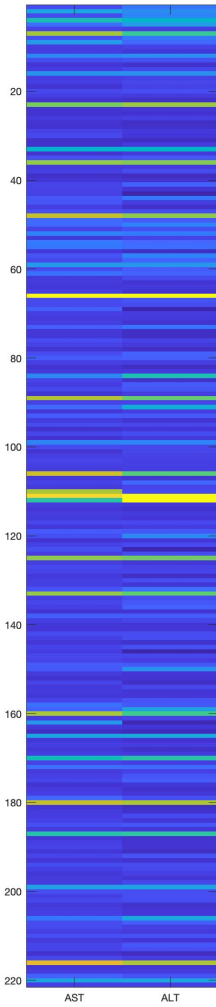
$$y = b_1 x_1 + \dots + b_n x_n + b_0 + \epsilon$$



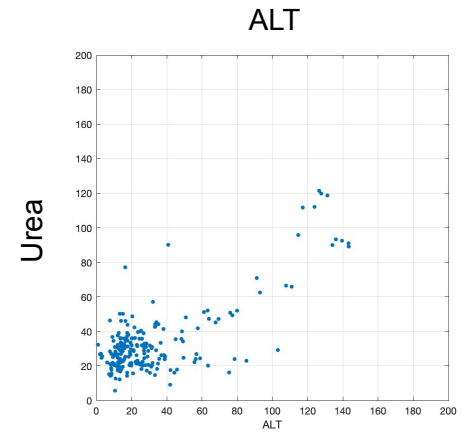
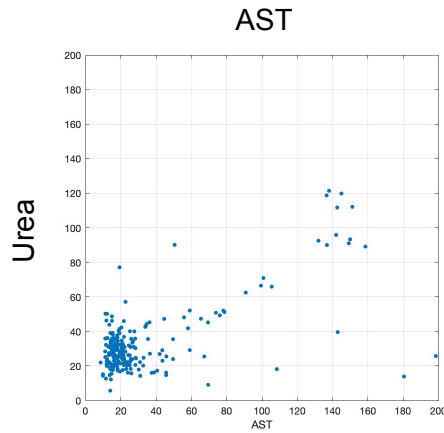
Linear regression (multivariate)



Urea



AST, ALT

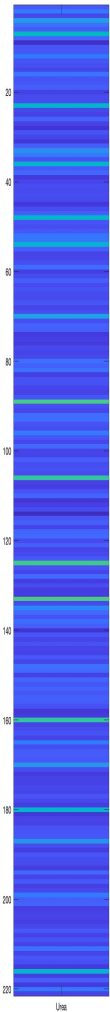


A Machine Learning Approach to Predict Creatine Kinase Test Results

Zehra Nur Canbolat ^a, Gökhan Silahtaroğlu ^{a,b}, Özge Doğuç ^{a*}, Nevin Yılmaztürk ^b

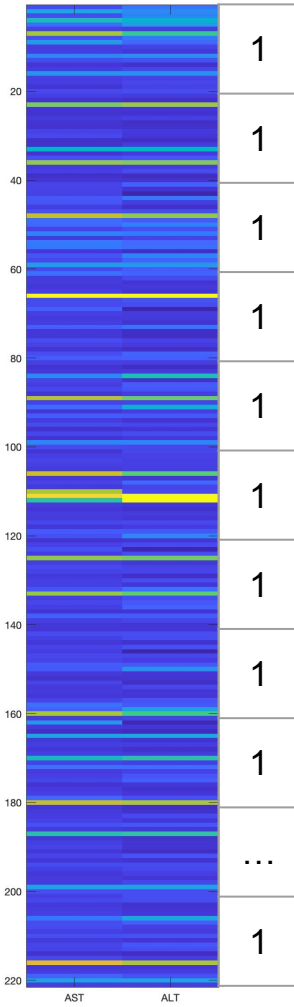
Age, Sex, AST, ALT, Urea, Glucose, etc. from 222 patients, hospitalized in Istanbul, Turkey, 2017-2019

Linear regression (multivariate)



