

# Statistical Learning and Linear Regression

Almost any attempt to clearly define a difference between statistics and machine learning - it's wrong.

Classically: Statistics (statistical inference)  
formulating probabilistic models about variables in order to infer properties/params

machine learning: using algos to predict output given input

## Statistical (machine) learning (SML)

build / understand ML approaches using prob. models

---

Broadly there are two categories of SML problems

### (I) Supervised Learning

supervised meaning we have some examples to "train"  
our approach

idea: want to predict  $\underline{Y}$  from  $\underline{X}$   
and we have example pairs  $(Y, X)$

w/in supervised problems two types of problems:

- (1) "regression" = predicting a continuous outcome
- (2) "classification" = predicting a categorical outcome

Ex. regression problems:

$$Y \in \mathbb{R}$$

- predicting stock market performance from economic indicators

$X$

- predicting adult height from childhood height

$$Y \in \mathbb{R}$$

$X$

classification problems:

- predict if individual will default on a loan given credit score

$$Y = \begin{cases} 0 & \text{don't default} \\ 1 & \text{do default} \end{cases}$$

- predict racial identity from genomic data

$$Y = \begin{cases} \text{white} \\ \text{black} \\ \text{green} \end{cases}$$

## II Unsupervised Learning

Don't have a clear prediction problem.

We want to learn/summarize important trends/info/rels about the vars.

In this class: first half = supervised problems  
second half introduce unsupervised problems.

# Mathematical Setup for supervised Learning

We have some var.  $Y$

called: outcome, predicted var, dependent var

the thing we want to predict

We're going to try to predict  $Y$  given some other vars

$$\underline{X} = (X^{(1)}, X^{(2)}, \dots, X^{(p)})$$

$\uparrow p = \# \text{ vars}$

called: predictors, features, covariates,  
independent vars

they are the things we use to predict.

Predict?

Want to come up w/ a fn  $\hat{f}$  so that

$$Y \approx \hat{f}(\underline{X}) = \hat{f}(X^{(1)}, \dots, X^{(p)}) \stackrel{\text{def}}{=} \hat{Y}$$

Course goals:

① Methods: how do we construct  $\hat{f}$ ?

② Evaluation: how do we determine if  $\hat{f}$  is good.

For supervised problems

We assume we have training data to construct  $\hat{f}$  of the form

$$(y_n, \underline{x}_n) \text{ for } n=1, \dots, N$$

← sample size

where  $y_n \in S$  ~  $S = \text{space of possible } y_n\text{'s}$   
 $S = \mathbb{R}$  for regression  
 $S = \{c_1, c_2, \dots, c_K\}$

$$\underline{x}_n = (x_n^{(1)}, \dots, x_n^{(p)})$$

-----SML-----

Assume that these training data are sampled from some joint dist  $p(\underline{x}, y)$  — typically i.i.d.

What is a general way of constructing a  $\hat{f}$ ?

We can construct a measure  $L$  (loss fn) that tells us, for some putative  $f$ , how good  $f$  does at predicting  $Y$  from  $\underline{X}$ .

[Given  $(\underline{X}, Y)$  come from  $p(\underline{X}, y)$ ]

Idea!  $L(Y, f(\underline{x})) = \begin{cases} \text{large if } Y \neq f(\underline{x}) \\ \text{small if } Y \approx f(\underline{x}) \end{cases}$

Examples: regression context

$$L(y, f(x)) = (y - f(x))^2 \quad \text{Squared loss}$$

$$L(y, f(x)) = |y - f(x)| \quad \text{abs. loss}$$

classification context

$$L(y, f(x)) = \begin{cases} 0 & y = f(x) \\ 1 & y \neq f(x) \end{cases} \quad \begin{matrix} 0-1 \\ \text{loss} \end{matrix}$$

either context:  $L(y, f(x)) = -\log p(x, y)$

neg. log-likel. loss

How do we use a loss to build a  $\hat{f}$ ?

We'd like to choose  $\hat{f}$  so that it incurs a small loss.

