

RAW DATA VS FEATURE EXTRACTION:

Advantages: No domain specific knowledge required

Disadvantages:

- Highly redundant, large dimensional spaces
- Unknown discriminability

• Attempt to capture discriminant information in the data

• Lower dimensionality and complexity

• Domain specific knowledge required.

MACHINE LEARNING: Set of methods that can automatically detect patterns in data and then use the uncovered patterns to predict future data or to perform other kinds of decision making under uncertainty.

It is used when: { There is a pattern
We cannot pin it down mathematically
We have data on it

TYPES:

- ① PREDICTIVE or SUPERVISED LEARNING: labelled data set, $D = \{(x_i, y_i)\}_{i=1}^N$.
training set, we seek to find a mapping from x to y :
regression, classification...
prediction of quantity \uparrow YES/NO target
- ② DESCRIPTIVE or UNSUPERVISED LEARNING: given a data set $\{x_i\}_{i=1}^N$ we seek to find out more about their structure: density estimation, clustering, dimensionality reduction...
- ③ REINFORCEMENT LEARNING: given an external system upon which we exert control action a and receive percepts p , a reward signal r indicating good performance, find mapping from $P \rightarrow A$ that maximizes some long-term measure of r .

ML 1:

SIMULATE EXPLOTATION DATA:

Given a data set D , we divide it in 3 sets $D = D_{\text{train}} \cup D_{\text{test}} \cup D_{\text{validation}}$

We use D_{train} to learn different models.

We use $D_{\text{validation}}$ to decide which model better adapts the problem

(less accuracy)

We use D_{test} to compute a performance of the selected method.

We apply the selected method to new data and hope it predicts correctly.

ML 2: ABOUT DATA:

HANDLING MISSING DATA:

Do missing data encode any kind of information?

YES → Convert it to a meaningful number or create its own category

NO → Is missing data set small and randomly distributed?

YES → Remove sample with missing data.

NO → Check the histogram: is the distribution of data simple?

YES → Use regressor to infer a plausible value.

NO → Gaussian: use the mean

Non-gaussian: median or randomly draw a sample from the empirical distribution.

INPUTING: process of replacing missing data with another value.

Ways of dealing with missing data:

1) Deletion: Remove data sample with missing values

Pair-wise deletion: keep all possible data in each analysis.

2) Single substitution: Replace with mean/median (INPUTING)

Dummy variable (new variable indicating the value is missing) and impute missing data with mean.

Regression imputation: regression on available data to impute values.

3) Model-based methods: Multiple imputation: by sampling independent values from a distribution.

CATEGORICAL DATA: ONE-HOT ENCODING / DUMMY VARIABLES:

Encode one feature into $k-1$ new features, where k is the amount of values the original feature has.

If we add k new feature we create a metric where all the distances between categories (it is always 1)

More features does not necessarily translate in more accuracy!

NUMERICAL DATA PREPROCESSING: NORMALIZATION TECHNIQUES:

Standardization creates a metric with regular distance computation.

$$z = \frac{x - \mu_x}{\sigma_x}$$

Min Max Scaler:
$$z = \frac{x - \min(x)}{\max(x) - \min(x)}$$

All normalization techniques are influenced by OUTLIERS \Rightarrow It is a good practice to detect outliers and potentially remove them.

observation which deviates so much from the other observations as to arouse suspicion that it was generated by a different mechanism

Approaches:

- 1) Statistical description: infer parameters of the generating distribution and detect points with small probability of belonging to it
- 2) Geometric considerations: convex hulls to detect points on the external boundary of the dataset
- 3) Distance based: distance to neighbors: K-NN.

*) Percentile prc of the data as outlier and normalize according:

$$z = \frac{x - \text{prc}(x, \theta)}{\text{prc}(x, 100 - \theta) - \text{prc}(x, \theta)}$$

θ is the percentile we consider as outlier.

