# Overview of Supervised Learning.

## 1
– Quick Gothrough. : Terminology & Variable Types.

$$\text{Inputs} \xrightarrow{\text{predict}} \text{Outputs}$$

Stat        predictors                 independent vars.

Pattern         features                     responses
Recognition.

–

Variable types :   quantitative measurements ,   qualitative vars. ,   ordered categorical vars.

          categorical / discrete.
                    ↑

Output ·                 Regression              classification.        Chap 4  ?

– Two Simple Approaches to prediction.

    Least Squares & Nearest Neighbors.

                                    Training set $T$.

\* Linear models and least squares                    low variance  Decision boundary  smooth

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^{p} X_j \hat{\beta}_j$$       include $1$ in $X$.       $\hat{Y} = X^T \hat{\beta}$   Rely heavily on assumption.
                                    $\Rightarrow$                          potentially high bias.
        ↑                       $\hat{\beta}_0$ in $\hat{\beta}$
    intercept /
    bias (in ML).

Least squares : pick the coefficients $\beta$ to minimize the " residual sum of squares"

$$RSS(\beta) = \sum_{i=1}^{N} (y_i - x_i^T \beta)^2 .$$    –   quadratic function of the parameters.

$$= (\vec{y} - \vec{X}\beta)^T (\vec{y} - \vec{X}\beta)$$   $\xrightarrow[\text{w.r.t. } \beta]{\text{Differentiate}}$   $\vec{X}^T(\vec{y} - \vec{X}\beta) = 0$
            ↓                    ↑
        N-vector              N×p.                        $\Downarrow$  $\vec{X}^T\vec{X}$ is non-singular·

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

                        Adaptive (.no stringent assumptions ).   low bias
                        Decision boundary wiggly & unstable.      high variance.

\*. Nearest – Neighbor Methods.         $N_k(x)$ – neighbourhood of $x$ defined by $k$ closest
                                    points $x_i$ in $T$.
KNN :   $\hat{Y}(x) = \frac{1}{k} \sum_{x_i \in N_k(x)} y_i$ .

    – Effective # of para para : $N/k$.

A large subset of the most popular techniques are variants of the two simple procedures

E.g.
- Kenerl methods : use weights. decrease smoothly to zero, with distance.
  v.s. "0 / 1" weights in KNN.
  - In high-dimensional spaces. the distance kernels are modified to emphasize some vars. more than others.
- Local regression
- Linear models fit to a basis expansion.
- Projection pursuit & Neural network. → consist of sum of non-linearly transformed models.

2. Statistical Decision Theory.

Let $X \in R^p$. $Y \in R$. with joint distribution $Pr(X, Y)$.

We seek a function $f(x)$ for predicting $Y$ given values. of $X$.

⤷ require a loss function. $L(Y, f(x))$. for penalizing errors in prediction.

$$L(Y, f(x)) = (Y - f(x))^2 \quad - \quad \text{squared error loss.}$$

⟹ Expected prediction error (EPE). as a criterion for choosing $f$.

$$EPE(f) = E(Y - f(x))^2.$$
$$= \int [y - f(x)]^2 Pr(dx, dy).$$

$P(X, Y) = P(Y|X) P(X).$

$$= E_X E_{Y|X}([Y - f(x)]^2 | X).$$

suffice to minimize EPE pointwise.

$$f(x) = \arg\min_c E_{Y|X}([Y - c]^2 | X = x).$$

⟹ $f(x) = E(Y | X = x)$ ↝ i.e., the conditional expectation.
aka. the regression function.

*linear regression.*

$$f(x) = x^T \beta.$$

*model-based approach.*
*solve $\beta$ theoretically*
$$\beta = [E(xx^T)]^{-1} E(xY).$$

★ Thus the best prediction of $Y$. at any point $X = x$ is the conditional mean, when best is measured by average squared error.

───

NN methods. $\hat{f}(x) = Ave(y_i | x_i \in N_k(x))$. — Expectations is approx. by averaging over sample data.

under mild regularity. $N, k \to \infty$ s.t. $k/N \to 0$.
$$\hat{f}(x) \to E(Y | X = x).$$
$p \uparrow$. rate of convergence ↓. "curse of dimension"

— conditioning at a point
↓ relaxed to
on some region "close" to target point.

Both KNN and least squares approx. conditional expectations by average.   #3.

Diff :   · Least squares assume $f(x)$ is well approx. by a global linear function.

   · k - NN assumes $f(x)$ is well approx. by a locally constant func.

Other Loss function :   L1 : $E |Y - f(x)|$.  $\Rightarrow$ conditional median  $\hat{f}(x) = median (Y|X=x)$

<span style="color:green">discontinuities in derivatives.   hinder the widespread use.</span>
<span style="color:green">$\hookrightarrow$ squared error is analytically convenient, popular.</span>

Categorical varibles :

   Output :  categorical variable G.
                                         $\rightarrow$ cardinality.
   Loss function:  $L_{K \times K}$ ,   $K = card(G)$.

$$\begin{bmatrix} 0 & L(k,l) \\ & \ddots \\ & & 0 \end{bmatrix}$$   $L(k,l) \geqslant 0$ ,   zero - one loss function. (most often)
                         price paid for classifying an obs. belonging to $G_k$ as $G_l$.