Urea

==

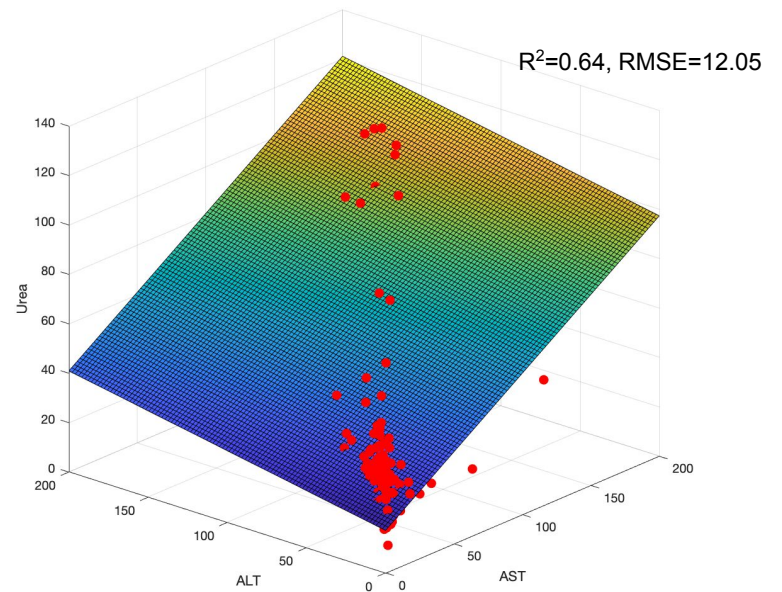
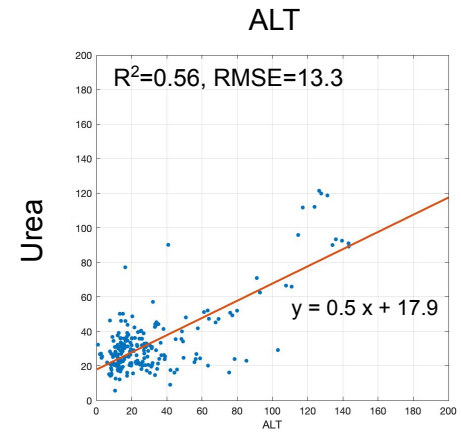
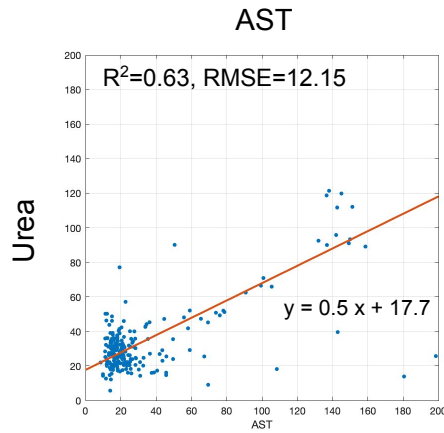


AST, ALT, 1

×

.40
.12
17

coeffs



Linear regression (multivariate)

- Regularization: antagonizing “weak” features
- Feature selection: choosing the variables that matter
- Reducing complexity \Rightarrow more generalization

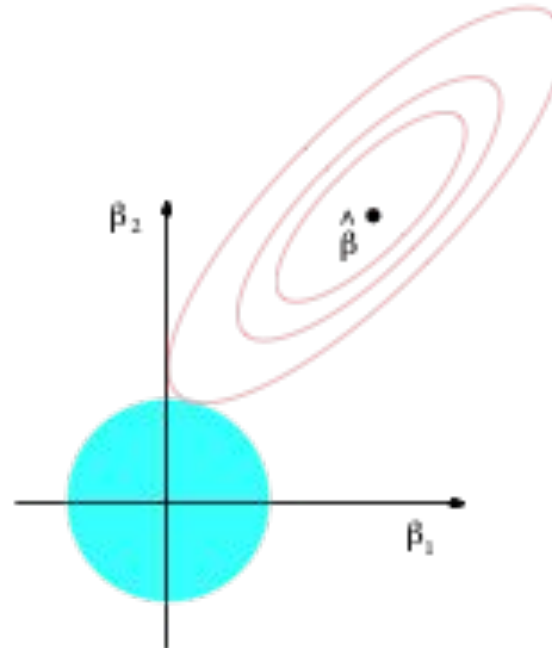
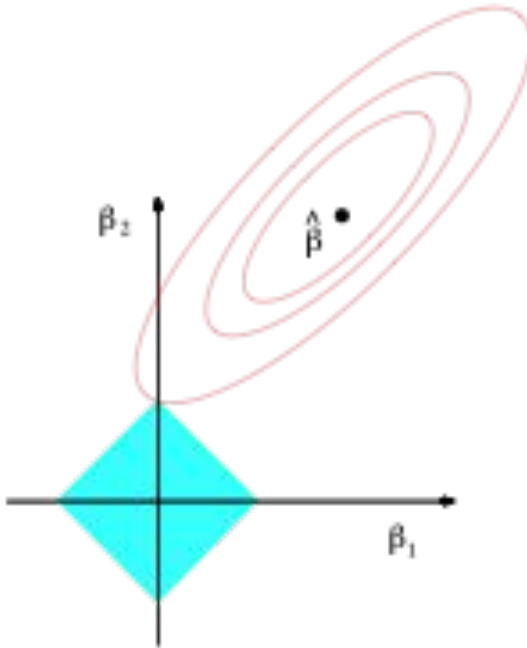
$$y = b_1x_1 + \cancel{b_2x_2} + \cancel{b_3x_3} + \dots + b_nx_n + b_0 + \epsilon$$

