



Machine Learning Basic Algorithms

Unsupervised Algorithm: K-means ...

Supervised Algorithm: Linear Regression/Regression

Supervised Algorithm: Decision tree ...

Supervised Algorithm: KNN ...

Meta Algorithm: Ensemble ...

K-MEANS?!



Y NOT A-MEANS

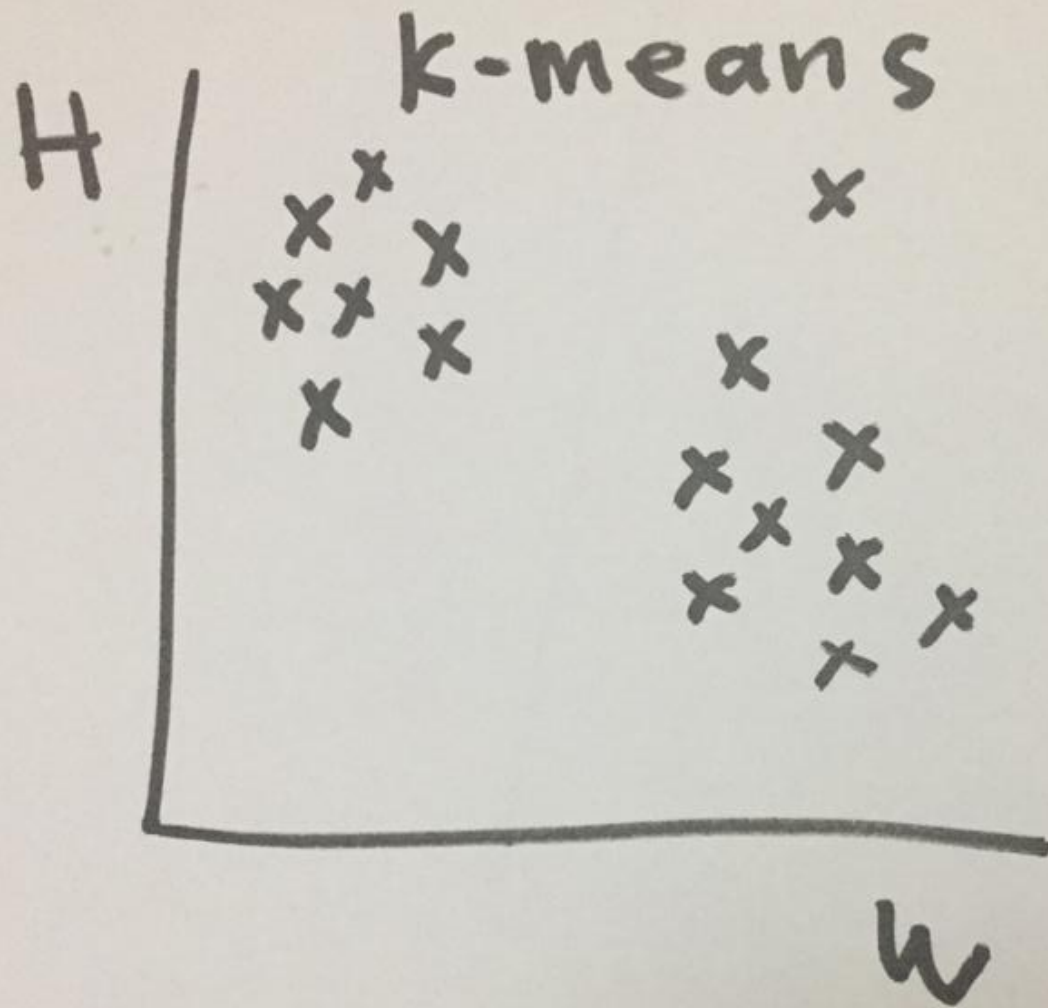


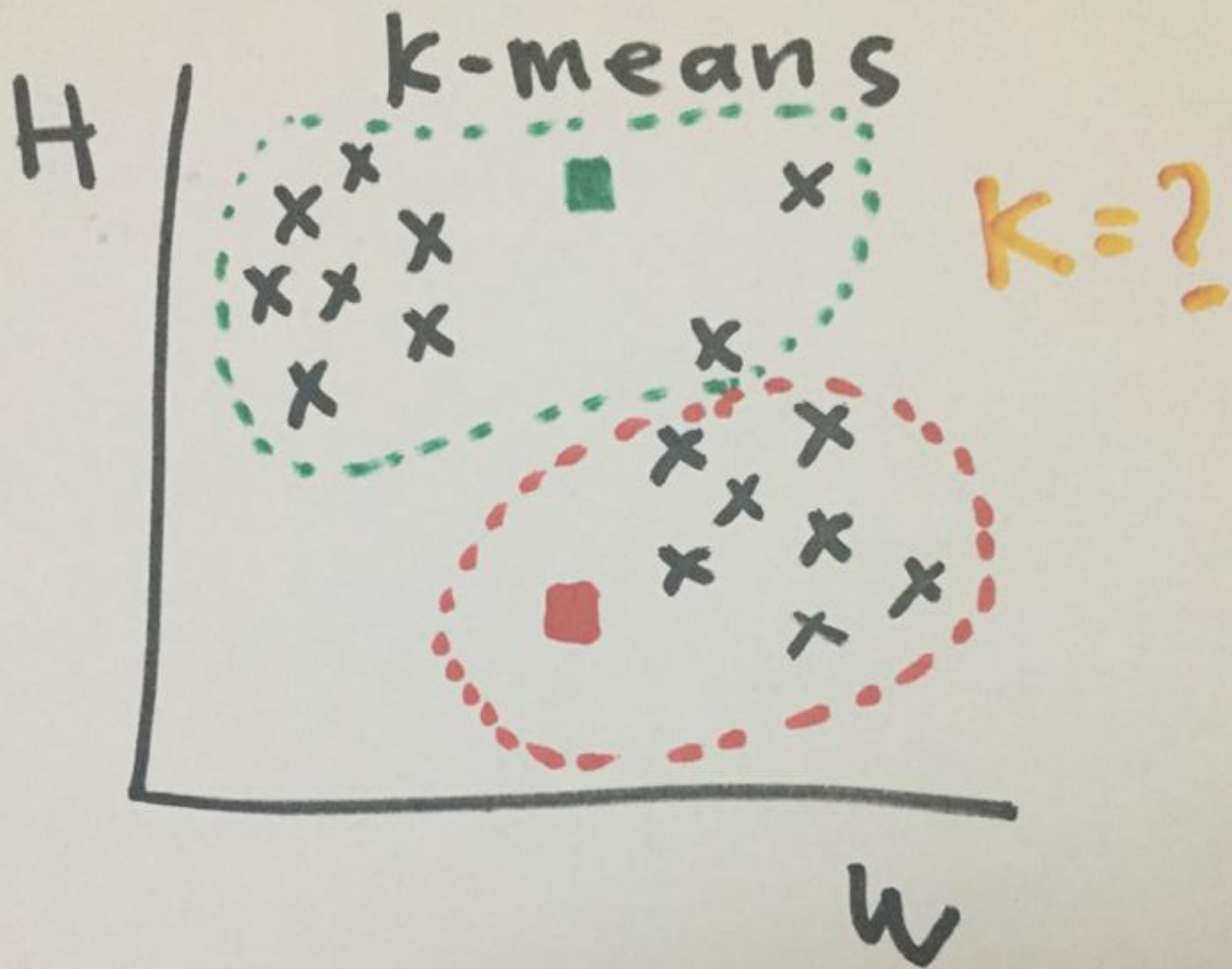
K-means

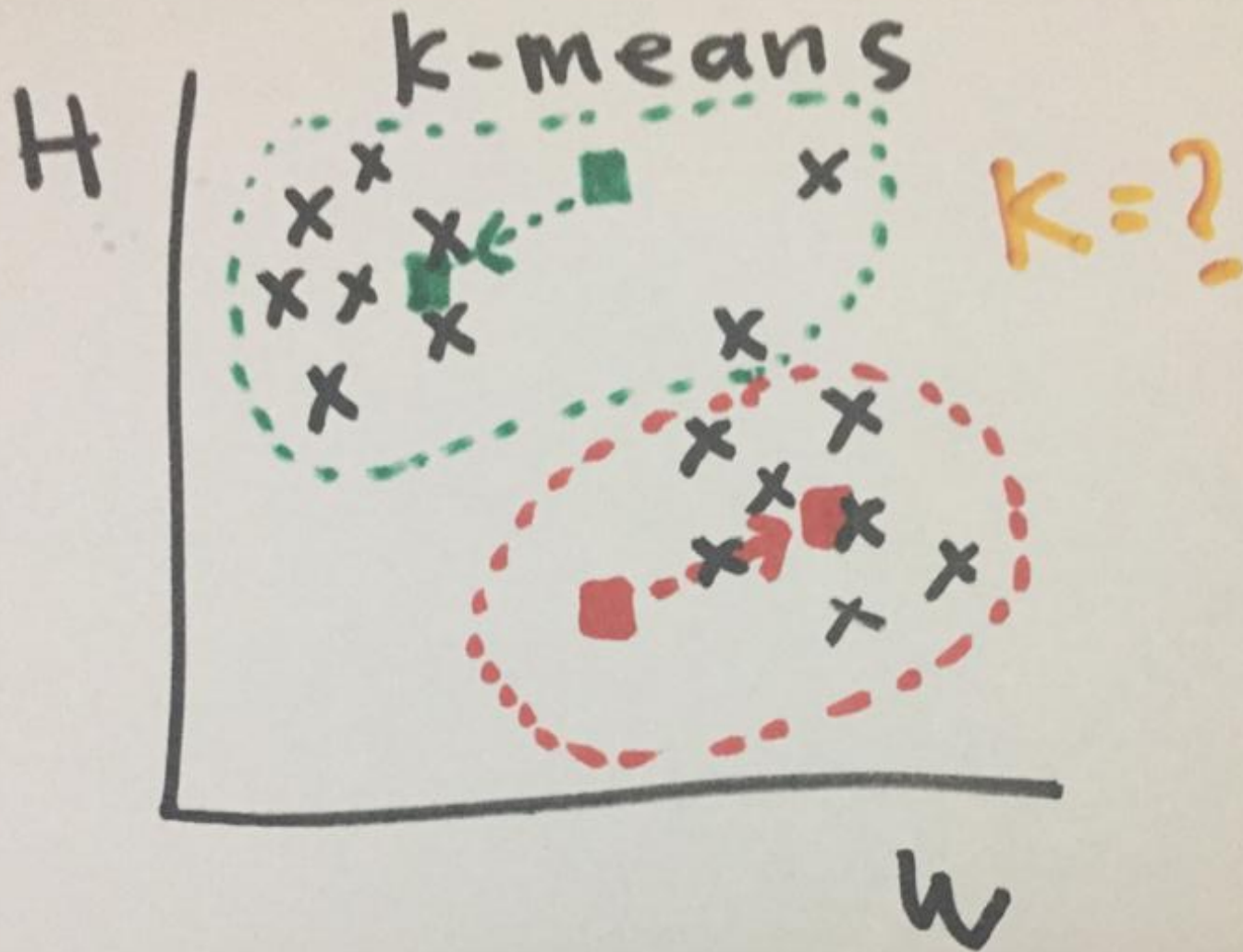
Adalah clustering teknik dimana observasi dikelompokkan ke dalam cluster yang berbeda berdasarkan nilai rata-rata (means) dari tiap cluster.

