$\boxed{\text{AI/ML Overview}}$ ①

② General Classes of problems:

①  Supervised (i.e. Predictive) Learning

②  Unsupervised Learning

Supervised: Goal is to learn a $\boxed{\text{mapping}}$
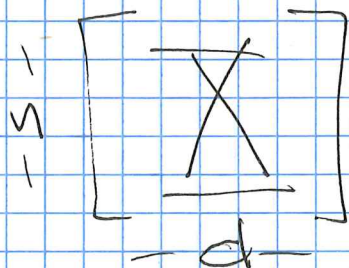
From inputs: $X \xrightarrow{\text{To}} Y$

outputs (labels)

With supervised learning we are given $\boxed{\text{labels}}$.

$$D = \{(x_i, y_i)\}_{i=1}^{n} \quad \text{where} \quad x_i \in \mathbb{R}^d$$

$X$ usually denotes the $\boxed{\text{design matrix}}$

where $X$ $n \times d$.

When $y_i \in \mathbb{R}$, the problem
$\underset{\text{labels}}{\smile}$

archetype is usually referred to as $\boxed{\text{Regression}}$.

When $y_i \in \{1, ..., K\}$ the problems
$\underset{\substack{\text{Discrete} \\ \text{classes}}}{\smile}$

is referred to as $\boxed{\text{Classification}}$.

E.g. Predict expected income from
education level (Regression)

E.g. Predict whether digital image contains
a pedestrian (Binary Classification).

Goal: $Y = f(X) + \varepsilon$

$\underset{\text{"True" mapping}}{\smile}$ $\underset{\text{error}}{\smile}$

Why estimate $f$?

① Prediction
② Inference

$\hat{Y} = \hat{f}(x)$ (Our estimate)

③

Flowchart for a standard AI/ML Algorithm

WANT: $f \lesssim \hat{f}$ ← (our estimate of $f$)

"True" relationship b/w X & Y

How To Quantify The proximity of $f$ To $\hat{f}$?

Use a Loss Function.

Examples of Loss Functions:

① 0-1 Loss (for Classification)

$$L(f(x), \hat{f}(x)) = \begin{cases} 0 & \text{if } f(x) = \hat{f}(x) \\ 1 & \text{if } f(x) \neq \hat{f}(x) \end{cases}$$

② Quadratic Loss

$$L(f(x), \hat{f}(x)) = \left( f(x) - \hat{f}(x) \right)^2$$

(4)

Flow Chart

Partition into Test & Training sets

I Collect Data : $D = \{(x_i, y_i)\}_{i=1}^{n}$

$D_1$ : Trning Set
$D_1 \subset D$

$D_2$ : Test Set
$D_2 \subset D$

$(D_1 \cap D_2 = \emptyset)$

Note!

II Train a Model: $\hat{f}$ using $D_1$
(Regression, NN, SVM)

III Evaluate Model with Loss function, using $D_2$

(*) Big Ideas: The smaller (Total) Loss on The Test set, The better The model (ideally).

We use The results on The Test set to generalize approximate how well The model will generalize to new data.

With [unsupervised learning] we are just given data without labels.

In this case we are interested in discovering "interesting structure" in the data; This is sometimes called [knowledge discovery] or [cluster analysis].

Note: [Reinforcement Learning], offers a 3rd problem class in AI/ML — where an "agent" learns how to act or behave when given occasional reward or punishment signals (e.g. Alpha GO (2016), Atari w/ Deep Q-learning (2014))

[Parametric Models vs. Non-Parametric Models]

Parametric Models consist of a finite (a fixed #) of parameters : $\overline{\Theta} = \langle \Theta_1, \Theta_2, ..., \Theta_N \rangle$.

Idea: Using the training data, we "learn" values for the parameters.

Ex. Polynomial Regression: Fit a polynomial curve to a data set (e.g. using OLS, etc.)

Linear Regression

$$\hat{f}(x) = \theta_0 + \theta_1 x$$

→ "Learning" this model entails finding plausible values for $\theta_0, \theta_1$

Quadratic Regression

$$\hat{f}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$$

Model Parameters: $\vec{\theta} = \langle \theta_0, \theta_1, \theta_2 \rangle$.

Polynomial Regression

$$\hat{f}(x) = \sum_{i=0}^{d} \theta_i x^i \quad (d+1 \text{ parameters})$$

Note: If we use a model with a small number of parameters it is usually easier to train (requires less time & data).