Ignoring x_2, x_3 is like setting $b_2=0, b_3=0$

Linear regression (multivariate)

$$\arg \min_{\beta} \left(\underbrace{\|y - X\beta\|^2}_{\text{Error}} + \underbrace{\lambda_1 \sum |\beta_j|}_{\text{Lasso}} + \underbrace{\lambda_2 \sum |\beta_j|^2}_{\text{Ridge}} \right)$$

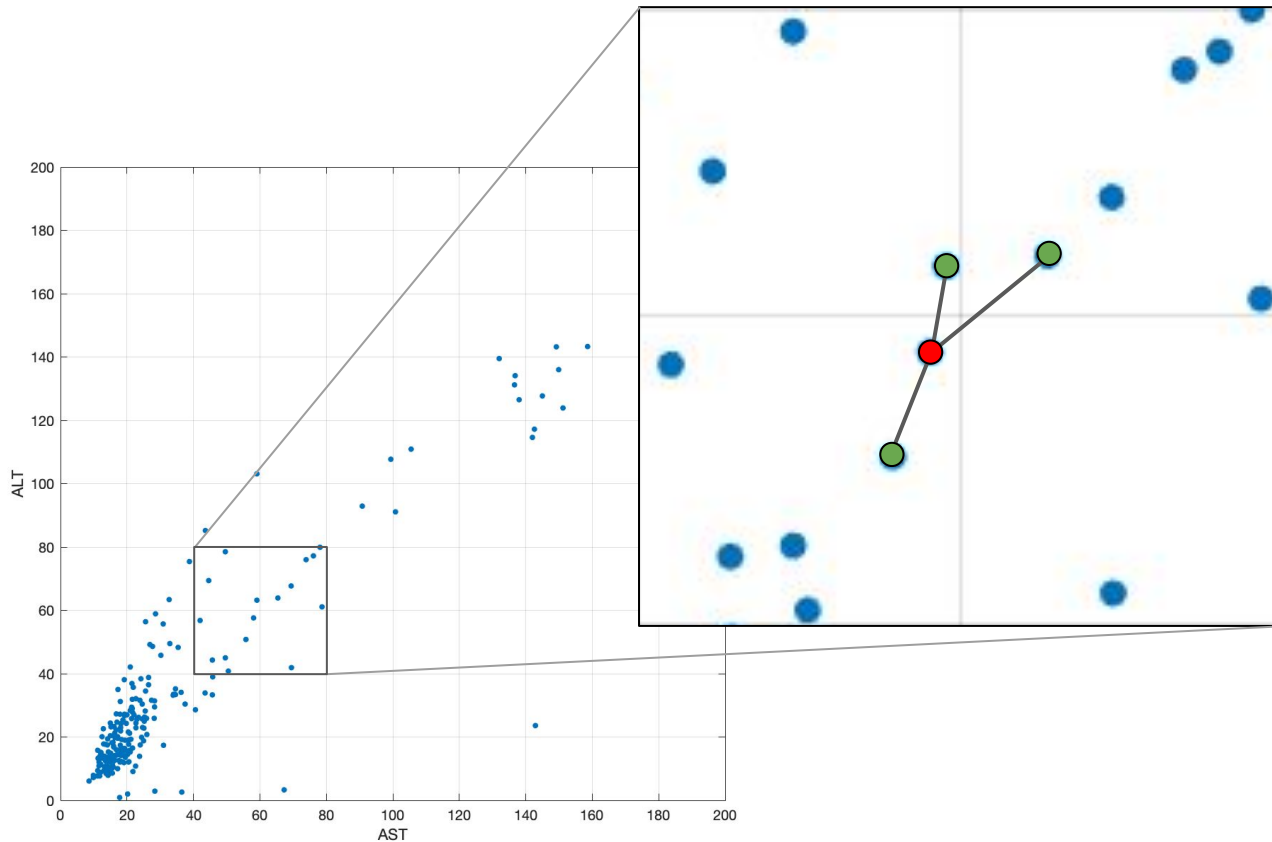
← Elastic net →



Cohort-based models

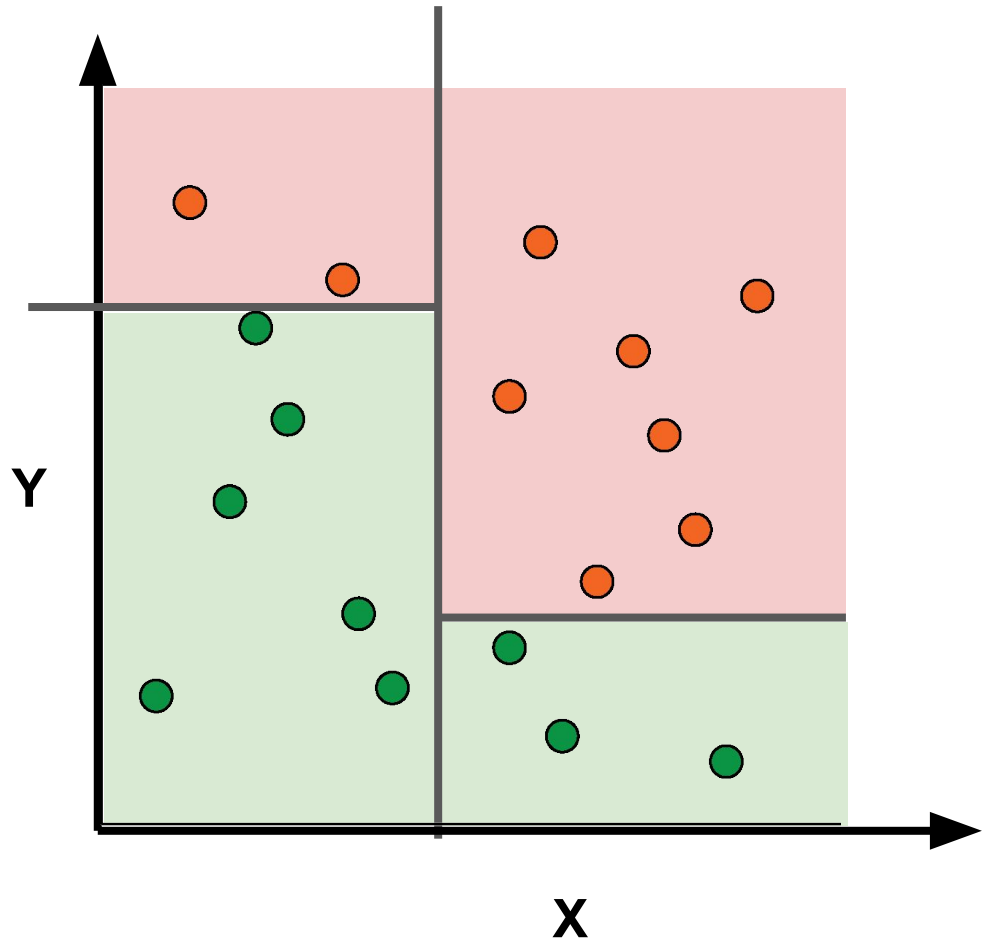
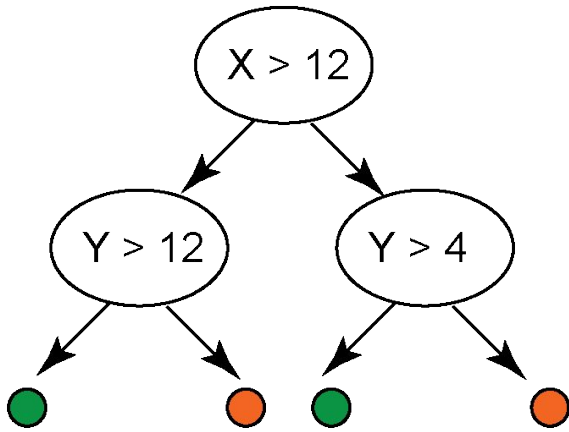
K-nearest neighbors regression models