FEATURE HASHING: encoding using a hashing function that given the input category returns the value of the index where the one is located.
(if the category is not in that position we conclude it was not on the list)

Hashing functions must be

- 1) Random: randomly distributes data among all possible index
- 2) Consistent: constant and well defined
when two different samples are assigned to the same position we have a collision.

If data is not uniformly distributed, we must design a function that generates uniformly distributed indexes from data values.

In general, the scheme for hashing data with long character strings, is to break the input into a sequence of small units (bits, words...) and combine all the units $b[1], \dots, b[m]$ sequentially.

When using hash numbers we have to create the hashing function that distributes all the values equiprobably.

Hashing advantages:

- Vectors will usually be very sparse: we can store them efficiently
- Increasing complexity of hashing function may prevent collisions. (bloom filter)
- Don't need to prepare dictionaries or structures ^{online} (real time friendly)
- Distribution of the hashed data tends to be uniform.

Hashing disadvantages:

- Metric ⁱⁿ original space disappears (Locality Sensitive Hashing)
- We have to set in advance the dimensionality of the embedding space.

* The more irrelevant information we add to the data set the worse accuracy we get. We are in fact introducing noise to the model.

ML 3: PERFORMANCE MEASURES:

MODEL SELECTION.

Way of comparing models by simulating the exploitation stage: splitting the data set in TRAINING, VALIDATION and TEST sets.

In order to actually get representative results it is advisable to do this process many times in a way that every point will be used for testing purposes.

CROSS-VALIDATION techniques.

LEAVE ONE OUT (LOO): $\begin{cases} \text{maximizes the amount of data used for training} \\ \text{+ produces good estimation of validation error} \end{cases}$

- i) Take 1 sample of the data set, x_i
- ii) Train the classifier with all the data except x_i : $X_{\text{train}} = \{X \setminus \{x_i\}\}$
- iii) Test the classifier on x_i and store result
- iv) Repeat the process for all samples
- v) Array output with all the results ready for the computation of a performance metric.

K-FOLD:

- i) Split the data set in k disjoint subsets with same cardinality:
 $\{X\} = S_1 \cup \dots \cup S_k$ where $S_i \cap S_j = \emptyset$, $i \neq j$, $|S_i| \sim |S_j| \forall i, j$
- ii) Select one S_i as test set
- iii) Train the classifier on all except S_i : $X_{\text{train}} = \{X\} \setminus S_i$
- iv) Test the trained classifier with S_i and store each result for every sample in the subset
- v) Repeat with each subset
- vi) Array output with all individual results ready for the comp. of p. metric.

OPERATING POINT: precise threshold we select for a specific application by controlling the true positive rate (recall) vs the false positive rate ($1 - \text{specificity}$)

This is a way to control the amount of true and false positives.

Consider a BINARY CLASSIFIER, that stores probability of belonging to class A, i.e. $P(x \in A)$; then obviously $1 - P(x \in A)$ is probab. of belonging to class B.

The decision of belonging to A is given by $P(x \in A) > \text{thr} = 0.5$
 \uparrow
 usually

Imagine we lower the threshold:

$P(x \in A) > \text{thr}$, $\text{thr} < 0.5$, A: positive class

$\Rightarrow \uparrow \text{FP}$

$P(x \in B) < \text{thr} < 0.5$, $\Rightarrow \downarrow \text{FN}$

Recall the confusion matrix and the partial metrics.

Sensitivity / Recall = $\frac{\text{True Positive}}{\text{Real Positives}} = \frac{TP}{TP + FN}$ \nearrow we have more positive predicted, then since the real are fixed $TP \uparrow$

Specificity = $\frac{\text{True Negative}}{\text{Real negatives}} = \frac{TN}{TN + FP}$ \searrow

Precision = $\frac{\text{True Positive}}{\text{Predicted Positive}} = \frac{TP}{TP + FP}$ \searrow , $\text{NPR} = \frac{TN}{TN + FN}$ \nearrow

HL 4: SUPERVISED LEARNING

MACHINE LEARNING ALGORITHM COMPONENTS:

- 1) **MODEL CLASS/HYPOTHESIS SPACE:** family of mathematical models used.
Target decision boundary from this space (linear or non-linear models)
- 2) **PROBLEM MODEL:** formalizes and encodes the desired properties of the solution \Rightarrow in general an optimization problem, minimization of an error function (predicted model - target model)
In classification, the ideal error function is the 0-1 loss
cost function \swarrow
1 when incorrectly classify a training sample
0 otherwise (correct)
- 3) **LEARNING ALGORITHM:** optimization/search method or algorithm that given a model class fits it to the training data according to the error function. (min error or max probable model)

HL 3

RECEIVER OPERATING CHARACTERISTIC: ROC plot of the operating points by varying the threshold value (precision vs recall)

④ MACHINE LEARNING PROBLEM: $\begin{cases} \text{Model class} \\ \text{Cost function} \\ \text{Learning algorithm} \end{cases}$

MODEL CLASS: Linear models $H(w_0, w_1) = \hat{y} = w_1 x + w_0$, w_0, w_1 param.

COST FUNCTION: mean ^{cost} squares: $L(y_i, \hat{y}_i) = \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2$, \hat{y} prediction

The problem of learning is modelized as:

$$\min_{w_0, w_1} \frac{1}{N} \sum_{i=0}^N (y_i - (w_1 x_i + w_0))^2$$

In order to solve the optimization problem let's consider the Gradient Descent algorithm,

- i) $w^0 = (w_0^0, w_1^0)$
- ii) $w^t = w^{t-1} + \eta \Delta w$ with $\begin{cases} \eta: \text{learning rate} \\ \Delta w = - \nabla_w L \end{cases}$

Matricial form: Extended data $\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$, $\hat{y} = \sum_{i=1}^d w_i x^i + w_0 = \tilde{x}^T w$
 where, $x = (x_1, \dots, x_n)^T$, $N = \#$ data points.

Loss matrix: $L(y_i, \hat{y}_i) = \frac{1}{N} \sum_{i=0}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} (\tilde{y} - \hat{y})^T (y - \hat{y})$

Then, the opt. problem becomes: $\min_w \frac{1}{N} (y - \tilde{x}^T w)^T (y - \tilde{x}^T w)$

Note $\nabla_w L = -\frac{2}{N} \tilde{x} (y - \tilde{x}^T w)$

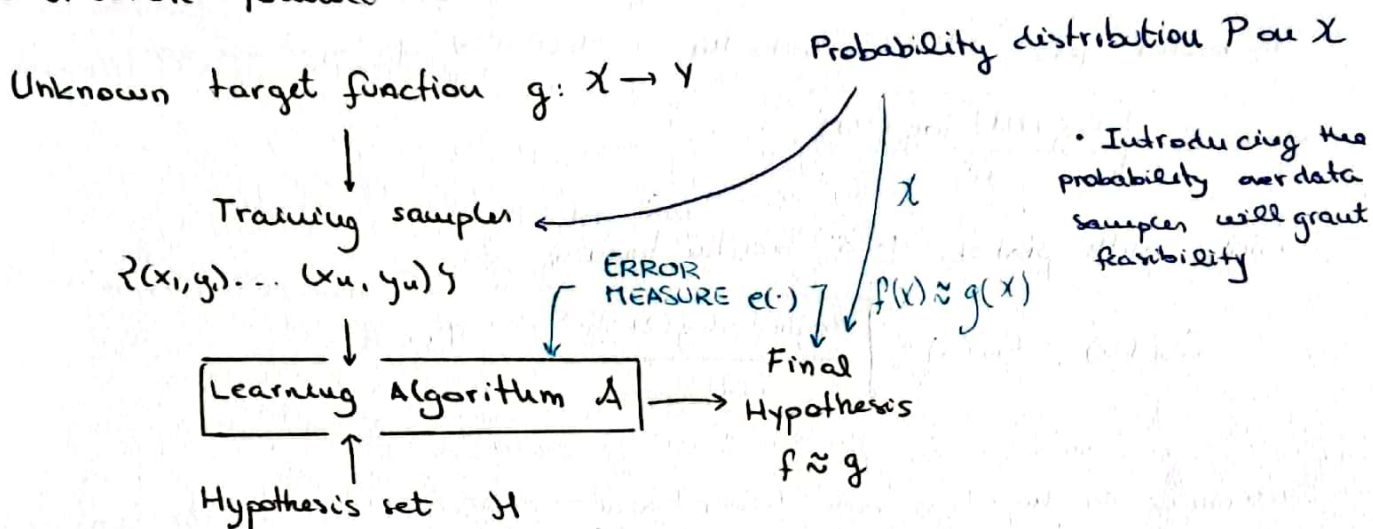
LEAST SQUARES LINEAR REGRESSION:

we have an extreme point $\Leftrightarrow \nabla_w L = 0 \Leftrightarrow w = (\tilde{X} \tilde{X}^T)^{-1} \tilde{X} y$.

ML: 5 FEASIBILITY OF THE LEARNING PROBLEM

A Machine learning problem consists of $\begin{cases} \text{Model class} \\ \text{Cost function} \\ \text{Learning algorithm} \end{cases}$

The structure followed is:



Given a model with an unknown function of probability, when training a learning model it would be possible that the training set represents an improbable situation, misleading the training process. However, although any configuration is possible, not all of them are equally probable. Then, in a big sample, large N , $\hat{\mu}$ (the supposed probability distribution) is probably close to μ (the real one) within ϵ :

$$\text{HOEFFDING INEQUALITY: } P[|\hat{\mu} - \mu| > \epsilon] \leq 2 \cdot e^{-2\epsilon^2 N/2}$$

$\hat{\mu}$: in sample error = freq. of hypothesis getting it wrong: $E_{in}(h)$

μ : out of sample error = expected error: $E_{out}(h)$

$$P(|E_{out} - E_{in}| > \epsilon) \leq 2e^{-\frac{N\epsilon^2}{2}}$$

LEARNING: process of assuming that the in sample error is an indicator of the out of sample error.

We choose the hypothesis with the smaller out of sample error.

Hoeffding doesn't apply to multiple bias!

Using the union bound $P(A_1 \cup \dots \cup A_n) = P(A_1) + \dots + P(A_n)$ we get:

$$P\left[\max_{h \in H} |E_{in}(h) - E_{out}(h)| > \epsilon\right] \leq \sum_{h \in H} P[|E_{in}(h) - E_{out}(h)| > \epsilon] \leq \underbrace{2|H|e^{-2\epsilon^2 N}}_{\delta} \quad (*)$$

The bound gets WORSE with the COMPLEXITY of hypothesis space.

$$\Rightarrow \epsilon = \sqrt{\frac{\log |H| + \log(2/\delta)}{2N}}$$

Then, with probab. $1 - \delta$, over the training set:

$$E_{out}(h) \leq E_{in}(h) + \sqrt{\frac{\log |H| + \log(2/\delta)}{2N}}, \quad \forall h \in H$$

We can use the hypothesis with lower bound \Rightarrow MODEL SELECTION.

We can get N : number of training set if we seek for ex. δ , and generalized error G_E :

$$N = \frac{\log |H| + \log(2/\delta)}{2(G_E)^2}$$

average of point wise errors \Rightarrow

• ERROR MEASURES on the final Hypothesis: $h \approx g$: $E(h, g)$

Different point wise definitions: Squared error: $e(h(x), g(x)) = (h(x) - g(x))^2$

Binary error: $e(h(x), g(x)) = \mathbb{I}[h(x) \neq g(x)]$

$$\square \text{ In sample error: } E_{in}(h) = \frac{1}{N} \sum_{i=1}^N e(h(x_i), g(x_i))$$

$$\text{Out of sample error: } E_{out}(h) = E_x[e(h(x), g(x))]$$

The error measure should be specified by the user/client.