However, a <u>low</u> dimensional Model

might not be sufficiently "complex" to capture

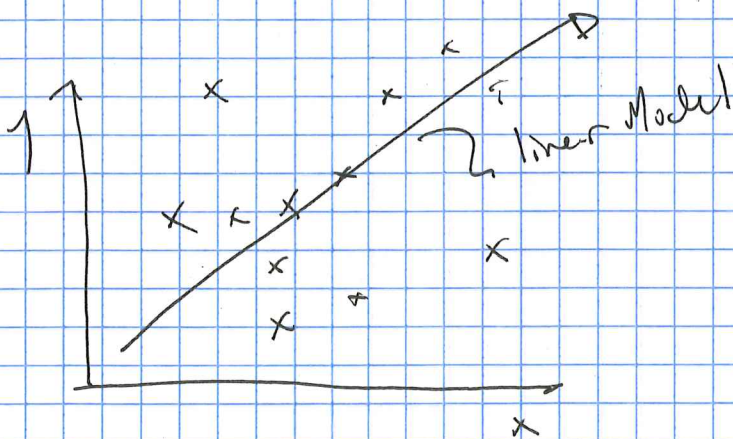all the interesting & useful patterns in our data!

(This is called [underfitting])

On the other hand, a large dimension/complex

model requires more computation & Time on

average; Moreover, an excessively complex model

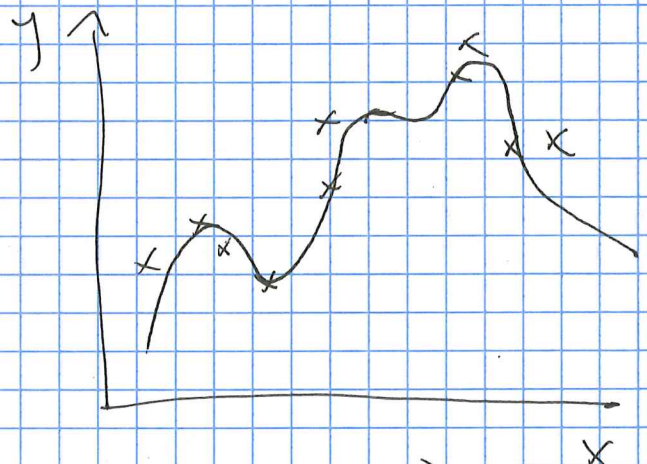will be "overly tuned" to the training data

(This is called [overfitting])

<u>Conclusion:</u> There is No "free lunch" in data science!



(underfitting)

low complexity model

(overfitting)

high complexity model

How do I know when I get it right?

This is the "art" of AI/ML/Data Science!

(In other words there is no simple answer)

In general, however, remember that we can assess our model accuracy with a Loss function:

Quad Loss : $MSE = \frac{1}{n} \sum_{i=1}^{n} \left( y_i - \hat{f}(x_i) \right)^2$

Mean-squared error

over test data

0-1 Loss : $Ave = \frac{1}{n} \sum_{i=1}^{n} L\left( y_i, \hat{f}(x_i) \right)$

Total Average

counts # of "mistakes"

