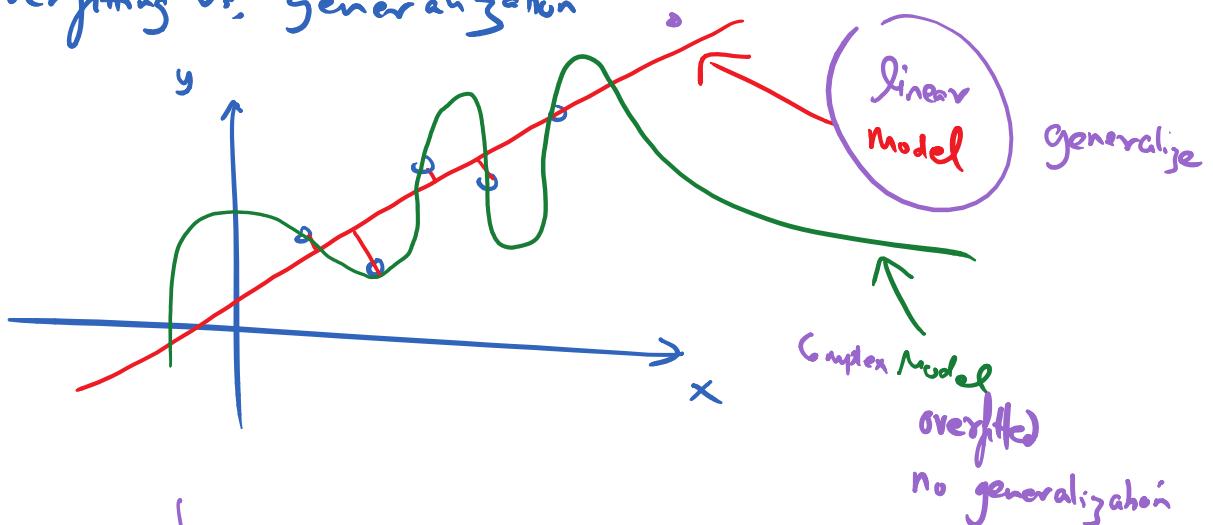


Bias - Variance Decomposition

Overfitting vs. generalization



Loss function : classifier M

$\hat{y} \leftarrow M(x)$
Predicted class

y vs. \hat{y}
true class prediction

Classification : 0-1 loss function

$$L(y_i, \hat{y}_i) = I(y_i \neq \hat{y}_i)$$

Expected loss:

$$1 - \dots - n$$

$$I(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{if } z \text{ is false} \end{cases}$$

Indicator Function

Expected loss:

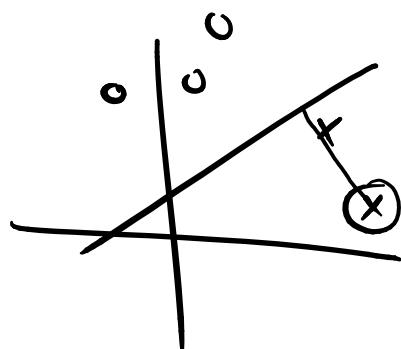
$$\begin{aligned} & \frac{1}{n} \sum L(y_i, \hat{y}_i) \\ &= \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i) \\ &= \frac{1}{n} \{ \# \text{errors} \} \\ &= \text{error-rate} \end{aligned}$$

Indicator Function

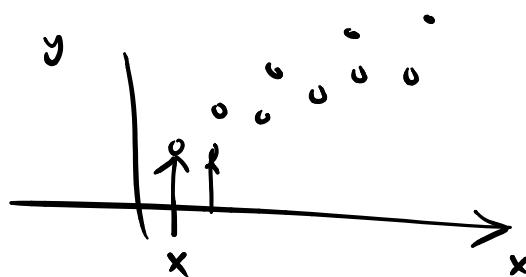
Squared loss:

$$L(y_i, \hat{y}_i) = (\vec{y}_i - \vec{\hat{y}}_i)^2$$

true value predicted value



Regression problem



Given the points
 $\vec{x}_i \in \mathbb{R}^d$
 predict a "real" variable
 $y \in \mathbb{R}$

Expected loss?

$$\begin{aligned} & E[L(y_i, \hat{y}_i)] \\ &= E[(y_i - \hat{y}_i)^2] \\ &= E[(y_i - E[y_i] + E[y_i] - \hat{y}_i)^2] \end{aligned}$$