ERROR CONCEPTS: CONFUSION MATRIX

| | | |
|----------|--------------------------------------|---|
| | Gold Standard | |
| | Positive | Negative |
| Positive | TP | FP \rightarrow Precision |
| Negative | FN | TN \rightarrow Negative Predictive Value. |
| | \downarrow Sensitivity (Recall) | \downarrow Specificity |

Metrics:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Partial metrics:

$$\text{Column-wise: Sensitivity} = \frac{TP}{\text{Real Positives}} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \frac{TN}{\text{Real Negatives}} = \frac{TN}{TN + FP}$$

$$\text{Row-wise: Precision} = \frac{TP}{\text{Predicted Positives}} = \frac{TP}{TP + FP}$$

$$\text{NPV} = \frac{TN}{\text{Predicted Negatives}} = \frac{TN}{TN + FN}$$

Another metrics:

$$\text{F1-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = 2 \cdot \frac{\frac{TP}{TP+FP} \cdot \frac{TP}{TP+FN}}{\frac{TP}{TP+FP} + \frac{TP}{TP+FN}} = \frac{2TP}{TP+FN+TP+FP}$$

$$\text{TPR} = \text{recall} = \frac{TP}{TP+FN}, \quad \text{FPR} = 1 - \text{specificity} = \frac{FP}{TN+FP}$$

Then, for a perfect operational/denification we have $FN = FP = 0 \Rightarrow \begin{cases} \text{TPR} = 1 \\ \text{FPR} = 0 \end{cases}$

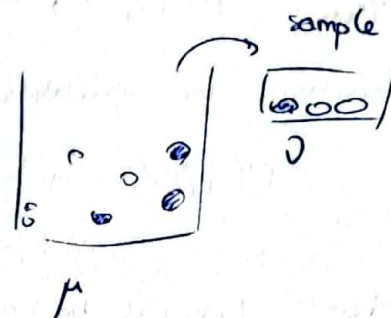
Example considered:

Learning: unknown function $g: X \rightarrow Y$

Each marble \bullet is a point $x \in X$

• Hypothesis right: $h(x) = g(x)$

• " wrong: $h(x) \neq g(x)$



D : in sample error = frequency of the hypothesis getting it wrong: $E_m(h)$
 μ : out of sample error = expected error: $E_{out}(h)$

$$\text{Then } P(|E_{out} - E_m| > \epsilon) \leq 2e^{-\frac{N\epsilon^2}{2}}$$

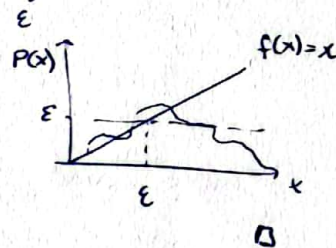
We seek to bound the training error, $P(|\frac{1}{N} \sum_{i=1}^N x_i| > \epsilon) \leq 2\epsilon e^{-\frac{N\epsilon^2}{2}}$

Markov inequality: x positive random variable, then $P(x \geq \epsilon) \leq \frac{E(x)}{\epsilon}$

Proof:

$$E(x) = \int_0^{\infty} x P(x) dx = \int_0^{\epsilon} x P(x) dx + \int_{\epsilon}^{\infty} x P(x) dx \geq \int_{\epsilon}^{\infty} x P(x) dx \geq \int_{\epsilon}^{\infty} \epsilon P(x) dx$$

$$\Rightarrow E(x) \geq \epsilon \int_{\epsilon}^{\infty} P(x) dx = \epsilon P(x \geq \epsilon) \Rightarrow P(x \geq \epsilon) \leq \frac{E(x)}{\epsilon}$$



Chebyshev inequality: $P(|x - \mu| \geq \epsilon) \leq \frac{\text{Var}(x)}{\epsilon^2}$, $\mu = E(x)$

Proof:

$$|x - \mu| \geq \epsilon \Leftrightarrow (x - \mu)^2 \geq \epsilon^2 \xrightarrow{\text{Markov}} P((x - \mu)^2 \geq \epsilon^2) \leq \frac{E[(x - \mu)^2]}{\epsilon^2} = \frac{\text{Var}(x)}{\epsilon^2}$$

$$\text{Then, } P(|x - \mu| \geq \epsilon) \leq \frac{\text{Var}(x)}{\epsilon^2}$$

□

Chernof bound: $P(|x - \mu| \geq \epsilon) \leq \min_{\lambda} E(e^{\lambda|x - \mu|}) e^{-\lambda\epsilon}$, $\lambda > 0$.