Note: Unfortunately, having a low training error (e.g. MSE) does not guarantee low Test error in general.
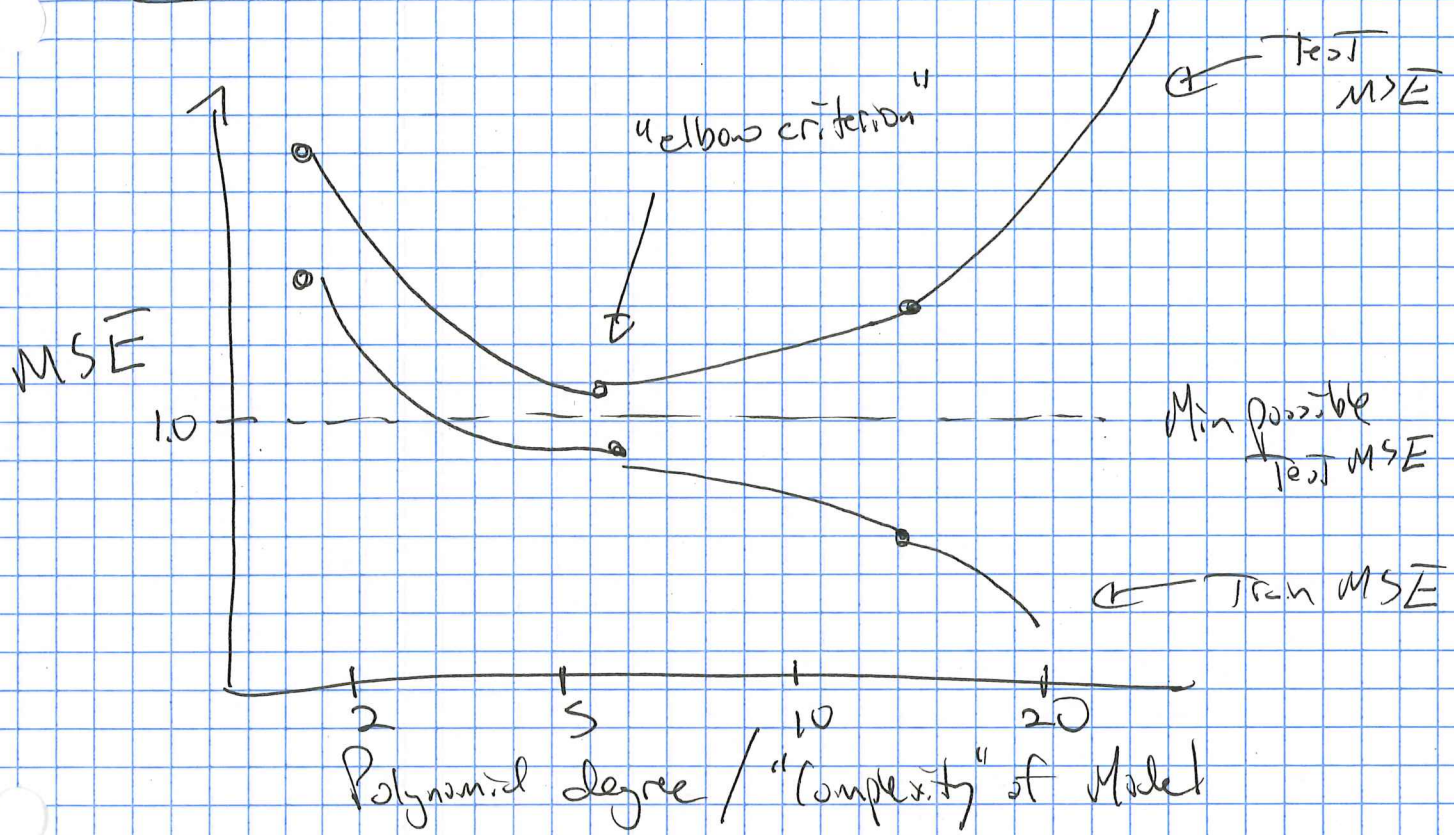
One common remedy for parametric models:

Train several models of varying complexity

(e.g. linear regression, quad, cubic regression), compute MSE for each Test set, choose model w/ lowest MSE.

Common "U-shape" error graph"



"elbow criterion"

← Test MSE

MSE

1.0 - - - - - - - - - - - - - - - - - - -    Min possible Test MSE

← Train MSE

2      5      10      20

Polynomial degree / "Complexity" of Model

---

## Bias-Variance Tradeoff

The "U-shape" phenomenon in the Test MSE is indicative of Two competing properties of learned models: Bias & Variance.

low-dimensional (simple models): high Bias & low variance

high dimens'l (complex/flexible models): low Bias & High Variance

More concretely,

The expected $\overset{\text{Test}}{\text{MSE}}$ for a given value $x_0$,

can always be decomposed into the sum

of ③ fundamental quantities:

$$\underbrace{E\left[(y_0 - \hat{f}(x_0))^2\right]}_{\substack{\text{expected Test} \\ \text{MSE}}} = Var\left[\hat{f}(x_0)\right] + \left[Bias(\hat{f}(x_0))\right]^2$$

$$+ \underbrace{Var(\varepsilon)}_{\substack{(\text{irreducible} \\ \text{error})}}$$

From above, we see that the ideal model will

simultaneously achieve $\boxed{\text{low variance \& low bias}}$.

