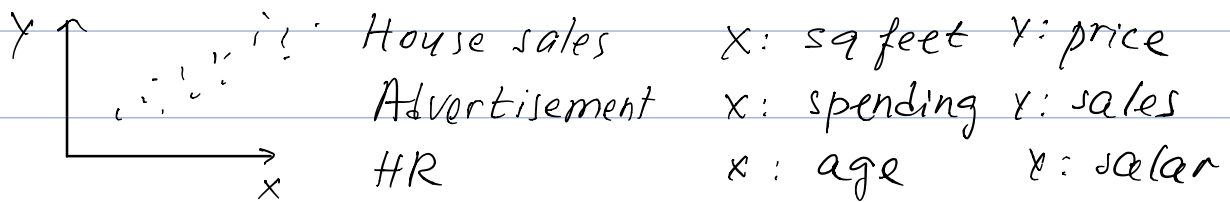


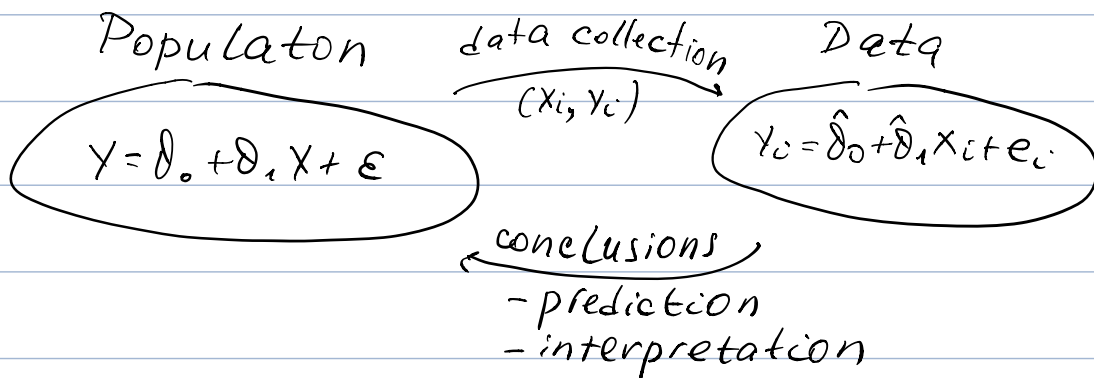
Day 1
9/7/18

Goal of supervised machine learning:

use labeled example data to make
and interpret predictions



Process of supervised ML:



Components: x_i - input, $x_i \in \mathcal{X} = \mathbb{R}^p$
 y_i - output; target, $y_i \in \mathcal{Y} = \mathbb{R}$
 $i = 1, \dots, N$
 (x_i, y_i) - training example
 $\{(x_i, y_i)\}_{i=1}^N$ - training set

Goal: develop function $h(x): \mathcal{X} \rightarrow \mathcal{Y}$
 $h(x)$ - hypothesis, model
 $h_\theta(x)$ - hypothesis depends on parameters θ

Process:

- (1) define problem space
- (2) collect data
- (3) specify model / hypothesis
- (4) develop learning algorithm
(i.e., find values of params)
- (5) make predictions

Terminology $h_{\theta}(x) : X \rightarrow Y$ - supervised ML

$Y = \mathbb{R} \rightarrow$ regression

$Y = \text{set of discrete values}$ - classification,
patterns of X - unsupervised ML

Simple (one-variable) linear regression

Step (3): $h_{\theta}(x) = \theta_0 + \theta_1 x = \theta_0 \cdot 1 + \theta_1 \cdot x$

$\underbrace{\theta_0, \theta_1}_{\substack{\text{parameters} \\ \text{weights} \\ \text{constants}}} \quad \underbrace{1, x}_{\substack{\text{fixed} \\ \text{basis}}}$

vector notation: $h_{\theta}(x) = \begin{bmatrix} 1 & x \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = \underset{1 \times 2}{X^T} \underset{2 \times 1}{\theta}$

Step (4): to estimate (i.e., determine from data) the values of θ_0 and θ_1 , minimize cost function (= loss function)

$$J(\theta) = \sum_{i=1}^N (y_i - h_{\theta}(x_i))^2 - \text{squared loss}$$

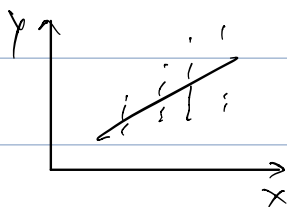
(many other loss functions are possible)

Params $\hat{\theta} = \arg \min_{\theta} J(\theta)$ - least squares estimates

Prediction: $\hat{y}_i = \hat{\theta}_0 + \hat{\theta}_1 x_i$

Probabilistic interpretation

Randomness in the data: consider probability distribution of Y given values of X



$P\{Y|X\}$ - class of such models are called discriminative models

Special case: Normal linear regression

Assume: $Y_i | X_i \stackrel{\text{indep}}{\sim} \mathcal{N}(\beta_0 + \beta_1 X_i, \sigma^2)$

Equivalently: $Y_i | X_i = \beta_0 + \beta_1 X_i + \varepsilon_i, \varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$
independent and identically distributed

Decomposition of the overall variation:

$$\begin{aligned} E\{(Y - \hat{Y})^2\} &= E\{(h_0(X) + \varepsilon - \hat{h}_0(X))^2\} = \\ &= E\{(h_0(X) - \hat{h}_0(X))^2 + 2E\{(h_0(X) - \hat{h}_0(X))\varepsilon\} + E\{\varepsilon^2\}\} \\ &= E\{(h_0(X) - \hat{h}_0(X))^2\} + \sigma^2 \\ &\quad \uparrow \text{reducible error} \quad \uparrow \text{irreducible error} \end{aligned}$$

Parameter estimation: maximum likelihood

$$f(Y_i | X_i, \beta) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2\sigma^2}(Y_i - \beta_0 - \beta_1 X_i)^2}$$

Probability distribution: $f(\text{data} | \text{params})$

Likelihood: $f(\text{params} | \text{data})$

Params that maximize the likelihood are maximum likelihood estimates

Bivariate Normal distributions:

regression vs correlation

$$f(Y_1, Y_2) = \frac{1}{2\pi \sigma_1 \sigma_2 \sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)} \left[\frac{(Y_1 - \mu_1)^2}{\sigma_1^2} - 2\rho \left(\frac{Y_1 - \mu_1}{\sigma_1} \right) \left(\frac{Y_2 - \mu_2}{\sigma_2} \right) + \frac{(Y_2 - \mu_2)^2}{\sigma_2^2} \right] \right\}$$

where ρ is the coef. of correlation

$$\rho = \frac{\text{Cov}(Y_1, Y_2)}{\sigma_1 \sigma_2} = \frac{E\{(Y_1 - \mu_1)(Y_2 - \mu_2)\}}{\sigma_1 \sigma_2}$$

Can show: $Y_1 | Y_2$ is Normally distributed

$$E\{Y_1 | Y_2\} = \left(\mu_1 - \mu_2 \rho \frac{\sigma_1}{\sigma_2} \right) + \left(\rho \frac{\sigma_1}{\sigma_2} \right) Y_2$$

$$\text{Var}\{Y_1 | Y_2\} = \sigma_1^2 (1 - \rho^2) \text{ - variance of the conditional distr. } \downarrow \text{ if correlation } \uparrow$$