One test case x_i

$$= E[(y_i - E[y_i]) + (E[y_i] - \hat{y}_i)]^2$$

$$= E\left[\frac{(y_i - E[y_i])^2}{a^2} + \frac{(\hat{y}_i - E[\hat{y}_i])^2}{b^2} + \frac{2(y_i - E[y_i])(E[\hat{y}_i] - \hat{y}_i)}{ab}\right]$$

$$= E[(y_i - E[y_i])^2] + E[(\hat{y}_i - E[\hat{y}_i])^2] +$$

true expectation prior expected true value
 Variance of y_i
 avg noise

$$2E(y_i - E(y_i))(-) + 2(E[\hat{y}_i] - E(y_i))(-)$$



$$E_{\hat{y}_i}[(\hat{y}_i - E[\hat{y}_i])^2]$$

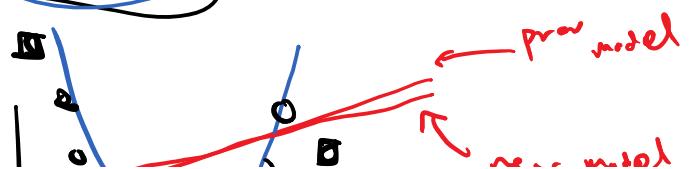
Over all possible test case

$$E\left[\frac{(\hat{y}_i - E[\hat{y}_i])^2}{b^2} + \frac{(y_i - E[y_i])^2}{a^2}\right]$$

$$\text{loss} = \left[E[(\hat{y}_i - E[\hat{y}_i])^2] \right] + E[(E[\hat{y}_i] - E[y_i])^2] + E[(y_i - E[y_i])^2]$$

avg bias of model
 avg noise of data

Simple model
linear classifier

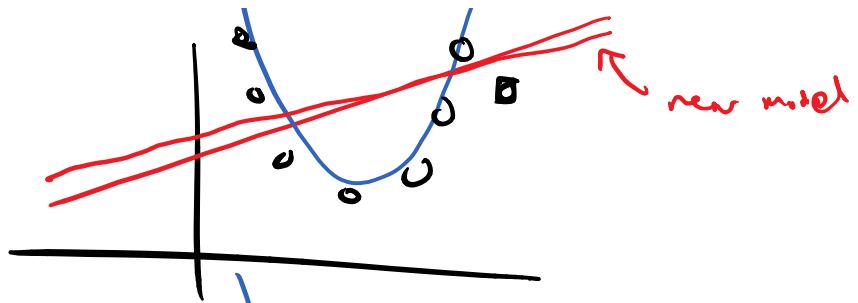


Simple model

linear classifier

→ high bias

→ low variance

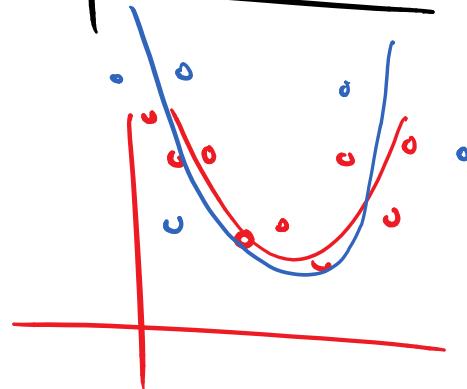


Complex model

quadratic decision model

→ low bias

→ high variance



bias-variance tradeoff

more complex the model → low bias → can

increase
variance

keep the simplest model with low variance

How to get low variance & low bias?