- Distances calculation
- Weight calculation



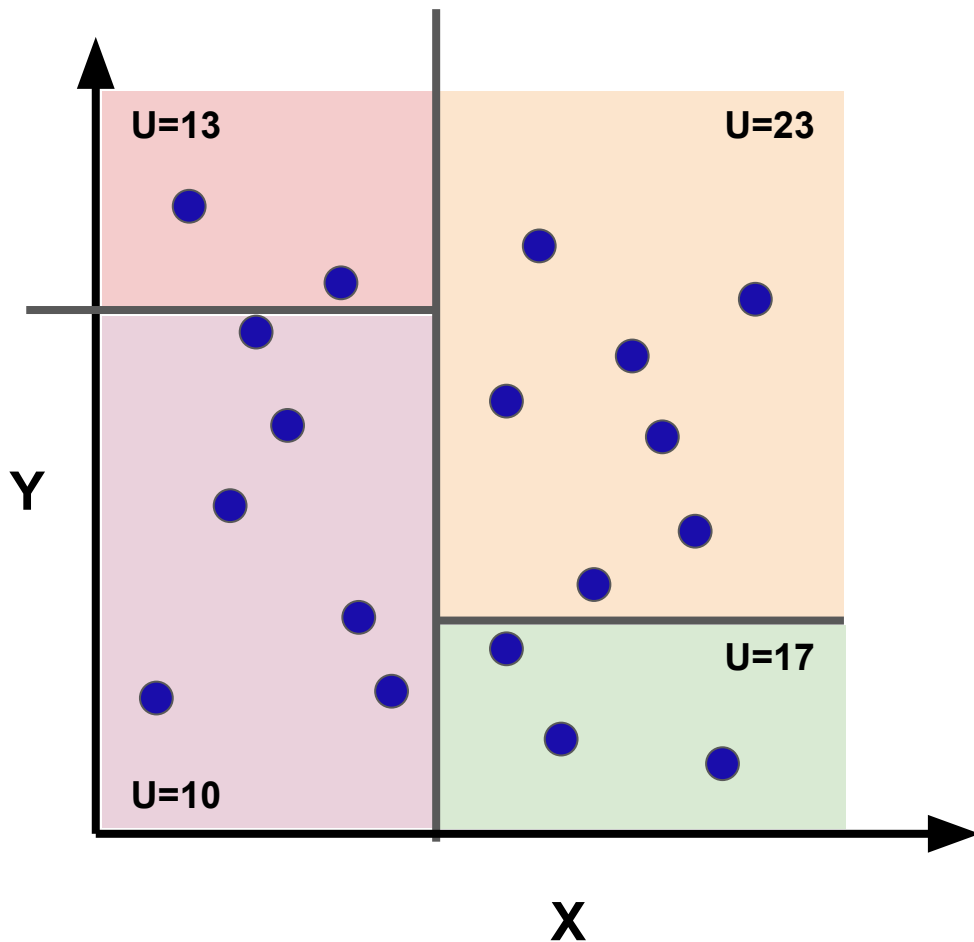
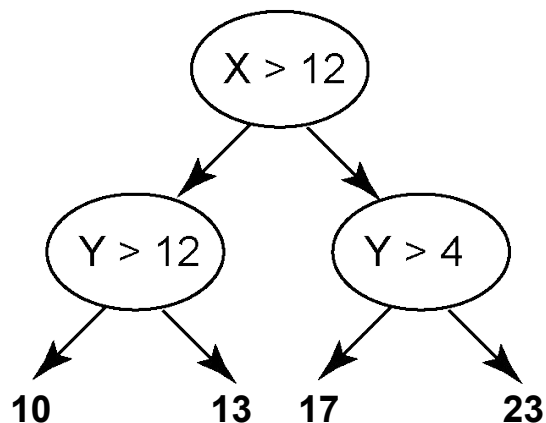
$$y_i = \frac{1}{K} \sum_{j \in N_i} y_j$$