In particular minimize expected loss!

ideal:  
Risk Minimization

$$f^* = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \mathbb{E}[L(Y, f(X))]$$

$\mathcal{F}$  = hypothesis space

$\mathbb{E}[L(Y, f(X))]$  = Risk = Expected loss

b/c we don't know  $p(x, y)$  to calc  $\mathbb{E}[\dots]$   
we instead approx. using a sample of training data.

Aside:  $E[z] \approx \bar{z} = \frac{1}{N} \sum_{n=1}^N z_i$  when  $z_i$  iid ...

So

$$\hat{f} = \underset{f \in \mathcal{F}}{\operatorname{argmin}} \underbrace{\frac{1}{N} \sum_{n=1}^N L(y_n, f(x_n))}_{\text{Empirical Risk}}$$

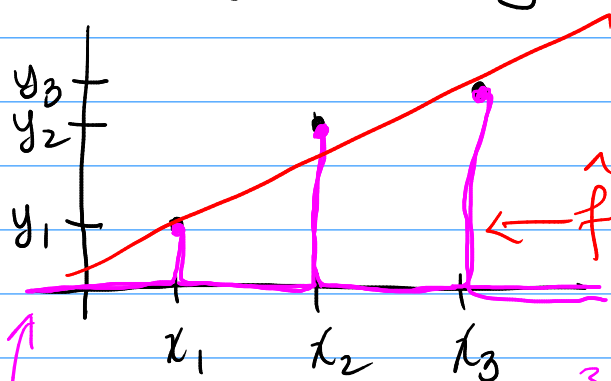
training data

Empirical Risk  $\rightarrow E[L(Y, f(X))]$   
as  $N \rightarrow \infty$ .

ERM: empirical risk minimization

Have a couple choices:  $L$  and  $\mathcal{F}$

One big problem: over-fitting



Consider  $L = \text{sq. err.}$

$\mathcal{F} = \text{all possible functions}$

too large

has  $ER = 0 = \frac{1}{N} \sum_{n=1}^3 (y_n - f(x_n))^2$

this won't generalize well.

Way to avoid this:

related

- ① Restrict  $\mathcal{F}$  (linear, quad, smooth, diff'able)
- ② Change  $L$  (penalization to avoid silly  $f$ )
- ③ get a better estimate of Risk. (test/validation, cross-validation)

# Linear Regression

Why look of LR?

- ① classic method - very well studied
- ② simple (good)
- ③ can be very powerful
- ④ LR is the basis for more complex methods

Properly: linear least-squares regression

$\mathcal{F} = \text{lin. fns}$   $L(y, f(x)) = (y - f(x))^2$   $\rightarrow \text{continuous } y$

$\rightarrow$  OLS: ordinary least-squares (regression)

Notation!  $Y, \underline{X}$  generic/RVs

$\underline{X} = (X^{(1)}, \dots, X^{(p)})$

Samples!  $y_n, \underline{x}_n = (x_n^{(1)}, \dots, x_n^{(p)})$

Setup!  $Y = f(X^{(1)}, \dots, X^{(p)}) = \beta^{(0)} + \sum_{j=1}^p \beta^{(j)} X^{(j)}$  (linear)

If  $\underline{X} = (1, X^{(1)}, X^{(2)}, \dots, X^{(p)}) \in \mathbb{R}^{p+1}$

$\xrightarrow{\text{design}}$  and  $\beta = (\beta^{(0)}, \beta^{(1)}, \dots, \beta^{(p)}) \in \mathbb{R}^{p+1}$

then

$$Y = \underline{X}^T \beta$$

linear in  $\beta$  b/c it's basically inner product

Idea: by assuming a linear form of  $f \in \mathcal{F}$   
so that

$$f(\underline{x}) = \underline{x}^T \beta \rightarrow \beta \in \mathbb{R}^{p+1}$$

we restrict  $\mathcal{F}$  to that of linear fns.

$\mathcal{F}$  (potentially infinite dim'l)

OLS  $\rightsquigarrow \mathcal{F} = \mathbb{R}^{p+1}$

b/c  $f = f_\beta \in \mathcal{F}$