Ensemble classifiers

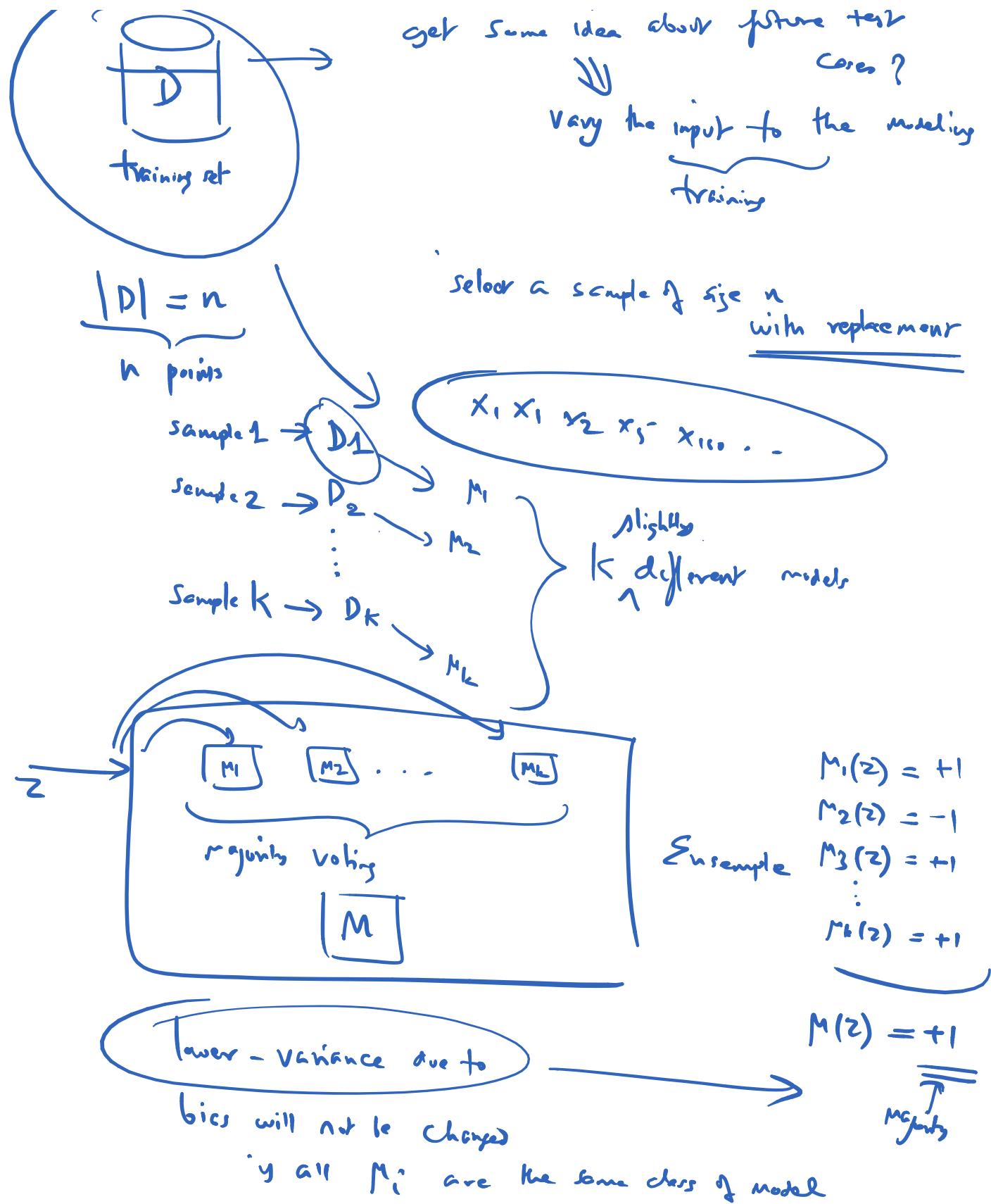
① Bagging → Bootstrapped Aggregation

→ lower variance due to averaging.

Bootstrap sampling



get some idea about future test
cores?

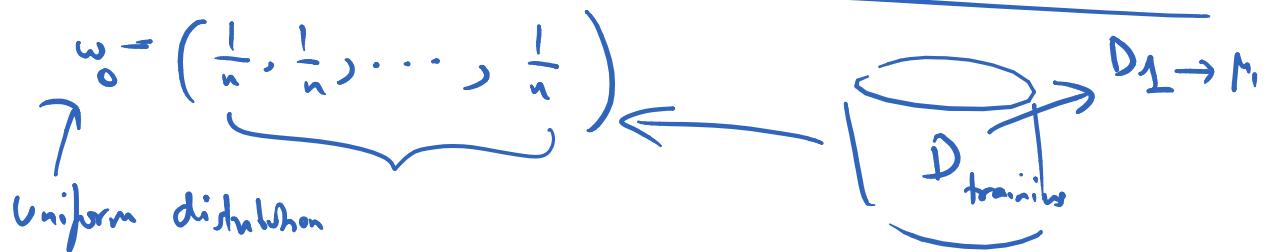


Boosting \rightarrow weighted sampling + quality of classifier

↓

boost or bump up
the weight for
Misclassified points

$d_i \nearrow$
weighted vote.



