

Recap lecture

No assumptions.
non-parametric

hardest
Regression

find F
simplest

Pattern
recognition.
next simplest

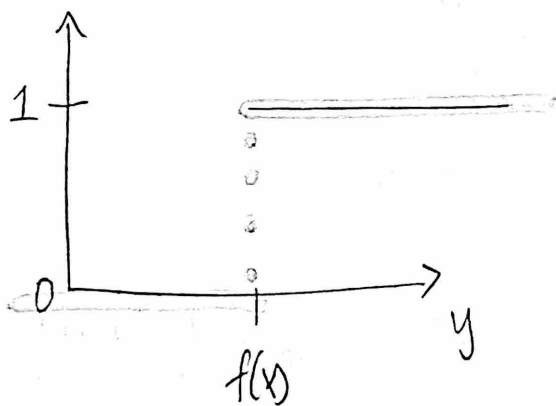
Find f :

$Y = f(x)$ exact relationship.

$$F := \{ F_{y|x}(y|x) = \mathbb{1}_{y=f(x)}(y) \}$$

distribution
func

for a given $f(x)$

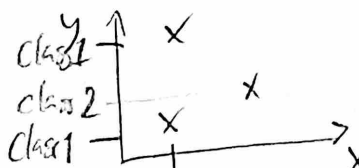


Pattern Recognition

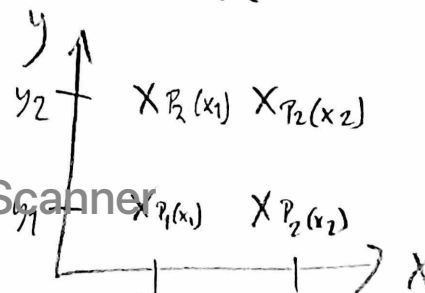
$F := \{ F(y|x) \text{ is discrete, } y_1, \dots, y_k \text{ are the classes} \}$

it has a # data points belonging to classes.

all x belong to the same classes



this would not work.
The x 's need to be labeled
See classes



this would work.
Need data that looks like this.

The prob. depend on the x 's.

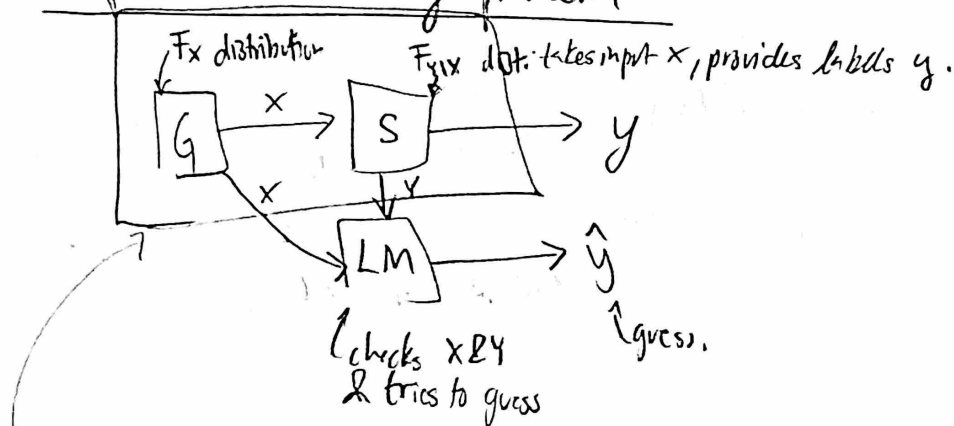
$$P(Y = y_i | X = x) = P_i(x)$$

These are non-parametric because the function $F(x)$ could be any function. And we make no assumptions on how $F(x)$ goes.
So two ways that determine if it is non-parametric.

Regression: see lecture notes.

Assume regression & y-variable has finite 2nd moment.

Supervised learning problem



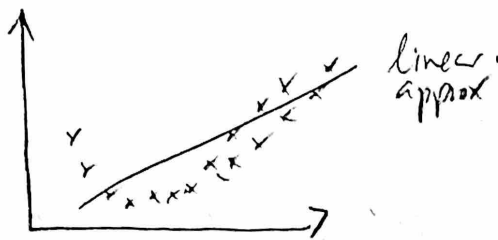
Statistical model
 $F_{xy} \in \mathcal{F}$

Model space (where LM searches)
 \mathcal{M}
LM

Risk (loss) (how LM searches)

Statistical model \neq Model space.

We are trying to find a model space with our approximation. We want to find a good approx to reality. But that doesn't mean we have found the true model space. We have just found an approximative model space.



The ^{approx} linear model space \neq the actual model space.

Teatplogg

- individuella
- group assignments
- lax itds
- problem solving sessions.