Regression trees

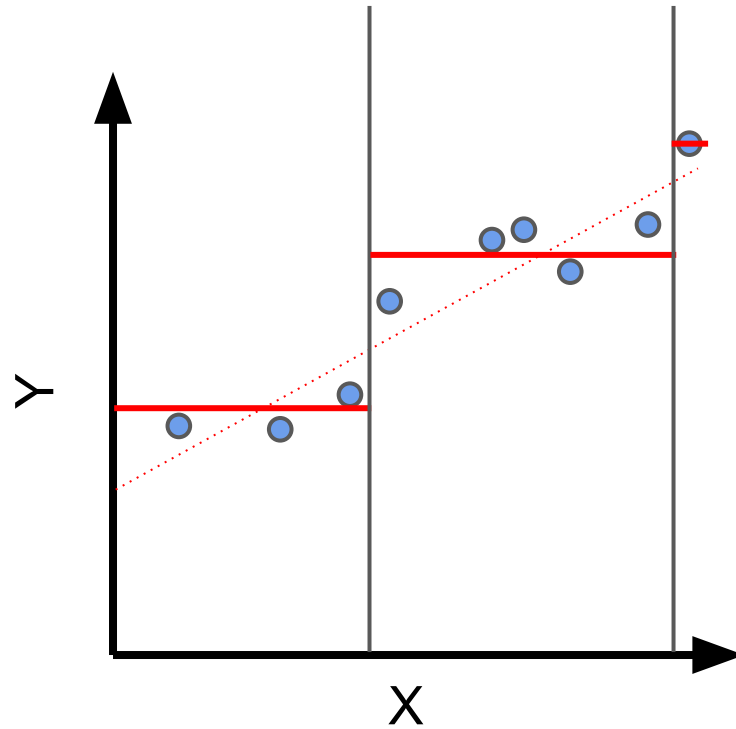


Regression trees

Predict a constant value at each leaf

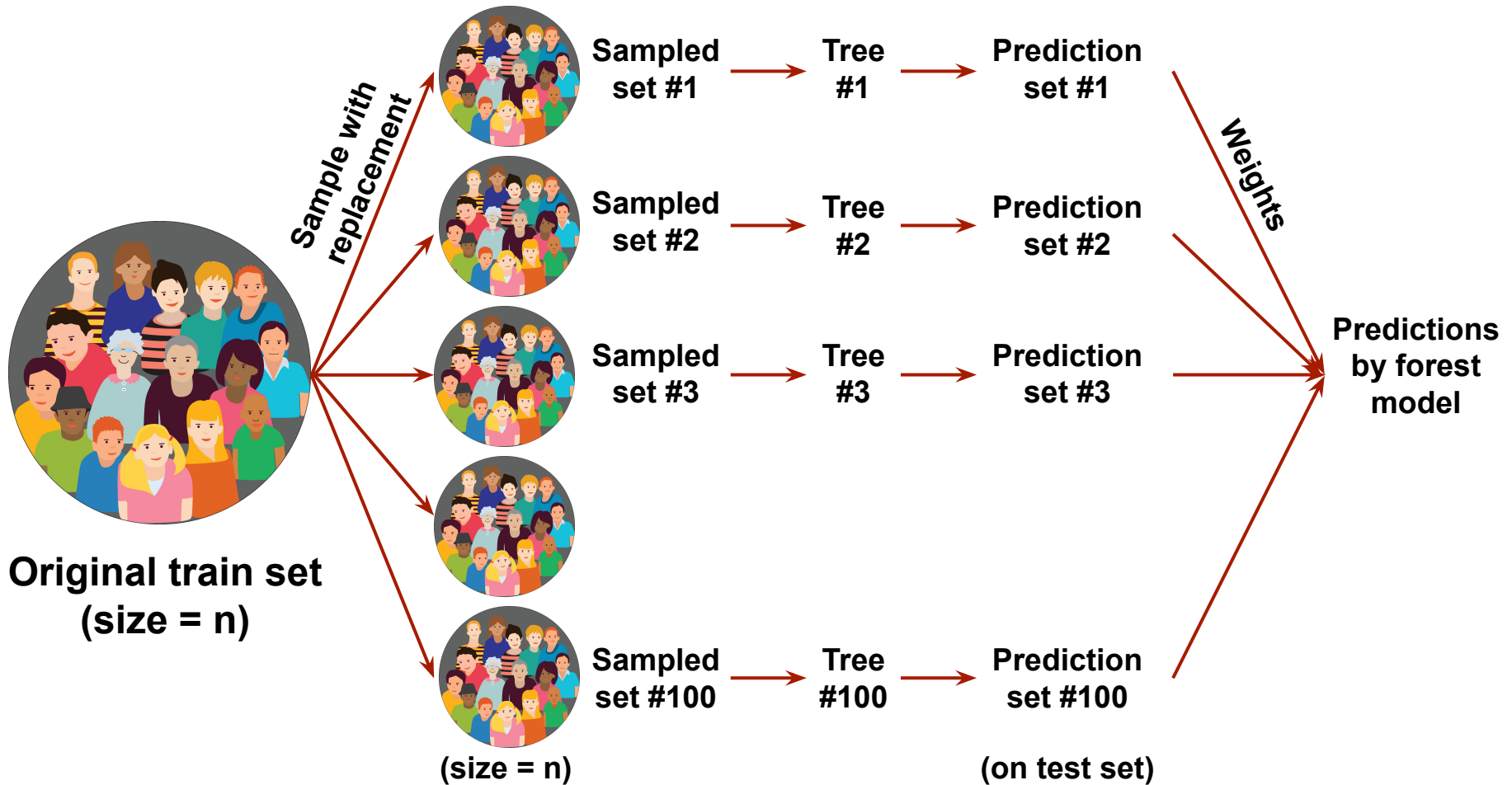


Regression tree (univariate)