① pick a sample of size n from D_{train} using \underline{w} as the prob vector
(with replacement)

② Build classifier M_1

③ $\varepsilon_1 = \text{error rate of } M_1$

$$0 < \varepsilon_1 < 0.5^-$$

← error rate

④ $S_1 = \{ \text{set of all misclassified points, } \hat{y}_i \neq y_i \}$

$$\alpha_1 = \ln\left(\frac{1}{\varepsilon_1} - 1\right)$$

← weight of classifier

$$\text{⑤ } \underline{w}'_i = \begin{cases} w_i \cdot y \quad y_i = \hat{y}_i & \ln\left(\frac{1}{0.5} - 1\right) \\ w_i \cdot e^{\alpha_1} \cdot y \quad y_i \neq \hat{y}_i & \ln(2 - 1) = \ln(1) = 0 \end{cases}$$

new weight vector

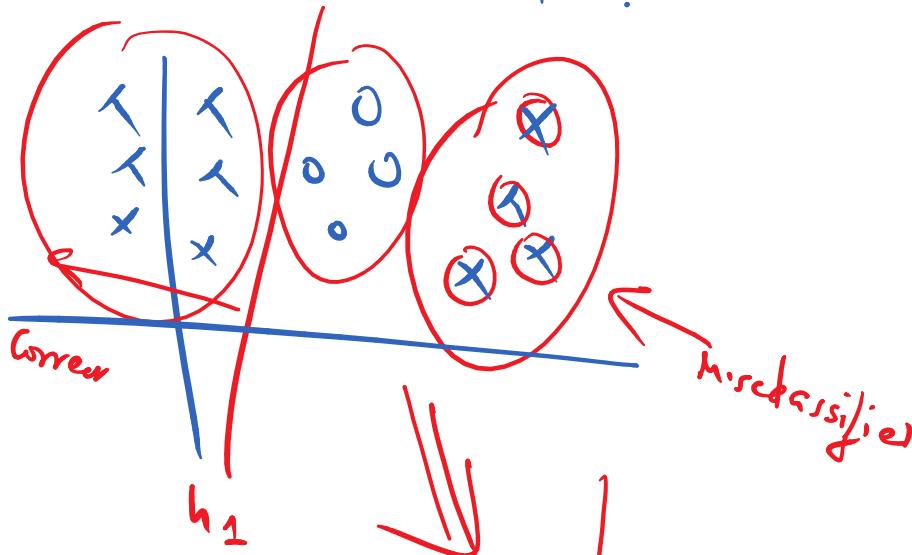
$$\underline{w}'_i = \begin{cases} w_i & \ln\left(\frac{1}{0.5} - 1\right) \\ w_i \cdot e^{\alpha_1} & \ln(2 - 1) = \ln(2) \end{cases}$$

⑥ normalize to be prob. vector

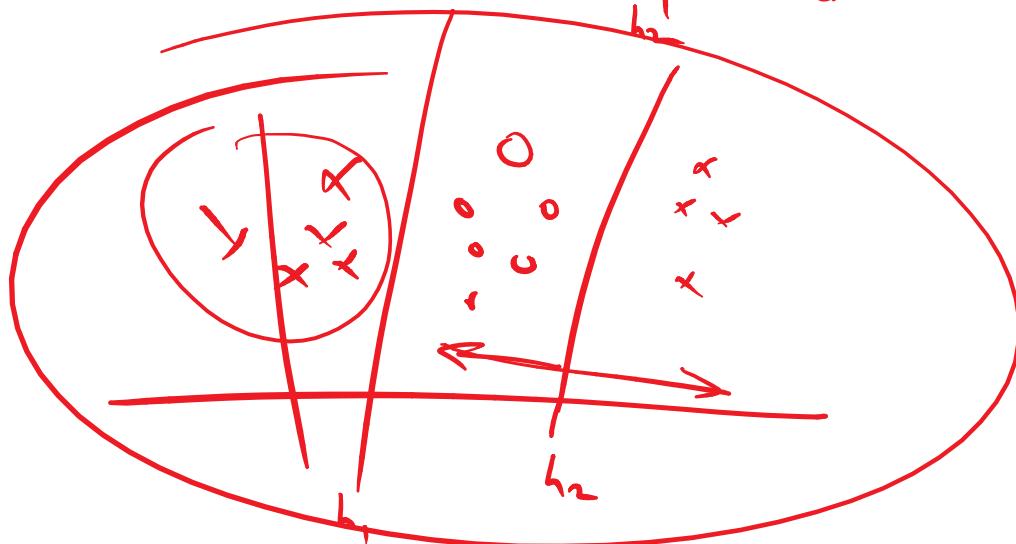
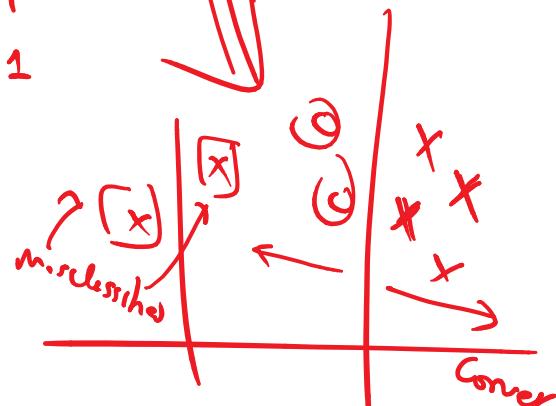
$$w'_i = \frac{w'_i}{\sum w'_i}$$