Risk

Risk = expected loss. (L)

$$L(g, x, y) \quad g \in \mathcal{M}$$

"an error"

$$R(g) = E[L(g, \underbrace{x, y}_{\text{Random}})]$$

expectation
to these
random
inputs/variables

LM wants to find the g^* which is the function that minimizes the expected loss.

$$g^* = \arg \min_{g \in \mathcal{M}} R(g)$$

$$R(g^*) = \min_{g \in \mathcal{M}} R(g)$$

LM cannot find $R(g^*)$. It cannot do that as it doesn't have access to the true distribution. It only sees the data so instead it minimizes \hat{R}_n . ^{LM} sees $(x_1, y_1), \dots, (x_n, y_n)$

$$\hat{R}_n(g) = \frac{1}{n} \sum_{i=1}^n L(g, x_i, y_i)$$

$$\arg \min_{g \in \mathcal{M}} \hat{R}_n(g) = \hat{g}_n^* \quad \text{hopefully } \hat{g}_n^* \sim g^* \text{ or at least}$$

$$R(\hat{g}_n^*) \text{ is small.}$$

^{the true}
^{n.b. = train-test}
Estimate $R(\hat{g}_n^*)$ using test set.

Estimation

Examples:
empirical variance, mean etc.

Two concepts {
* statistic: a function of the data.
* Estimator: = point estimator.
= "a statistic that estimates" some underlying value.

$$F := \{ F(x; \lambda) \mid \lambda \in \Lambda \}$$

← by lambda

$$\Theta(F) \rightarrow \mathbb{R}$$

parameter map

We use the parameter map to map the non-parametric world.

Def. Parametric = finite number of parameters.

The estimator T

$$T(x_1, \dots, x_n) \sim \Theta(F^*) = \Theta^*$$

the expectation

Ex 1:

$$\Theta(F) = \int x dF(x)$$

Ex 2:

$$\Theta(F) = \inf_{g \in \mathcal{H}} R(g)$$

usually we use these two as estimators.

In Ex 1:

what is a good estimator T ?

$$T = \frac{1}{n} \sum_{i=1}^n x_i$$

Ex 2:

$$T = \min_{g \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(g, x_i, y_i)$$

$$\text{Ex 3: } \Theta(F) = \arg \min_{g \in \mathcal{H}} R(g) = g^*$$

$$(x_1, \dots, x_n) = \arg \min_{g \in \mathcal{H}} \hat{R}_n(g) = \hat{g}_n^*$$

Properties of estimation:

- * bias (unbiased)
- * standard error
- * Mean squared error (MSE)
- * asymptotically unbiased
- * asymptotically consistent

Bias

$$E[T] - \theta^* : \text{bias is diff expected value of} \\ \underbrace{\text{the estimator}}_T - \text{true value}$$

negative bias:

$$\text{Standard error (se): } \sqrt{V[T]}$$

measures the variability of your estimator T .

MSE

$$E[(T - \theta^*)^2] = \text{bias}^2 + \text{se}^2$$

if have BB data: use MSE.

Annotations:
- Just minimize se, bias might go up.
- Just minimize bias, se might go up.

Asymptotically unbiased

as # samples $\rightarrow \infty$, bias $\rightarrow 0$.

Asymptotically consistent

as # samples increases, the estimator T converges to the true parameter.

$$T \xrightarrow[n \rightarrow \infty]{P} \theta^*$$

Use if $MSE \rightarrow 0$.

Empirical dist. function
samples are x_1, \dots, x_n

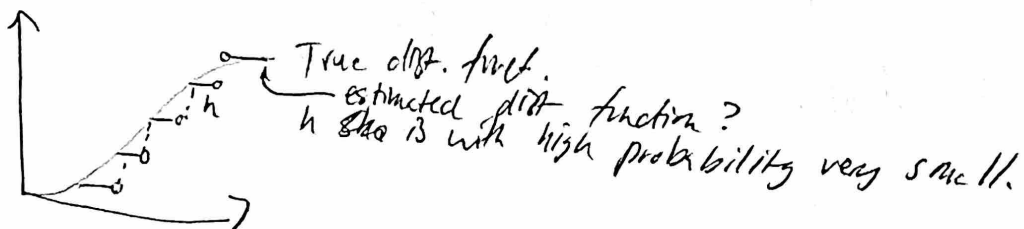
$$\hat{F}_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq x}$$

$$E[\underbrace{\hat{F}_n(x)}_T] = \frac{1}{n} \sum_{i=1}^n E[\underbrace{\mathbb{1}_{x_i \leq x}}_{\substack{TP(x_i \leq x) \\ \text{dist. func} \\ 'F(x)'}}] = F(x) \quad \leftarrow \text{This is our } \theta^*$$

DKW-inequality

= The thing that corresponds to the Hoeffding bound for empirical dist. functions.

$$= TP(|\hat{F}_n(x) - F(x)| > \epsilon) \leq 2ne^{-2n\epsilon^2}$$



This actually gives confidence intervals for the quantiles.