

Regression: What is the best \hat{f} ?

Statistical Learning: (X, Y) have some joint distribution
and $Y = f^*(X)$ $p(x, y)$
↑ unknown

⊗ Want \hat{f} that is an estimate of f^*
so that $Y \approx \hat{f}(X)$

⊗ we construct \hat{f} by defining a Loss $L: \mathbb{R}^P \rightarrow \mathbb{R}$

$$\hat{f} = \underset{f}{\operatorname{argmin}} \overbrace{E[L(f)]}^{\text{Risk}}$$

so far $L(f) = (Y - f(X))^2$

Reality: training
 $\{(x_n, y_n)\}$

minimize
empirical risk

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

MSE

What is theoretically \hat{f} ? (for sq. loss)

Quantity of interest:

$$E[L(f)] = E[(Y - f(X))^2]$$

$$= \iint (y - f(x))^2 p(x, y) dx dy$$

↑ density fn

Iterated Exp.

$$E[X] = E_X[E_{Y|X}[X]]$$

$$E[XY] = E_X[E_{Y|X}[XY]]$$

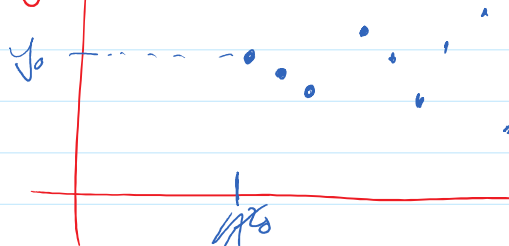
$$= \iint (y - f(x))^2 p(y|x) p(x) dx dy$$

density f_n

$$(*) = \int \left[\int (y - f(x))^2 p(y|x) dy \right] p(x) dx = \mathbb{E}_x \left[\mathbb{E}_{y|x} [(Y - f(x))^2] \right]$$

$$\hat{f} = \underset{f}{\operatorname{argmin}} \mathbb{E}[L(f)]$$

picking a fn $f: \mathbb{R}^p \rightarrow \mathbb{R}$
 $y = f(x)$ build a fn



building \hat{f} can be
done one x at a time
 $\{f(x)\}_x$

Aside

just a large table

x	$f(x)$
1	3
1.1	4
1.01	7.3
π	e

$$(*) \int \left[\int (y - f(x))^2 p(y|x) dy \right] p(x) dx = \int A(x) p(x) dx$$

$A(x) = \mathbb{E}[(Y - f(x))^2 | X=x]$

depends on $f(x) \in \mathbb{R}$

want to minimize
(chOOSE f to minimize this)

① I can choose $f(x)$
independently for each x

(Chosen minimize $T(x)$)

(2) So $A(x)$ are independently set (accady to chosen $f(x)$) at each x

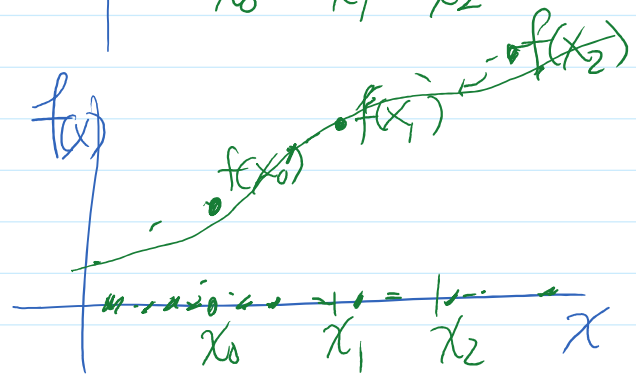
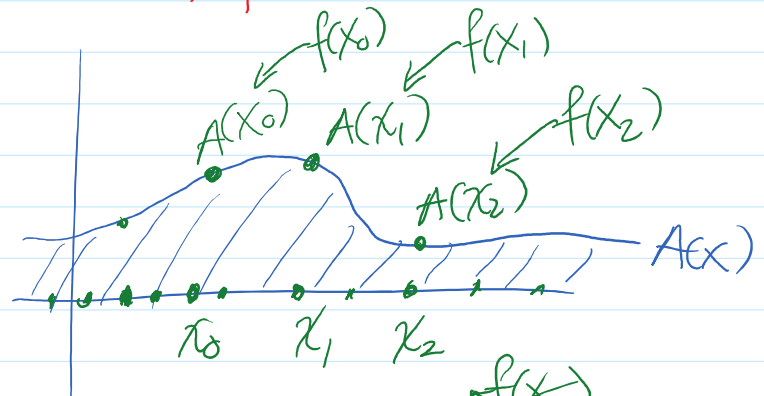
Want to minimize