$$v_i = \frac{v_i}{\sum v_i}$$

repeat K times.



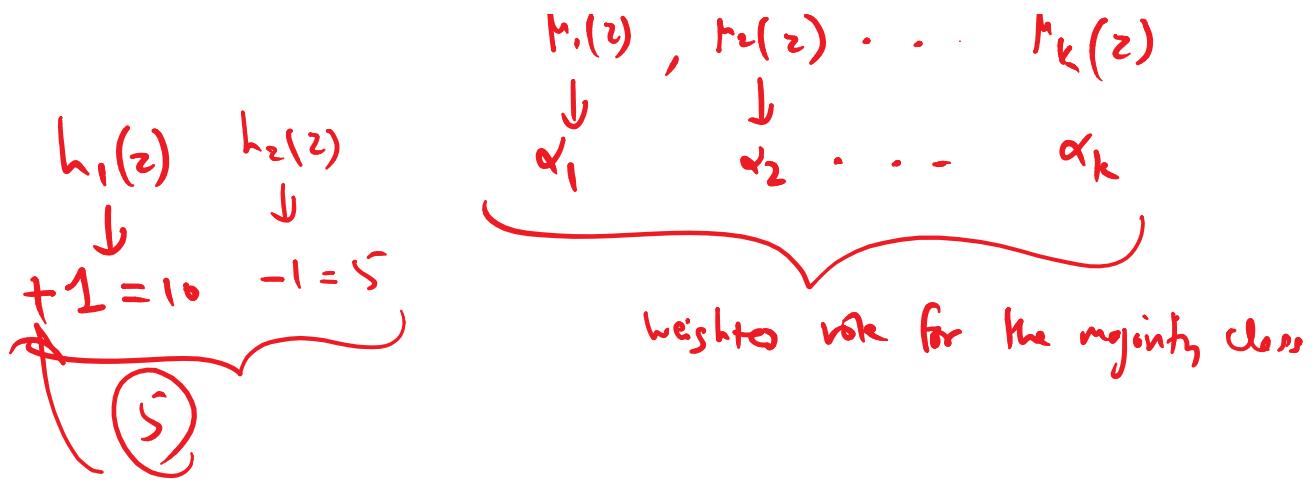
Linear SVM



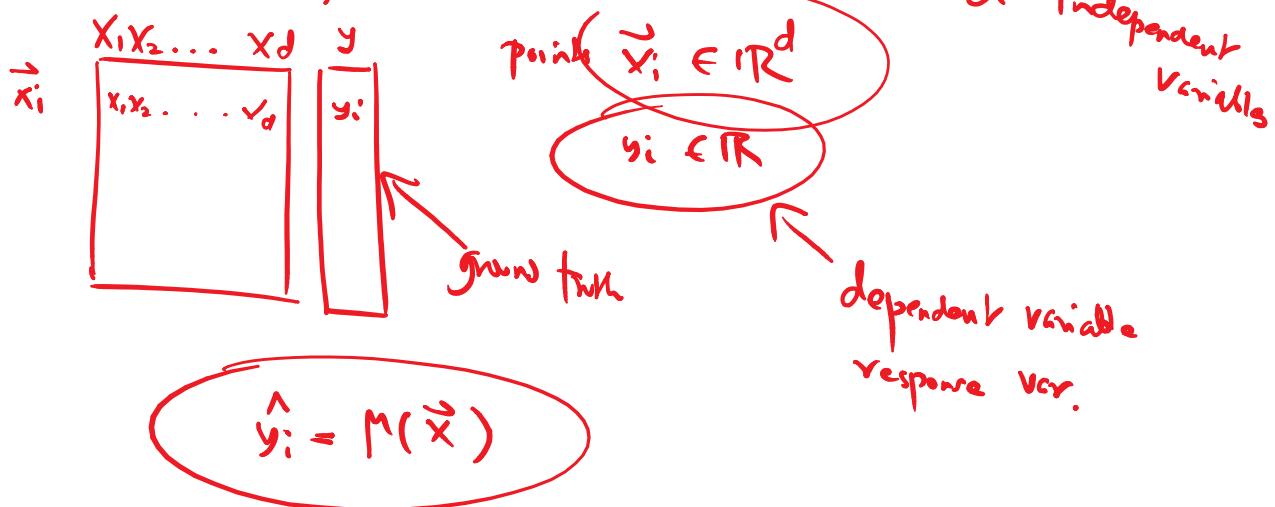
$$h_1 + h_2$$

$$\alpha_1 \quad \alpha_2$$

Ensemble: \rightarrow



Regression



linear regression.

$$\hat{y}_i = w_0 + \underbrace{w_1 x_1 + w_2 x_2 + \dots + w_d x_d}_{\text{bias}} = w_0 + \overrightarrow{w} \cdot \overrightarrow{x}$$

dot product

$$\vec{x}_i = (1 \ x_1 x_2 \dots x_d)$$

$$\hat{y}_i = \overrightarrow{w} \cdot \overrightarrow{x}$$

$$w = (w_0 \ w_1 \ w_2 \dots \ w_d)$$

$d+1$ dimensions

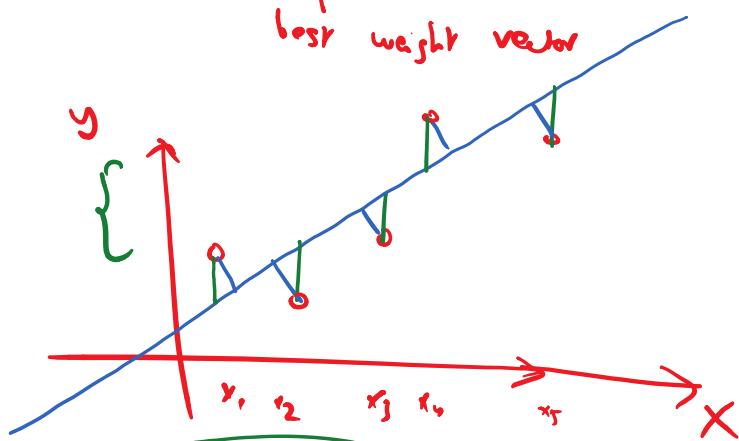
weight vector

$$J(w) = \min \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

Squared error

$$J(\omega) = \min_{\omega} \left(\sum_{i=1}^n (y_i - \hat{y}_i)^2 \right)$$

Squared error



$\hat{y}_i = \vec{\omega}^T \vec{x}_i$

scalar dot product

$\vec{\omega} \in \mathbb{R}^{d+1}$

$\vec{x}_i \in \mathbb{R}^{d+1}$

PCA:
best line
orthogonal sq. error

Min error along all attributes