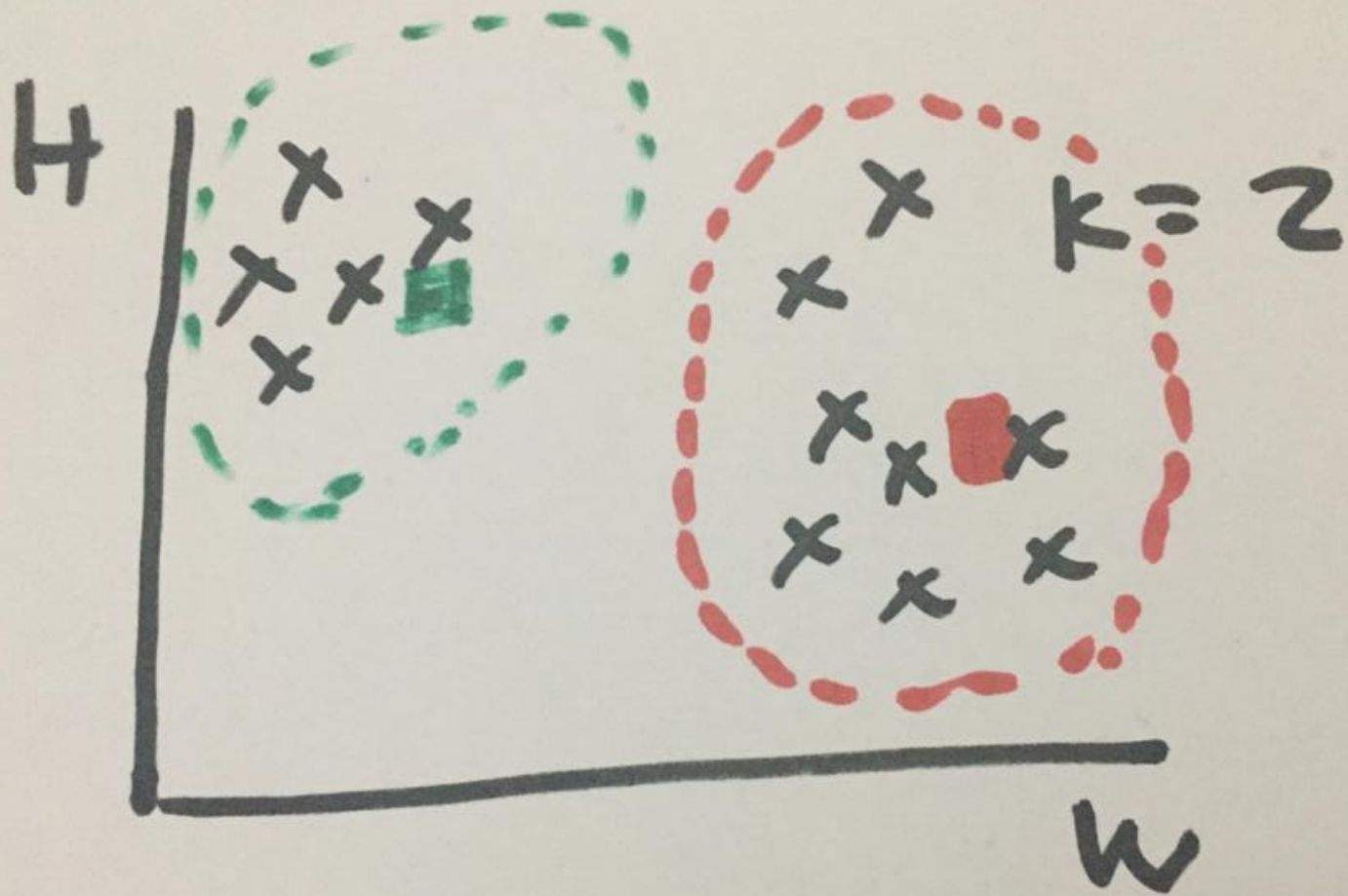
Problem: Seorang pegawai KPPN diminta atasannya untuk membuat kelompok profil dari satker -satker mitra untuk kebijakan baru

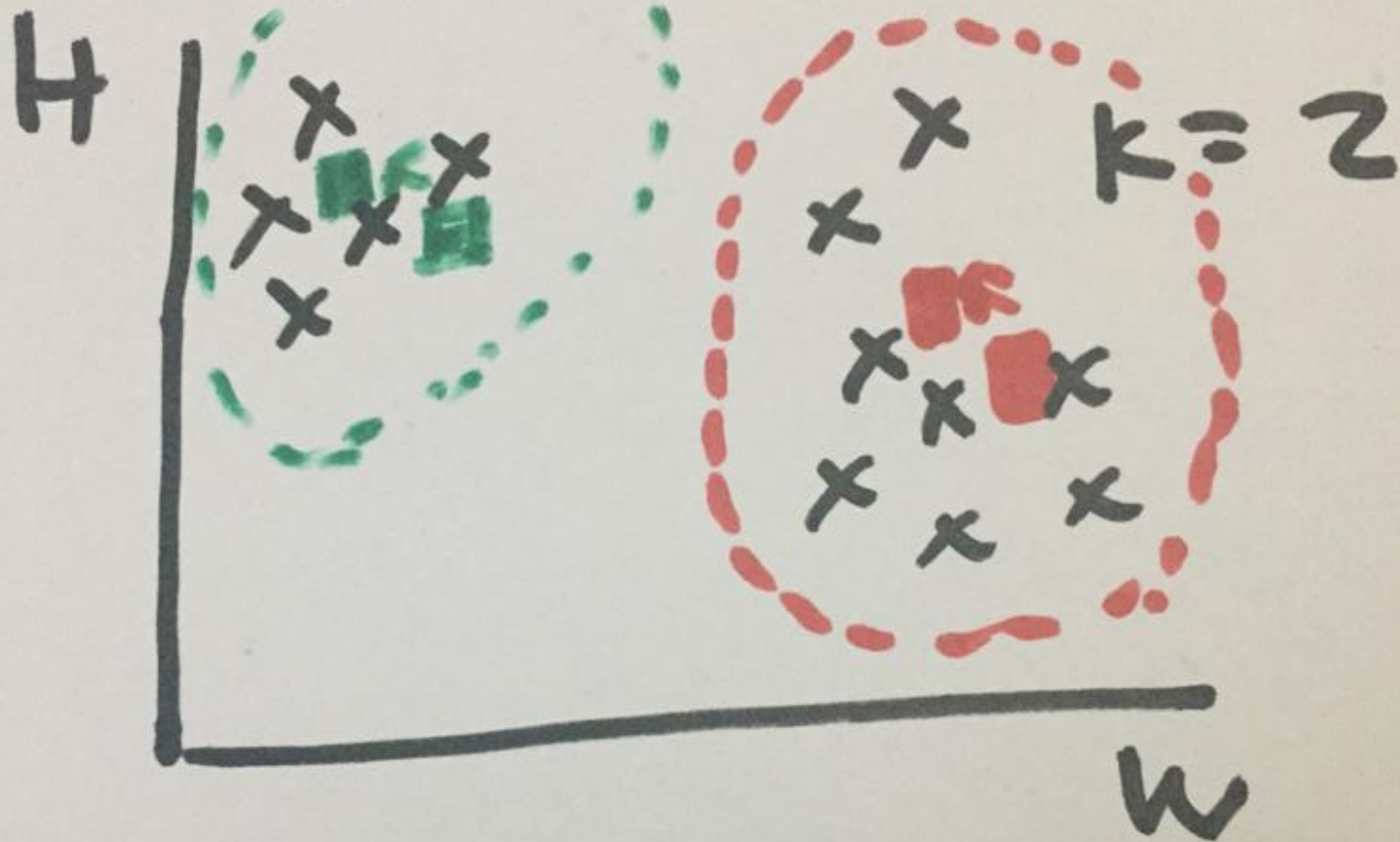
Hasil yang diharapkan: kelompok satker yang memiliki kesamaan













Unsupervised K-Means

When to stop?