$\boxed{\text{Ex.}}$ Unsupervised Learning $\overset{(\text{superscript} \atop \text{indicates datum \#})}{}$

Suppose we have $D = \left\{ \left( x_1^{(i)}, x_2^{(i)} \right) \right\}_{i=1}^{u}$

with no class labels                $(\text{2-d data})$
(i.e. no y values).
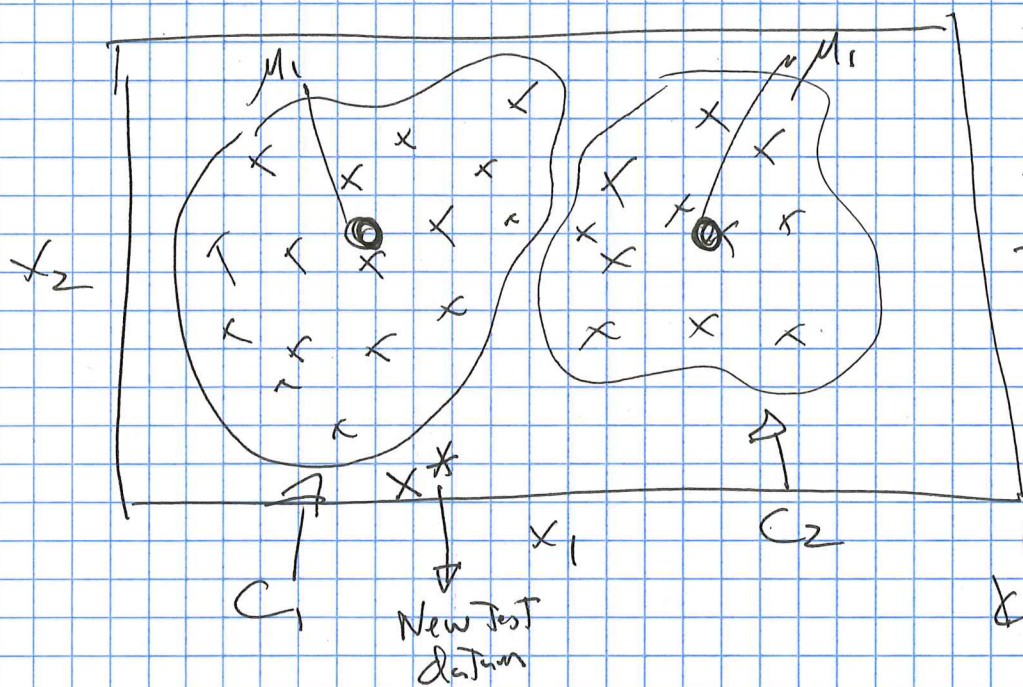
We will use a clustering method to first cluster the data ( let $k = \#$ of clusters) &

Then classify a new datum based on a nearest centroid criterion — This is usually called $\boxed{k\text{-means}}$.

Idea: Suppose we obtain biometric data from men & women $\left( \text{so } k=2 \right)$, where $\vec{x} \in \mathbb{R}^2$.



We use the k-means algorithm to identify each Training datum as belonging to either cluster 1 or cluster 2. (Don't worry about k-means details now)

$$\mu_1 = \frac{1}{|C_1|} \sum_{C_1} (x_1, x_2) \longrightarrow \boxed{\text{centroid of } C1}$$

$$\mu_2 = \frac{1}{|C_2|} \sum_{C_2} (x_*, x_2) \longrightarrow \boxed{\text{centroid of } C_2}$$

How to classify new Test datum $(*)$? $\longrightarrow$ i.e. classify $x^*$ based on nearest cluster center.

Class for $\boxed{x^* = \arg\min_i \| x^* - \mu_i \|_2}$

Note That The previous examples are considered parametric models, because The # of parameters is fixed & once we "learn" These parameter values, The training data can be "discarded" when we perform Test prediction/classification.
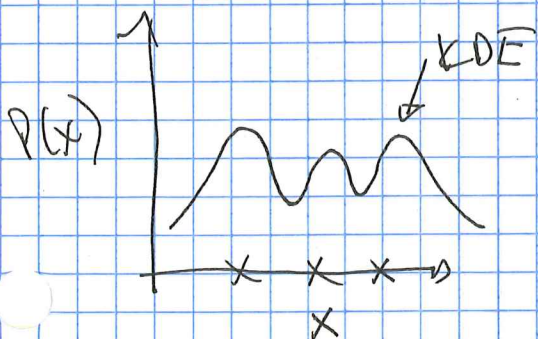
What would a non-parametric model look like?