Proof:

$$\text{Since } \lambda > 0, |x - \mu| \geq \epsilon \Leftrightarrow \lambda|x - \mu| \geq \lambda\epsilon \Rightarrow e^{\lambda|x - \mu|} \geq e^{\lambda\epsilon}$$

Then if it holds $\forall \lambda$, it holds for the minimum. Applying Markov:

$$P(e^{\lambda|x - \mu|} \geq e^{\lambda\epsilon}) \leq \frac{E[e^{\lambda|x - \mu|}]}{e^{\lambda\epsilon}} = E[e^{\lambda|x - \mu|}] e^{-\lambda\epsilon} \Rightarrow$$

$$P(|x - \mu| \geq \epsilon) \leq \min_{\lambda} E[e^{\lambda|x - \mu|}] e^{-\lambda\epsilon}, \text{ since } \lambda > 0.$$

□

Then, $P\left(\left|\frac{1}{N} \sum_{i=1}^N x_i - \mu\right| \geq \varepsilon\right) \leq 2e^{-N\varepsilon^2/2}$

with some variable change, this is equivalent to

$$P\left(\left|\frac{1}{N} \sum x_i\right| \geq \varepsilon\right) \leq 2e^{-N\varepsilon^2/2}$$

Proof: Let x_i be a Rademacher random variable, $x_i = \begin{cases} 1 & P(x_i=1) = 1/2 \\ -1 & P(x_i=-1) = 1/2 \end{cases}$

Chernof bound: $P(|x| \geq \varepsilon) \leq \min_{\lambda} \mathbb{E}[e^{\lambda x}] e^{-\lambda \varepsilon}$

Similarly, $P\left(\sum_{i=1}^N x_i \geq N\varepsilon\right) \leq \min_{\lambda} \underbrace{\mathbb{E}[e^{\lambda \sum x_i}]}_{\substack{\text{Xi random variable} \\ \text{independent}}} e^{-\lambda N\varepsilon}$

$$\mathbb{E}[\prod e^{\lambda x_i}] = \prod \mathbb{E}[e^{\lambda x_i}]$$

But for any i , we have $\mathbb{E}[e^{\lambda x_i}] = \frac{1}{2} e^{\lambda} + \frac{1}{2} e^{-\lambda} = \sum_i \frac{\lambda^{2i}}{(2i)!} \leq \begin{cases} (2i)! = 1 \cdot \dots \cdot i \cdot i \cdot \dots \cdot 2i \\ i! \cdot 2^i \\ \geq i! \cdot 2^i \end{cases}$

$$\leq \sum_i \frac{\lambda^{2i}}{i! 2^i} = \sum_i \frac{(\lambda^2/2)^i}{i!} = e^{\lambda^2/2}$$

since $\begin{cases} e^{\lambda} = \sum \frac{\lambda^i}{i!} \\ e^{-\lambda} = \sum \frac{(-\lambda)^i}{i!} \end{cases}$

Then, $P\left(\frac{1}{N} \sum_{i=1}^N x_i \geq \varepsilon\right) \leq \prod_{i=1}^N \underbrace{\mathbb{E}[e^{\lambda x_i}]}_{\leq e^{\lambda^2/2} \text{ which does not depend on } x_i} e^{-\lambda N\varepsilon}$

$$= e^{N\lambda^2/2} e^{-\lambda N\varepsilon} = e^{N\lambda^2/2 - \lambda N\varepsilon}$$