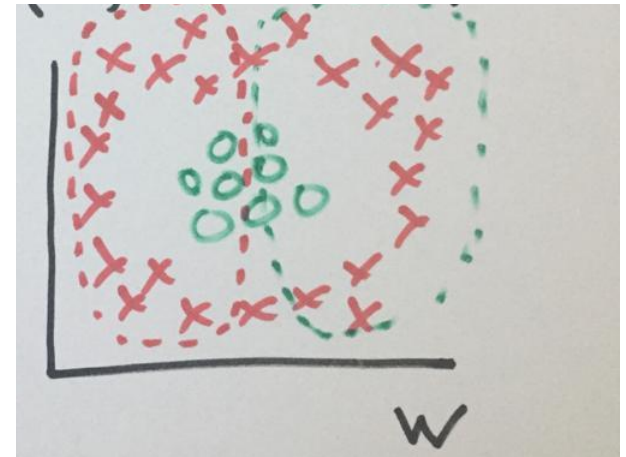
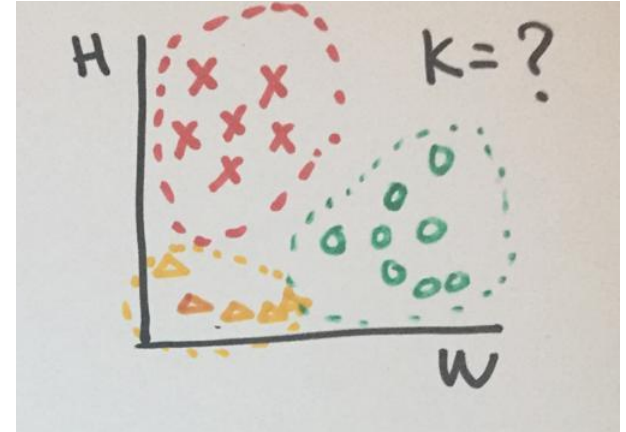
- Sampai centroid (kotak biru dan kotak hijau) tidak bergerak
- Sampai jumlah iterasi yang ditetapkan sebelumnya tercapai



Unsupervised K-Means

Plus : intuitif dan sederhana

Negatif: menentukan jumlah k ; terlalu sederhana



OH MY GOD

Linear Regression



I MADE A STRAIGHT LINE!



Supervised-Linear Regression

Regresi Linear adalah tehnik yang memprediksi variabel dependent (y) yang bersifat continuous berdasarkan variabel independen ($x_1, x_2 \dots x_n$) berdasarkan fungsi linear.

Problem: Seorang pegawai KPKNL diminta atasannya untuk memperkirakan nilai sewa kantin berdasarkan luas ruangan

Hasil yang diharapkan: nilai sewa kantin dalam rupiah



Supervised-Linear Regression

Task: memperkirakan nilai yhat

$$\text{yhat} = b + WX;$$

Dimana yhat: perkiraan nilai sewa kantin

X: vector luas ruangan kantin

W: Parameter weight/vector berupa koefisien

b: bias term.



Supervised-Linear Regression

Performance (P) : kita dapat hitung error untuk tiap observasi i sebagai:

$$e_i = \text{abs}(\hat{y}_i - y_i)$$

More popular: **Mean Squared Error (MSE)**

$$MSE = 1/2m \sum_i (\hat{y}_i - y_i)^2$$



Regularization (*remember bias vs variance?*)

Misal sebuah model yang kompleks (model dengan banyak variable) digunakan. **Low bias high Variance.**

Regularization : Reduce model complexity (reduce variance)

- Mengurangi Variance saat kita menggunakan banyak variable (complex model)
- Ridge : mengurangi koefisien/parameter shrinkage hingga variable yang tidak signifikan akan mendekati 0. Jumlah variabel yang digunakan oleh model tetap sama
- Lasso: mengurangi koefisien/parameter shrinkage sampai 0. Jumlah variable yang digunakan dapat berkurang
- Efek: **Higher Bias**



OLS (*remember bias vs variance?*)

- **Ridge** pada umumnya lebih cepat
- **Ridge** digunakan jika sebagian besar variable pengaruh yang sama terhadap response.
- **Lasso** umumnya digunakan saat hanya ada beberapa variabel yang signifikan/berpengaruh terhadap response



Supervised- Logistic regression

Problem : Pegawai DJP diminta untuk memprediksi apakah wajib pajak akan menyampaikan SPT pada akhir tahun

Kemungkinan yang akan terjadi: (1) wajib pajak menyampaikan (2) wajib pajak tidak menyampaikan

Probabilitas wajib pajak menyampaikan SPT: $0 \leq P \leq 1$



Supervised- Logistic regression

Pada dasarnya Logistic regression sama seperti linear regression

Logistic regression memprediksi probabilitas sebuah observasi berada dalam kategori positif

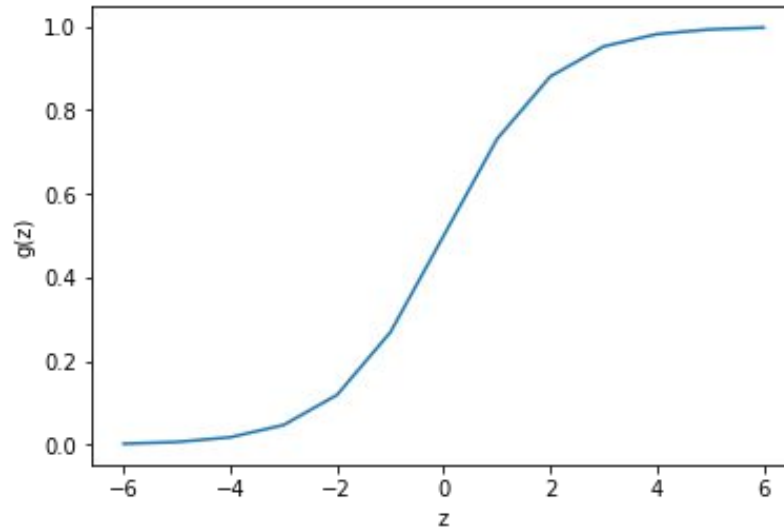
Target variabel pada Logistic regression adalah kategori (A atau bukan A)

Karena output merupakan probabilitas, gunakan sigmoid function



Supervised- Logistic regression

Sigmoid function: $g(z) = 1/(1+e^{-z})$





Supervised- Logistic regression

Jika sigmoid function dimasukkan ke dalam linear regression:

$$g(x) = \frac{1}{1 + e^{-(w_0 + w_f x)}}$$

Dimana $g(X)$ adalah probabilitas yhat merupakan kategori positive



Supervised- Logistic regression

Output logistic regression: angka probabilitas

Performance : Akurasi klasifikasi (\hat{y}) dengan threshold tertentu

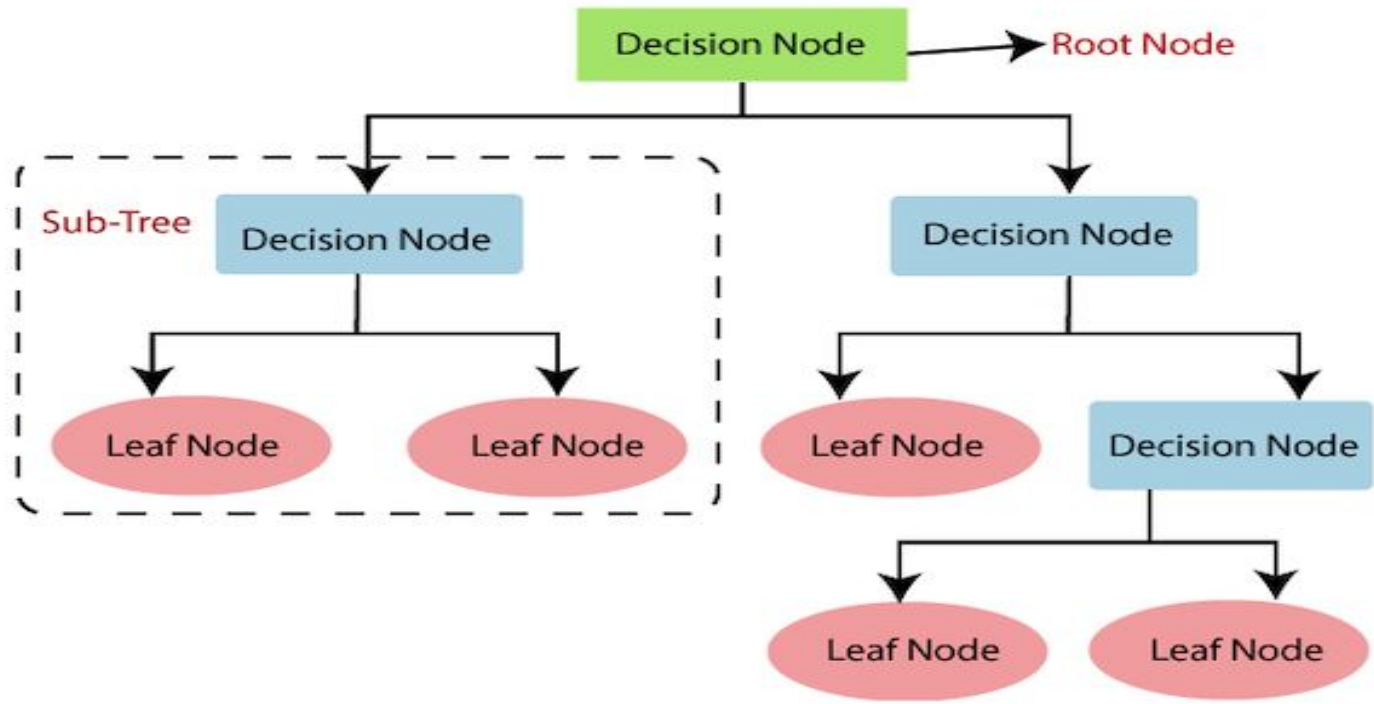
Training: mengoptimalkan akurasi

A photograph of a tree trunk with a dark, gnarled branch extending upwards and to the left. The text "Decision Tree" is overlaid in yellow.

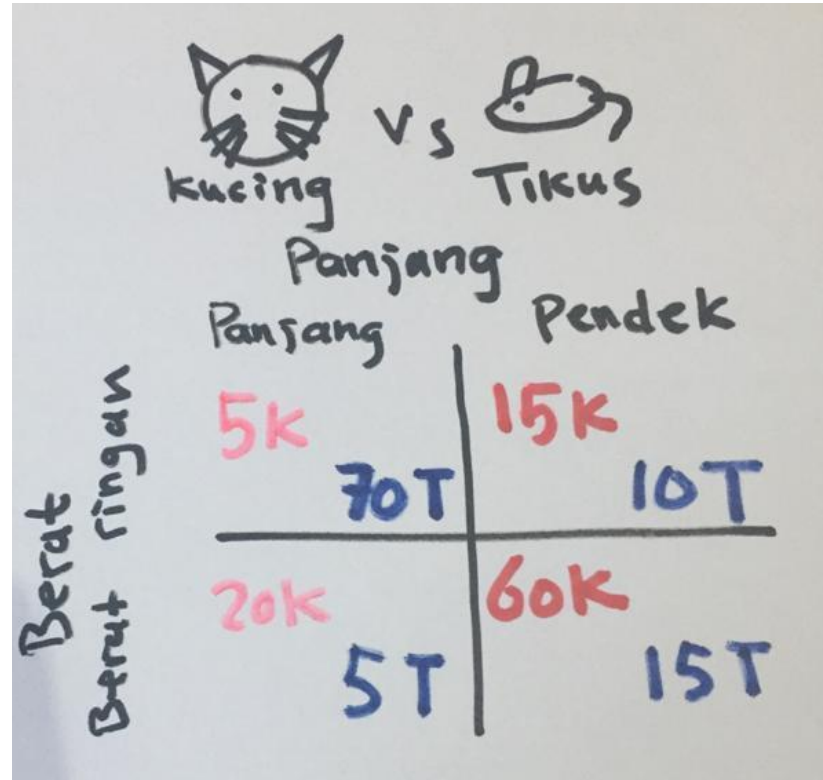
Decision Tree



Decision Tree

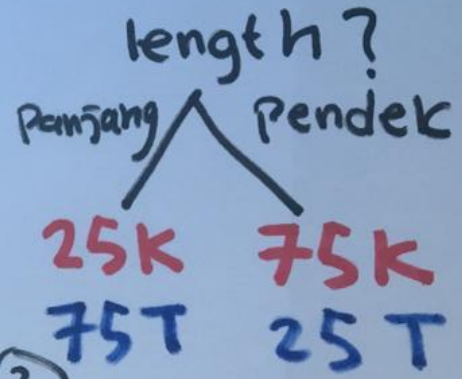
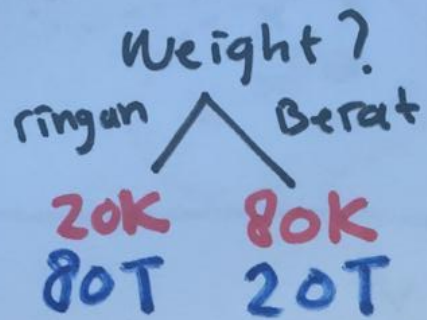


Tree (Decision Tree)