Two conditions: ① The # of parameters is not fixed, and typically grows as The training set gets larger;

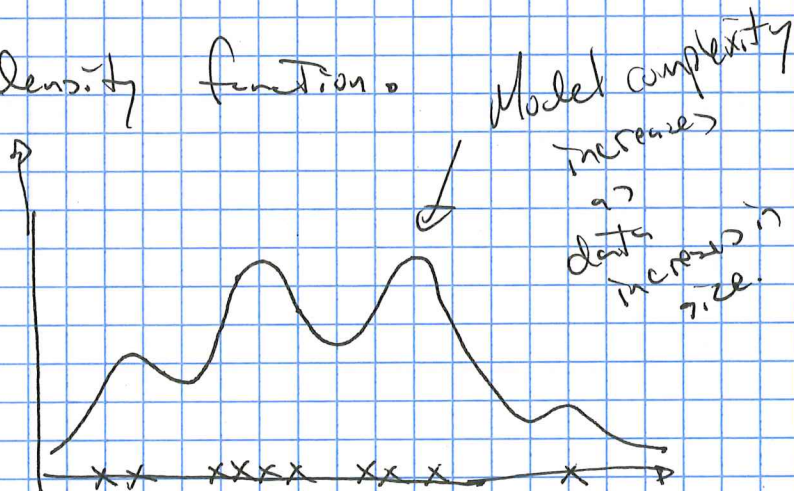② The training data is retained for prediction/classif. (necessarily)

Ex. Kernel Density estimation

Idea: fit a smooth/flexible "histogram" to data - produce a probability density function.

Model complexity increases as data increases in size.

$P(x)$

KDE



$(n=3)$                    $(n=10)$
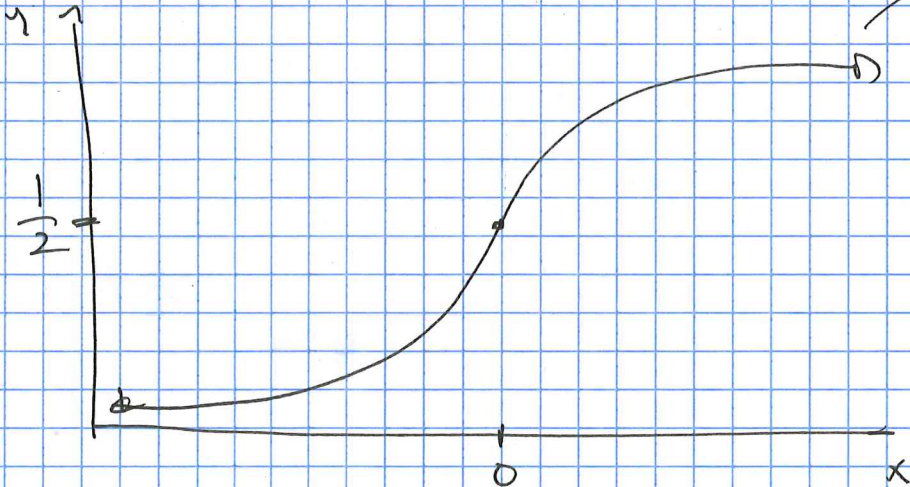
One more useful & common Model

**Ex.**

Logistic Regression : A parametric (Binary)
Classification Model.

**Def.** The sigmoid function :

$$\text{sigm}(n) = \frac{1}{1 + e^{-n}}$$

used for logistic regression & as "activation" function in NNs.

graph of sigmoid



Note: $\lim_{x \to +\infty} \text{sigm}(x) = 1$

$\lim_{x \to -\infty} \text{sigm}(x) = 0$

$\boxed{\text{Idea:}}$ Parameter Models $\vec{\Theta} = \langle \Theta_0, ..., \Theta_N \rangle$

$$P\left(y \mid \vec{x}, \Theta\right) = \underbrace{Ber}_{\text{Bernoulli}}\left(y \mid sigm\left(\Theta^T \vec{x}\right)\right)$$

Binary class

$y \in \{0,1\}$

datum for classification

Model parameters

Steps: ① With Training, linear Model parameters $\vec{\Theta}$.

② evaluate: $sigm\left(\Theta^T \vec{x}\right) = \dfrac{1}{1 + e^{-\Theta^T \vec{x}}}$

③ Apply a $\boxed{\text{decision rule}}$ (Threshold)

e.g. $\boxed{\hat{y}(x) = 1 \quad \longleftrightarrow \quad p(y=1 \mid x) > 0.5}$

$\boxed{\text{Toy example}}$: Training Data

Recall: $y \in \{0,1\}$

↳ Model Fit?

logistic Model

y

x x x

0 x x x

x

Data is one class

Data in zero class

y

logistic Model

Threshold

x will be classified as (+1)

0

x