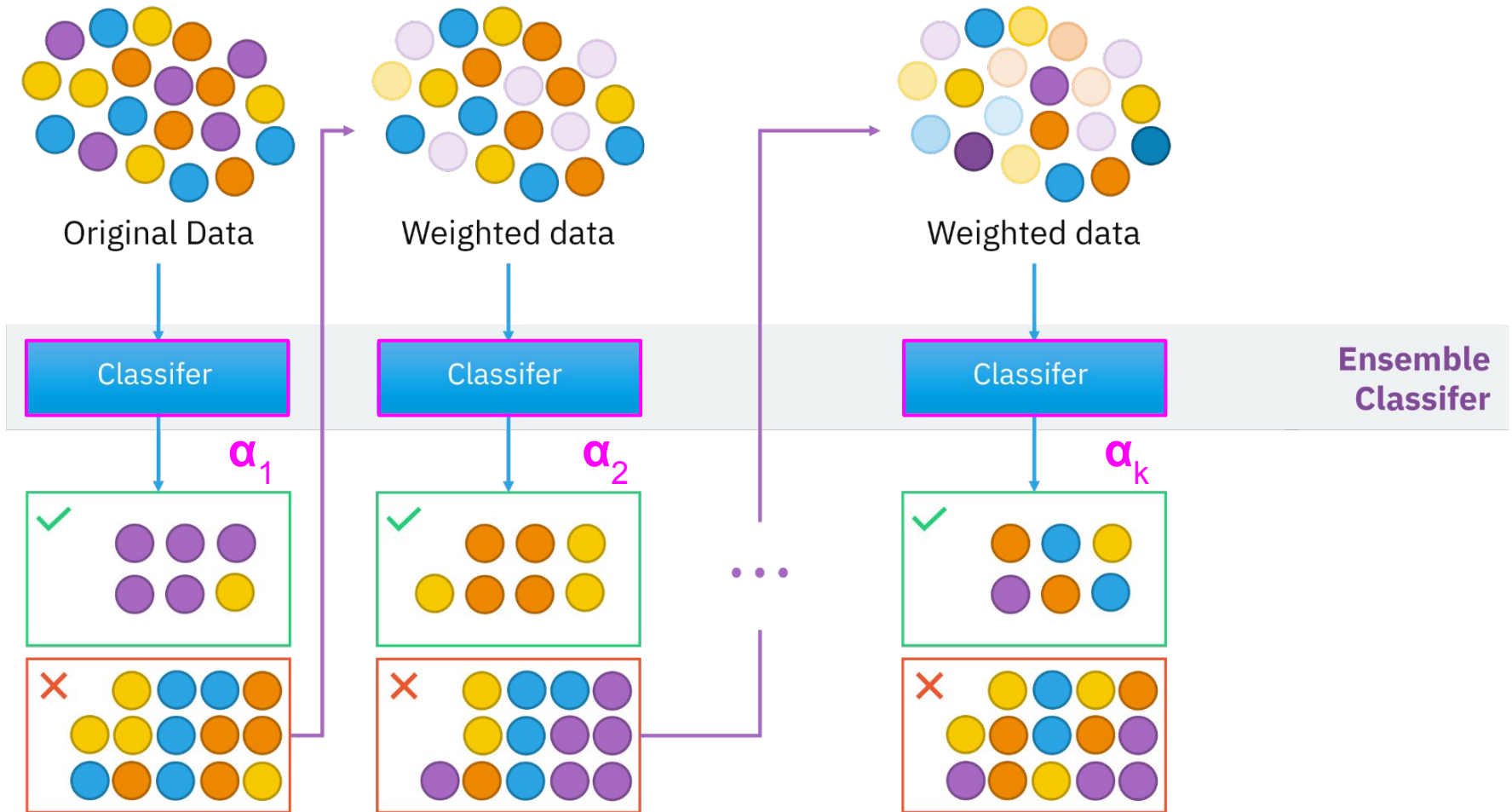


Ensemble models - forests

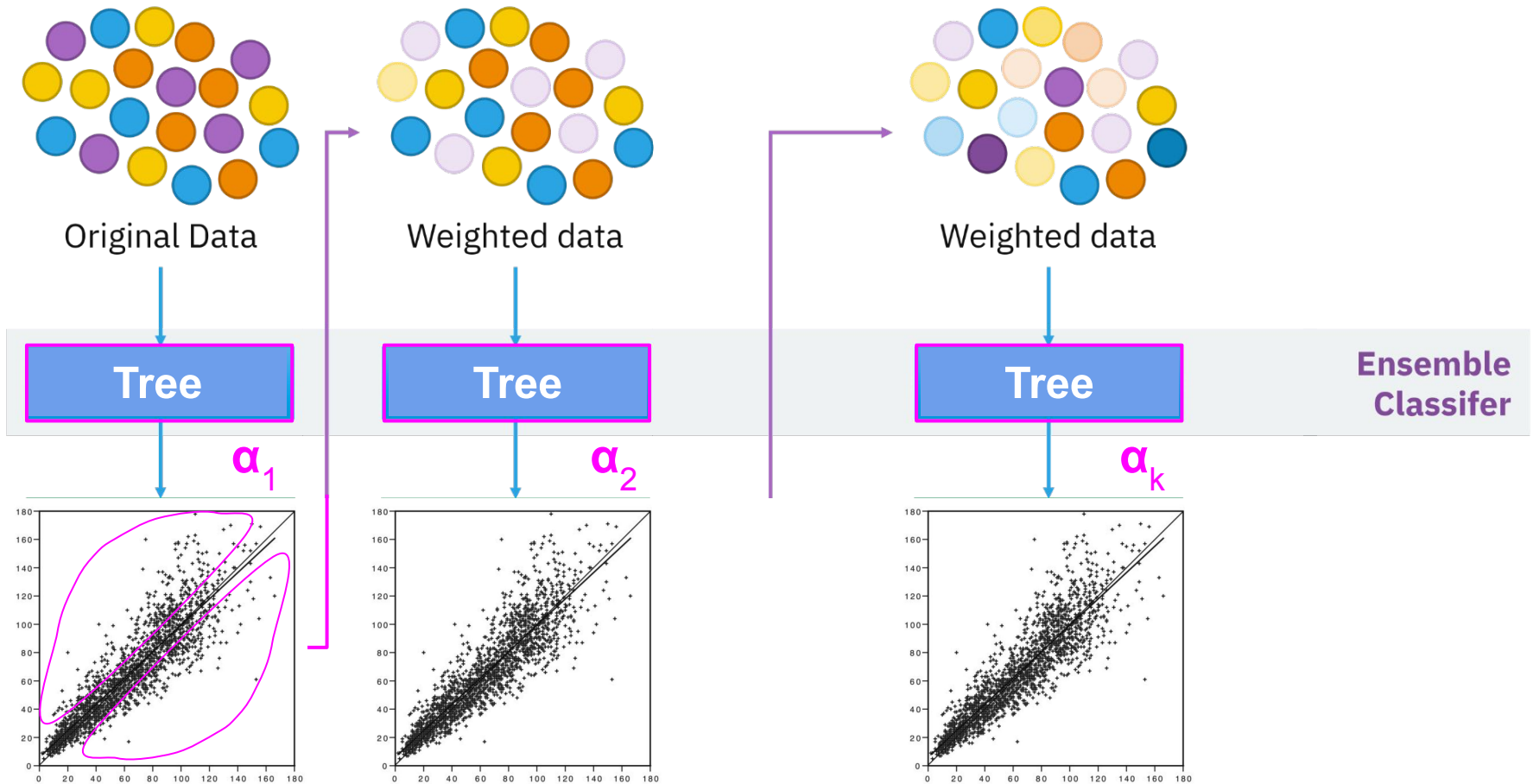
- Bootstrapping: sampling with replacement



Boosting (ensemble classifier)



Gradient boosting trees (XGBoost)



Non-linear regression

Log-transformed?

Squared?

Sigmoid

$BMI \sim \text{Weight}^2 / \text{Height}$

$\log(BMI) \sim 2 * \log(\text{Weight}) - \log(\text{Height})$

Syllabus

1. Introduction
2. Classification
3. Learning 1
4. AI in ophthalmology (Prof. Itay Chowers)
5. Learning 2
6. Regression
7. Clustering
8. Visualization (and dimensionality reduction)
9. Deep learning in image analysis (Prof. Leo Joskowicz)
10. Missing data, statistical dependencies
11. Natural language in medicine (Dr. Gabi Stanovsky)
12. Decisions (utility)
13. Longitudinal Data / Project