PT

ST



①

[P] Tikus $(T)^2 + (F)^2$
 $(\frac{4}{5})^2 + (\frac{1}{5})^2 = \frac{17}{25}$

[P] kucing $(T)^2 + (F)^2$
 $(\frac{4}{5})^2 + (\frac{1}{5})^2 = \frac{17}{25}$

$\frac{1}{2} \times \frac{17}{25} + \frac{1}{2} \times \frac{17}{25} = .68$

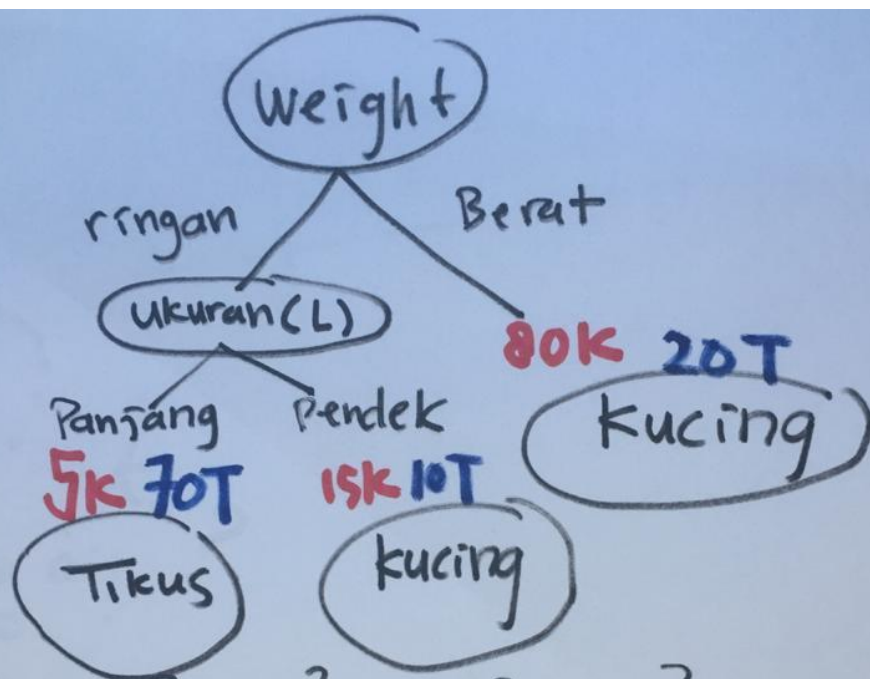
②

[P] Tikus $(T)^2 + (F)^2$
 $(\frac{3}{4})^2 + (\frac{1}{4})^2 = \frac{5}{8}$

[P] kucing $(T)^2 + (F)^2$
 $(\frac{3}{4})^2 + (\frac{1}{4})^2 = \frac{5}{8}$

$\Rightarrow \frac{5}{8} \times \frac{1}{2} + \frac{5}{8} \times \frac{1}{2} = .625$

P(B) P(P_d)



$$\left(\frac{70}{75}\right)^2 + \left(\frac{5}{75}\right)^2 + \left(\frac{15}{25}\right)^2 + \left(\frac{10}{25}\right)^2$$

$$\frac{75}{100} \cdot 87 + \frac{25}{100} \cdot 52$$

$$0.7825$$



Decision Tree

- Seperti namanya, Decision tree berbentuk cabang pohon
- Dapat digunakan baik untuk regression maupun classification problem
- Mampu memvisualisasi model dan keputusan secara jelas/eksplisit sehingga interetabilitasnya sangat tinggi
- Robust, **tidak rentan** terhadap outlier

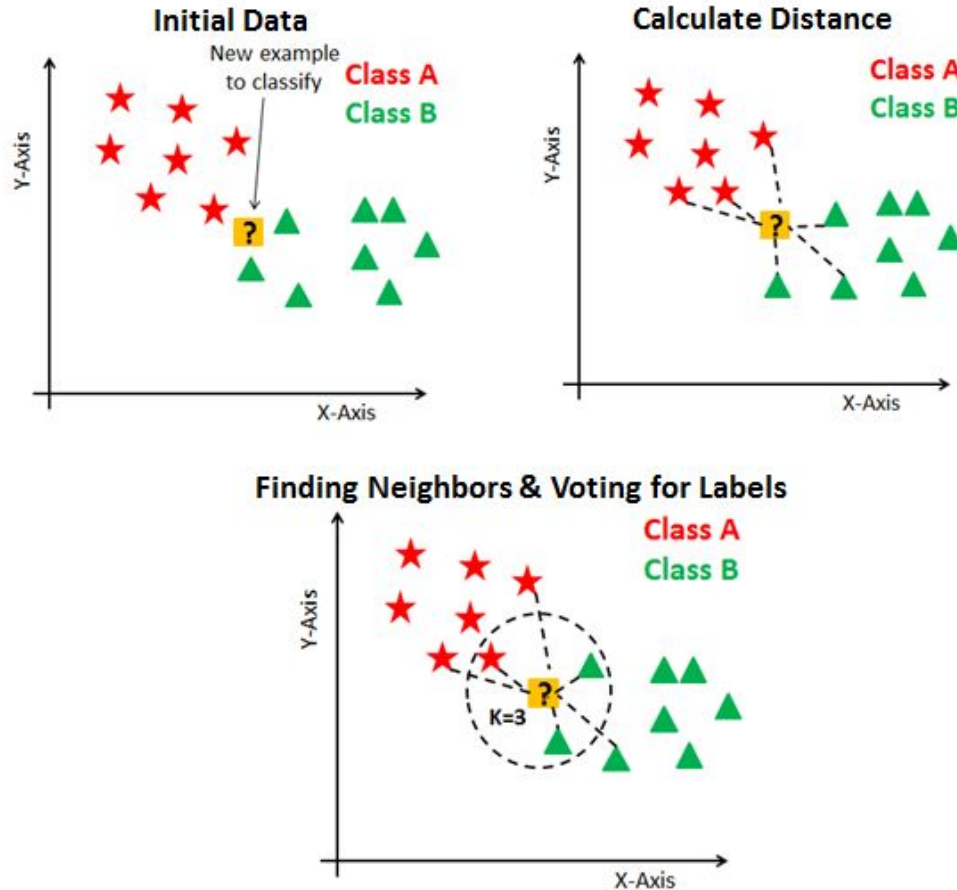




KNN- K nearest Neighbour

1. Menentukan sejumlah (k) observasi yang terdekat berdasarkan ukuran jarak yang ditetapkan sebelumnya (contoh nya euclidean distance)
2. Berdasarkan kategori “tetangga”, model memprediksi kategori suatu observasi dengan voting.
3. Jika hasil votingimbang maka kategori, dipilih secara random