   $EPE = E[L(G, \hat{G}(x))]$.

      $= E_x \sum_{k=1}^{K} L[G_k, \hat{G}(x)] Pr(G_k|x)$

   $\Rightarrow$ It suffices to minimize EPE pointwise.

      $\hat{G}(x) = \underset{g \in G}{argmin} \sum_{k=1}^{K} L(G_k, g) Pr(G_k|X=x)$.        <span style="color:green">Bayes classifier:</span>

         $\underset{0-1 \; loss}{=} \underset{g \in G}{argmin} [1 - Pr(g|X=x)]$. $\Leftrightarrow$    $\hat{G}(x) = G_k$  if  $Pr(G_k|X=x)$
                                                           $= \underset{g \in G}{max} Pr(g|X=x)$.

   * KNN classifier directly. approximate this

3. Statistical Models, Supervised Learning & Function Approximation

   Goal · Find a useful approximation $\hat{f}(x)$ to $f(x)$. that underlies the predictive. relationship
          between the inputs and outputs.

      — Squared error loss  $\rightarrow$ regression function $f(x) = E(Y|X=x)$ for a quantitative
                                                                              response
      — Nearest - neighbor methods  $\rightarrow$ direct estimates of the conditional expectation.
                         but fails in high dimension setting. & when special
                                                           (2.5).
                         structure exists. (2.7)
      — Other classes of models for $f(x)$.

— Discuss a framework for incorporating them into the prediction problem.

△ Statistical Model.

$$Y = f(x) + \varepsilon.$$

$\varepsilon$. random error.   $E(\varepsilon) = 0$.  ind. of $X$.  — Additive error model.

- Input – output pairs $(X, Y)$. not <u>deterministic</u>. $(Y = f(x))$.
  unmeasured variables.  through $\varepsilon$.  ↳. can be handled by techniques.

- Assumption. $\varepsilon$ iid.
  Simple modifications : to avoid independence assumption, e.g. $Var(Y | X = x) = \sigma(x)$

- Quantitative response.
  e.g.
  Other data types → Generalized Linear models.

△ Machine Learning. — Supervised Learning.
  Learn $f$ by example through a " teacher "
  ⇓
  training set.,   learning algorithm.

— Function Approximation.

  $\{x_i, y_i\}$  $(p+1)$– dimensional Euclidean space.   $X_{p \times 1} \xrightarrow{f} Y$
  $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad R^p \qquad\qquad R.$

  - Many approximations.  a set of parameters. $\theta$.
    - Linear models  $f(x) = x^T \beta.$,   $\theta = \beta.$

    - Linear basis expansions.
      $$f_\theta(x) = \sum_{k=1}^{k} h_k(x)\, \theta_k.$$
      $h(\cdot)$.  set of functions.
      e.g. polynomial  expansions
      trigonometric
      sigmoid. transformation.

○ Least squares to estimate the parameters $\theta$ in $f_\theta$. ⟺ maximum
  by minimizing the residual sum-of-squares.   $Pr(Y | X, \theta) = N(f_\theta(x), \sigma^2)$
  $$RSS(\theta) = \sum_{i=1}^{N} (y_i - f_\theta(x_i))^2.$$
  Let's look at it together
  $(p\ 31)$.

○ A more general principle for estimation : maximum likelihood estimation.
  Supp. $y_i$  $i = 1, \cdots, N$ $\sim Pr_\theta(y).$

  $$L(\theta) = \sum_{i=1}^{N} \log Pr_\theta(y_i).$$

− ~~Nearest~~ Nearest − neighboor Methods

( & other local ).

Face problems : <1> High dimensions → " curse of dimensionality ". More 2.5.

Consider a $p$- dimension unit hypercube.

Supp. We send out a hypercubical neiborhood. about a target point to capture fraction $r$ of the observations.

⇒. a fraction $r$ of the unit volume., the expected. edge length

$$e_p (r) = r^{1/p}.$$

$p = 10$ , $r = 0.01$ ~~$e_{10}(1)$~~

$$e_{10} (0.01) = 0.63$$

$r = 0.1$ $$e_{10} (0.1) = 0.80.$$

<2>. More structured approaches can make more efficient use of the data.

$$RSS (f) = \sum_{i=1}^{N} (y_i - f(x_i))^2. \quad any. \ f(·).$$

− minimize this → infinitely many solutions

if. multiple obs. $x_i$, $y_{il}$ $l=1, \cdots N_i$. → $N$ sufficiently large.

− obtain useful results for finite $N$. ⇒. restrictions

e.g. parametric representation.

or build into the learning method.

Classes of Restricted Estimators

＊. Roughness penalty and Bayesian methods.

~~RRS~~ $PRSS (f; \lambda) = RSS(f) + \lambda J(f).$

explicitly penalize $RSS (f)$ with a roughness penalty.

＊. Kernel methods and local regression.

explicitly providing estimates of the regression function or conditional expectation by specifying the nature of the local neiborhood.

＊ Basis function and Dictionary methods.