

Two types of learning problems:

- ① supervised problems
- ② un-supervised problems

← first half of course

①

Supervised learning:

→ we have some "training" examples we to learn from them

basically a predictor problem:

Want to predict some Y from X
↑ ↑
have examples

two flavors:

① "Regression": predicting a continuous outcome

Ex. predict stock price from economic indicators

Ex. predict adult height from childhood diet

② "Classification": predicting a categorical outcome

Ex. Will an individual default on a loan given credit score?

Outcome - $\{0 = \text{don't default}\}$

$$\text{Outcome} = \begin{cases} 0 = \text{don't default} \\ 1 = \text{will default} \end{cases}$$

Ex. Predict disease status from bloodwork

$$\text{outcome} = \begin{cases} 0 = \text{no-disease} \\ 1 = \text{have disease} \end{cases}$$

binary problems \nearrow

Ex. predict race from genomic data
"Multi-class problem"

$$\text{outcome} = \begin{cases} \vdots \end{cases}$$

② Unsupervised problems

Don't have "training data".

We want to learn/summarize important trends in the data.

(No specific quantity we are predicting.)

General Setup for Supervised problems

We have some variable Y

called: outcome, predicted var., dependent
the thing we want to predict

the thing we want to predict
we're going to try to predict Y using
other variables X_1, \dots, X_p # vars.

called: predictors, features, covariates,
independent variables.
the things used to predict Y .

Predict?

want to come up w/ a function \hat{f} so
that
$$Y \approx \underbrace{\hat{f}(X_1, \dots, X_p)}_{\hat{Y}}$$

Goal of this course:

(methods) (1) how do we construct \hat{f} ?

(evaluation) (2) how do we determine if \hat{f} is good?

Way supervised problems work:

we will use training data to construct \hat{f}

Statistical Learning

Assume we have some training data

training data size

$$(y_n, x_n) \text{ for } n=1, \dots, N$$

where

$$y_n \in S \text{ and } x_n \in \mathbb{R}^p$$

measured covariates for n^{th} obs.

$$\text{e.g. } x_n = (x_{n1}, x_{n2}, \dots, x_{np}) \in \mathbb{R}^p$$

$S = \mathbb{R}$ for regression problems

$S = \{c_1, c_2, c_3, \dots, c_k\}$ for classification

classes

for

Statistical learning

we think of the training data $\{(y_n, x_n)\}_{n=1}^N$

as having some (maybe complicated) joint distribution (probabilistic model)

We want to construct some \hat{f} , to do this we first construct a measure of "wrongness" called a Loss function $L: \mathcal{F} \rightarrow \mathbb{R}$

so that

same potential \hat{f}

space of possible f s

$$L(f)$$

how good a job f does
i.e. how close is $f(x_1, \dots, x_p)$
to y

ideally we would let

$$(*) \quad \hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}[L(f)]$$

Risk
= expected loss
(Risk minimization)

we often don't know the joint dist of training data exactly. So instead we let

we only have training data

$$\rightarrow \hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N L(f)$$

(empirical risk minimization)

Linear Regression:

- ① classic method (well-studied)
- ② Simple (this is good)
- ③ can be powerful
- ④ basis for many other methods

We want \hat{f} . Regression assumes

our estimate

$$Y = f(X_1, \dots, X_p) = \beta_0 + \sum_{j=1}^p \beta_j X_j$$

↑
true underlying model

↑
a linear relation
btwn Y and X s

If we let $\underline{X} = (1, X_1, \dots, X_p) \in \mathbb{R}^{p+1}$

If we let $\underline{x} = (1, x_1, \dots, x_p) \in \mathbb{R}^{p+1}$

and $\beta = (\beta_0, \beta_1, \dots, \beta_p) \in \mathbb{R}^{p+1}$

then we can simply write \leftarrow the β s assoc. w/ f

$$y = \underline{x}^T \beta$$

We say regression is linear b/c f is linear in β s (linear = inner product)

LR: makes simplifying assumption. Namely: linearity

\mathcal{F} could be infinite dimensional (intractable)

\Downarrow LR

$$\mathcal{F} \approx \mathbb{R}^{p+1}$$

(we only need to learn $p+1$ β s to determine \hat{f} or \hat{f})

To determine \hat{f} we need $\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p$

so we form a loss function

$$L(f) = L(\beta) = (y - (\beta_0 + \sum_{j=1}^p \beta_j x_j))^2$$

have good $\hat{y} \approx y$

and then

have $\hat{f}(x) \approx y$ and then

$$\hat{\beta} = \underset{\beta}{\operatorname{argmin}} L(\beta)$$

Subsequently, $\hat{f}(x) = \tilde{X}^T \hat{\beta}$

Practically: empirical risk to minimize

$$L(\beta) = \sum_{n=1}^N (y_n - \beta_0 - \sum_{j=1}^p \beta_j X_{nj})^2$$

If we let X be our "design" matrix

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} & \dots & X_{1p} \\ \vdots & X_{21} & X_{22} & \dots & X_{2p} \\ \vdots & & & & \\ 1 & X_{N1} & X_{N2} & \dots & X_{Np} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

and $y = (y_1, \dots, y_N) \in \mathbb{R}^N$

then our empirical risk to minimize is

$$L(\beta) = \|y - X\beta\|^2$$

note this takes the correct form.

turns out (calc III from hell)

$$\hat{\beta} = (V^T V)^{-1} V^T u$$

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

projector of y onto $\text{Col}(X)$.

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