KNN- K nearest Neighbour

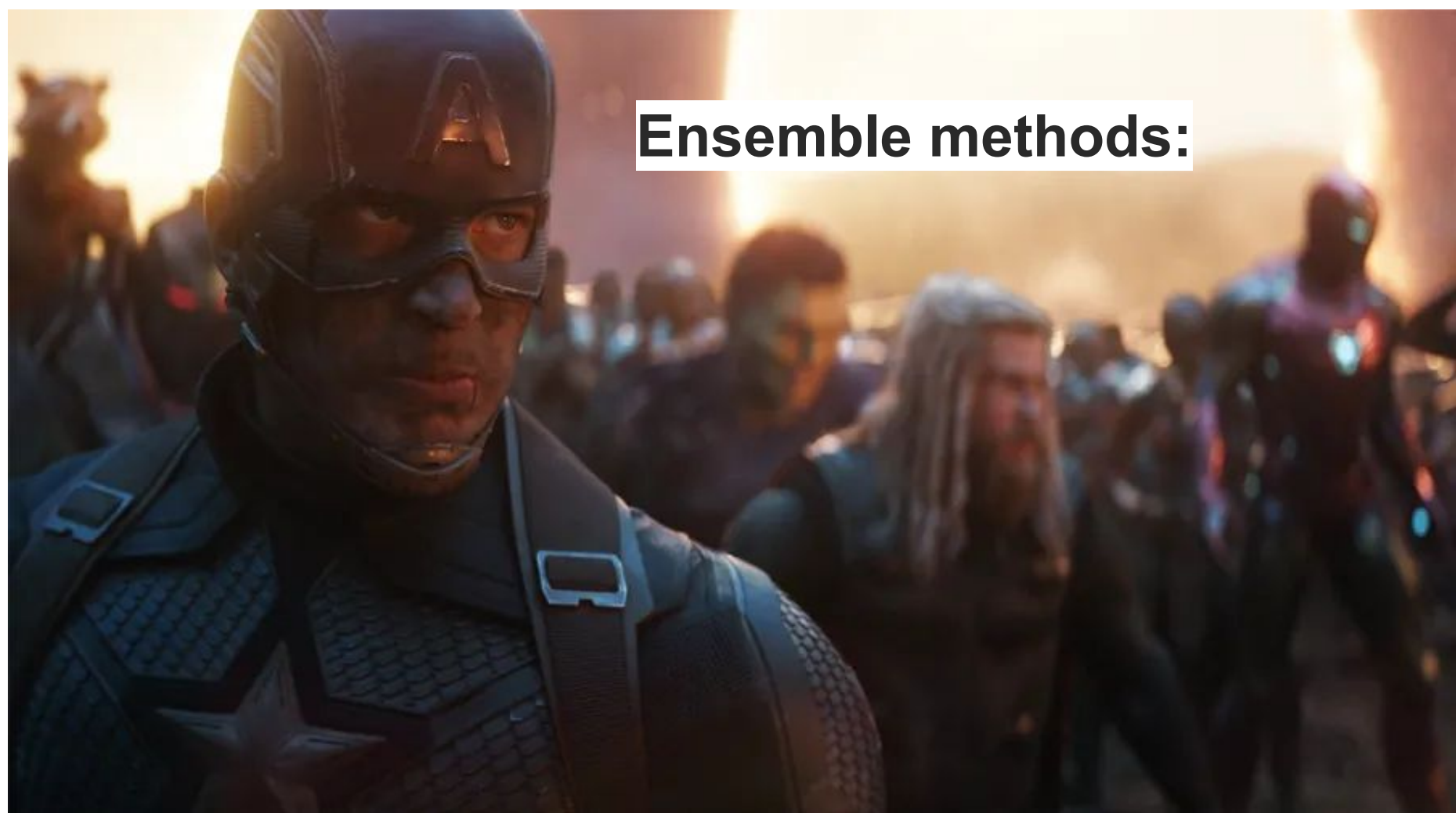




KNN

- + Robust, tidak terlalu terpengaruh oleh outliers
- Rentan terhadap satu variabel yang skalanya lebih besar daripada variabel lain

Ensemble methods:





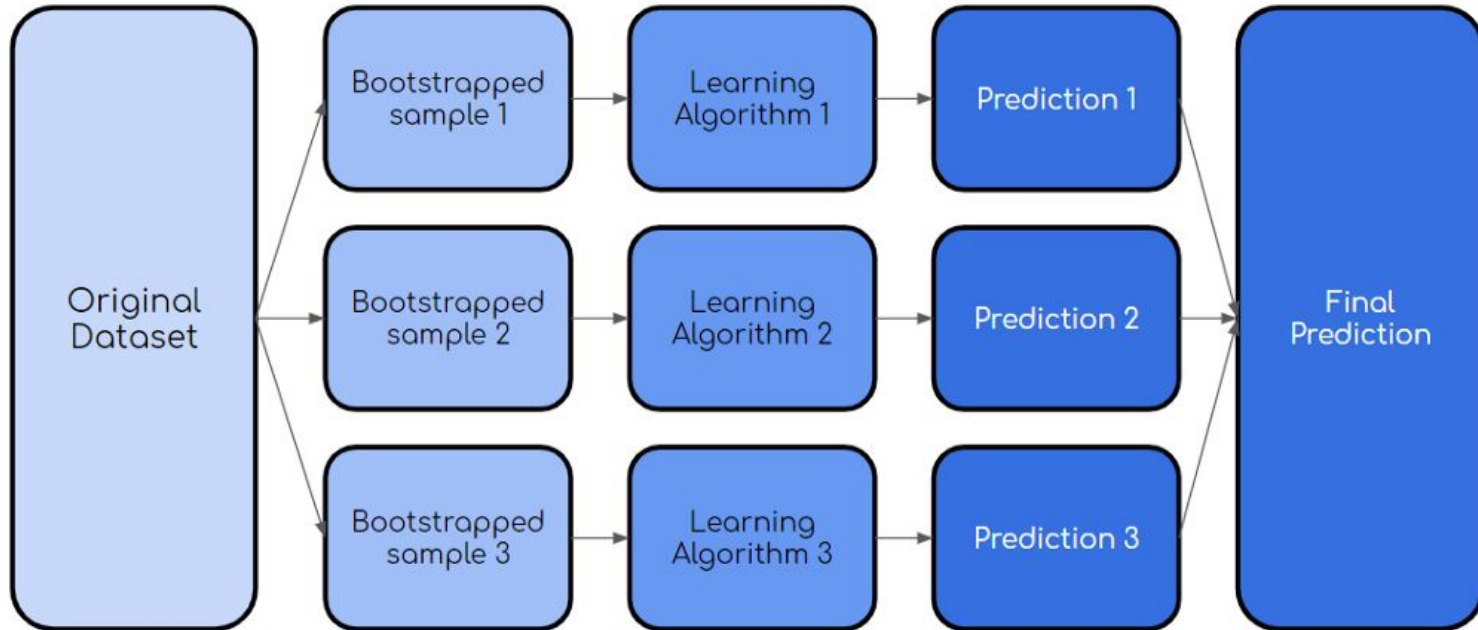
Ensemble methods: *bagging*, *boosting* and *stacking*

----- **weak learners**: Model-model yang memiliki bias atau variance yang tinggi namun **dapat** digunakan bersama untuk membentuk model yang lebih baik-----

- **Bagging** to decrease the model's **variance**;
- **Boosting** to decrease the model's **bias**;
- **Stacking** to increase the predictive force of the classifier.