Regression
"best line": along y only.
"vertical" sq. error

n -equations; one per point

Combine into one equation

$\hat{\vec{y}} = \vec{X} \vec{\omega}$

$$\begin{aligned} J(\omega) &= \min_{\omega} \left\| \vec{y} - \hat{\vec{y}} \right\|^2 \\ &= \min_{\omega} \left\| \vec{y} - \vec{X} \vec{\omega} \right\|^2 \\ &= (\vec{y} - \vec{X} \vec{\omega})^T (\vec{y} - \vec{X} \vec{\omega}) \end{aligned}$$

$X = D$

$n \times (d+1)$ matrix of points

$\vec{y} = n \times 1$ vector of predicted values

$\vec{\vec{y}} = n \times 1$ vector of true values

$\vec{\omega} = \mathbb{R}^{d+1}$

$$\min_{\omega} J(\omega) = \underbrace{\hat{y}^T \hat{y}} + \underbrace{\omega^T X^T X \omega}_{\text{散射矩阵}} - 2(X^T \omega)^T y$$

$$\frac{\partial J(\omega)}{\partial \omega} = \underbrace{2X^T X \omega}_{\text{散射矩阵}} - 2X^T y = 0$$

$$\Rightarrow X^T X \omega = X^T y$$

$$\hat{\omega} = (X^T X)^{-1} X^T y$$

Closed form solution

$$\hat{y} = \hat{m}(x) = \hat{\omega}^T \hat{x}$$

← regression model