$$\int A(x)p(x)dx$$

$p(x)$ given can't choose
choice of $f(x)$

Way to minimize is to minimize $A(x)$ independently at each x

by choosing $f(x)$ to minimize $A(x)$



(*) $\hat{f}(x) = \arg \min_{f(x)} A(x) = \arg \min_{f(x)} E[(Y - f(x))^2 | X=x]$

$A(x)$

Fact: $\hat{a} = \arg \min_a E[(z-a)^2]$ ply in $E[z]$

$$\Rightarrow E[(z-a)^2] = E[z^2 - 2az + a^2]$$

$$= E[z^2] - 2aE[z] + a^2$$

$$\frac{\partial}{\partial a} [\dots] = -2E[z] + 2a = 0$$

$$\boxed{\hat{a} = E[z]}$$

We said

$$\hat{f}(x) = \arg \min_{f(x)} E[(Y - f(x))^2 | X=x]$$

$$\boxed{\hat{f}(x) = E[Y | X=x]} \quad (*)$$

Fun fact: $L(f) = |Y - f(x)|$

$$\hat{f}(x) = \text{Median}(Y | X=x)$$

How does this relate?

$$\hat{f}(x) = E[Y|X=x]$$

In practice:
we approximate.

How? (1) Collect some data $\{(x_n, y_n)\}$

(2) Approx. $E[Y|X=x]$

How?

Way 1: $\hat{f}(x) = \text{avg value of } y_n \text{ for } x_n \text{ near } x$
(let's use K nearest neighbors to x)

$$= \frac{1}{K} \sum_{n: x_n \in N_K(x)} y_n$$

(KNN)