Bagging-mengumpulkan banyak weak learners (sejenis) kemudian di aggregate



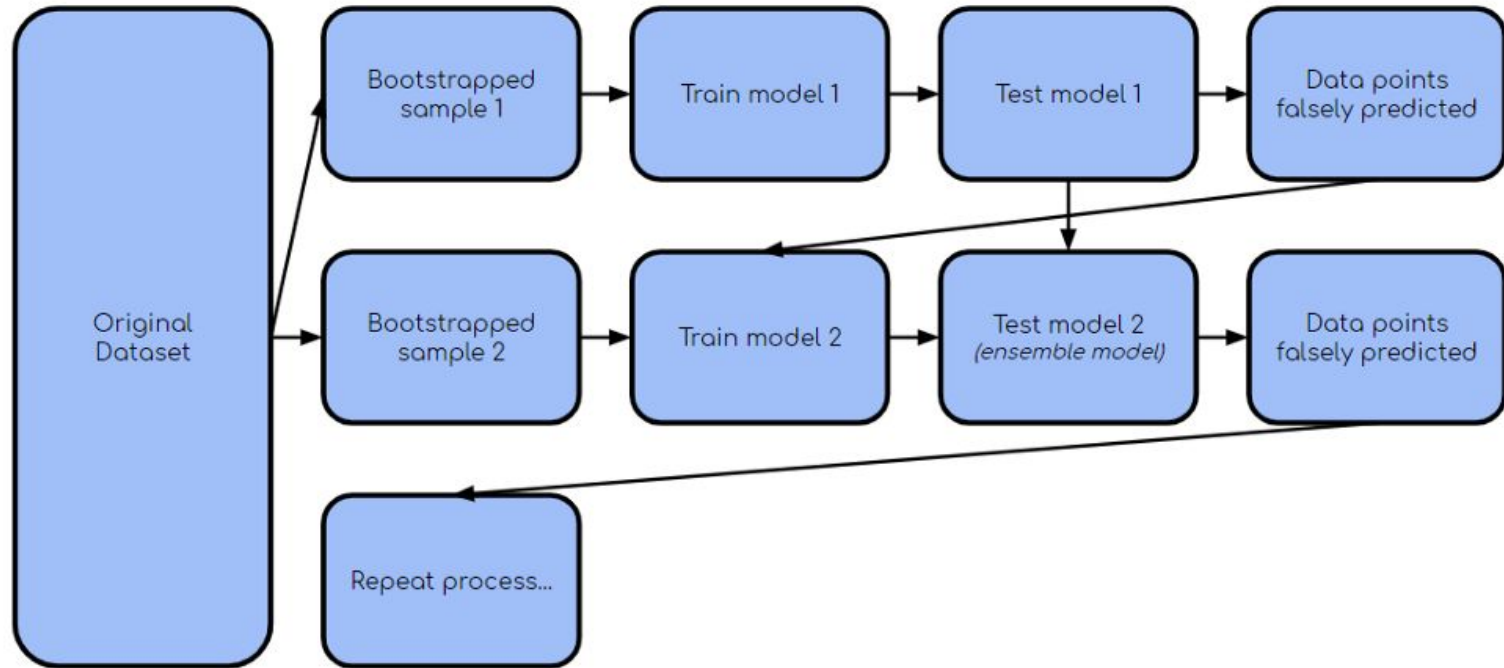


Example

- Bagged Decision Tree
- Random Forest



Boosting - mengembangkan weak learner secara sequential dan adaptive dengan weak learner lain (sejenis).

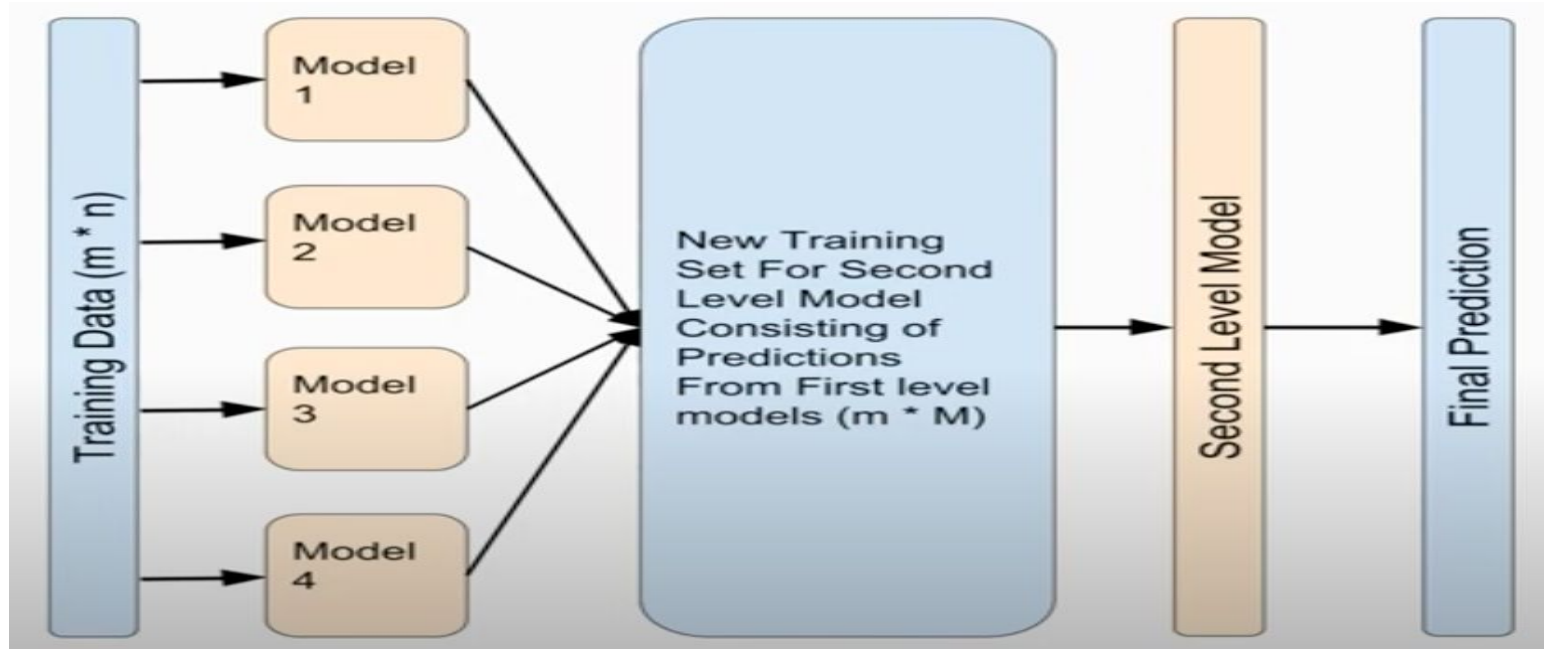




Example

- Adaptive Boosting
- Gradient Boosting

Stacking - menggabungkan weak learners dengan weak learner yang berbeda algoritma (heterogeneous)





Example

- Stacked Model