Way 2: Assume some ^{linear} structure of $E[Y|X=x] = x^T \beta$

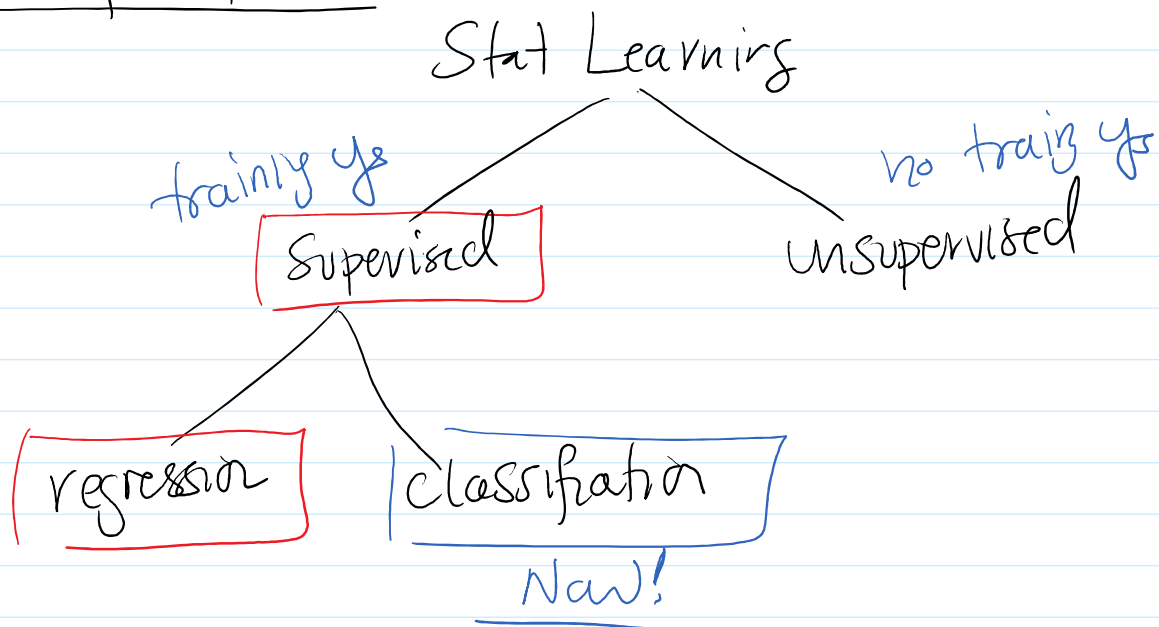
I estimate β s as $\hat{\beta}$

$$\hat{f}(x) = \hat{\beta}^T x$$

(Linear regression)

Classification:

Stat / learning



Classification: (Categorical y)

- predict song title from sounds multi-class
- predict tumor benign/malignant 2-class
- predict protein class:
 - ① α -helix ② β -helix ③ random coil
 - 3-class

Setup: $\{(y_n, x_n)\} : x_n \in \mathbb{R}^p ; y_n \in \mathcal{C}$

$\mathcal{C} = \{c_1, c_2, c_3, \dots, c_K\}$

↗ set of K classes

Goal: find \hat{f} so that $\hat{y} = \hat{f}(x) \approx y$

Optimal \hat{f} ?

Need a loss for classification: "0-1 loss"

$$L(f) = \mathbb{I}(f(x) \neq Y) = \begin{cases} 0 & f(x) = Y \\ 1 & f(x) \neq Y \end{cases}$$

$$\hat{f}(x) = \arg \min_{f(x)} \overbrace{\mathbb{E}[\mathbb{I}(Y \neq f(x) | X=x)]}^{\text{loss}}$$

Aside: $\mathbb{E}[\mathbb{I}(x \in A)] = \int \mathbb{I}(x \in A) p(x) dx$
 $= \int_A p(x) dx = P(X \in A)$

$$= \arg \min_{f(x)} P(Y \neq f(x) | X=x)$$

$$= \arg \min_{f(x)} 1 - P(Y = f(x) | X=x)$$

maximize this

$$\boxed{\hat{f}(x) = \arg \max_{c \in \mathcal{C}} P(Y = c | X=x)}$$

Bayes classifier

Idea: see! π

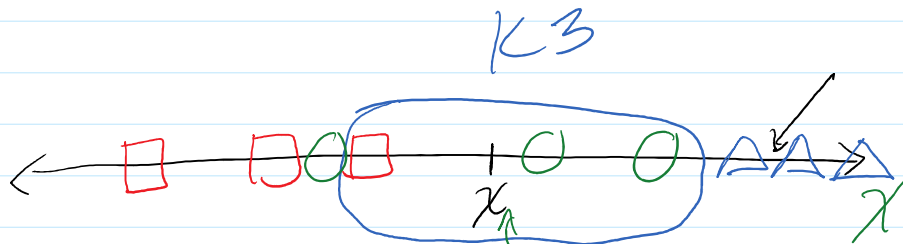
class probs

$$P(Y=c_1 | X=x), P(Y=c_2 | X=x), P(Y=c_3 | X=x), \dots$$

↑ pick largest.

KNN classification

classes: □ ○ △



$\hat{f}(x)$ = majority class among K nearest neighbors

$$\hat{f}(x) = 0$$

$$P_c = P(Y=c | X=x) \approx \frac{1}{K} \sum_{n: x_n \in N_K(x)} \mathbb{1}(y_n = c)$$

= % of K nearest neighbors of class c

↑ pick $\hat{f}(x)$ = class w/ largest.