Note that $\frac{\partial e^{(N\lambda^2/2 - \lambda N\varepsilon)}}{\partial \lambda} = e^{(N\lambda^2/2 - \lambda N\varepsilon)} [2\lambda N/2 - \varepsilon N] = 0 \Leftrightarrow \lambda = \varepsilon$

Therefore, the minimum for $\lambda > 0$ is given by $\lambda = \varepsilon$, then the Chernof bound is given by:

$$P\left(\frac{1}{N} \sum_{i=1}^N x_i \geq \varepsilon\right) \leq e^{-N\varepsilon^2/2}$$

I did 2 dou nrt (?)

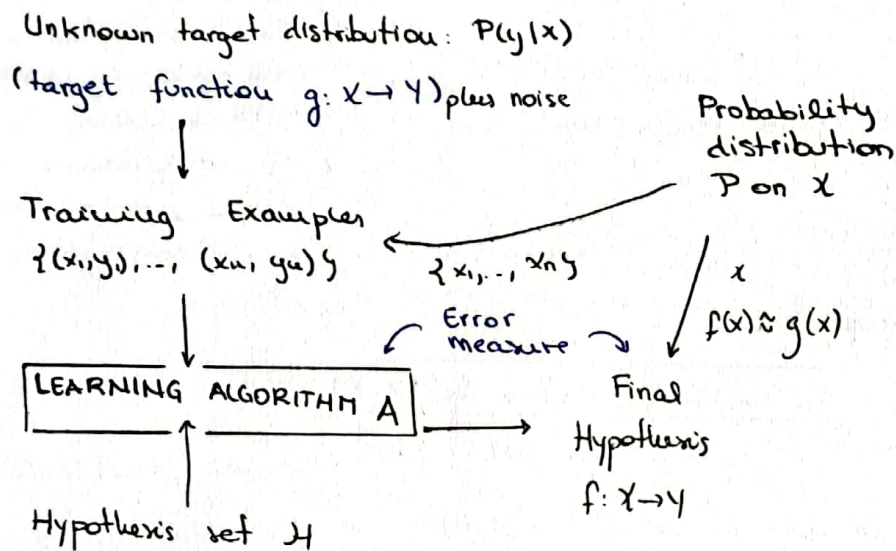
□

Sometimes the "target function" is actually a TARGET DISTRIBUTION: $P(y|x)$

- x still generated by $P(x)$ [we can use Hoeffding]
- y not deterministic with $x \Rightarrow y$ generated by $P(y|x)$.

Thus, (x, y) generated by the joint distribution $P(x) \cdot P(y|x) = P(x, y)$

We could consider $g(x) = E(y|x)$ with noise $y - g(x)$:



- LEARNING is FEASIBLE: Hoeffdings $\Rightarrow E_{\text{out}}(f) \approx E_{\text{in}}(f)$
We also need $E_{\text{out}}(f) \approx 0$, i.e. $f \approx g$

We say that a hypothesis set H SHATTERS a set of k points if at least one of the individual hypothesis in the set is able to generate all possible labels over some configuration of k number points.

VAPNIK - CHERVONENKIS DIMENSION: maximum number of points that a classifier can shatter.

Linear model: can shatter until 3 points but not 4:
(2D)

$$\dim_{VC}(L.M.) = 3$$



Circle classifier: can shatter until 4 points: $\dim_{VC}(C.M.) = 4$
(2D)

3D Linear classifier: $d_{VC}(\text{linear model } D) = D+1 \Rightarrow$
 $d_{VC}(L.M.3) = 4$

In ∞ dimensions we can solve any model with a linear classifier of ∞ dimensions.

FINAL RESULTS:

• $|d_{VC}| \approx 10$ free parameters:

$$E_{out}(h) \leq E_{in}(h) + \sqrt{\frac{d_{VC}(\log(2n/d_{VC} + 1) + \log(2/\delta))}{2n}}$$

MLG : OVERFITTING :

LEARNING consists of finding a model such that $E_{out} \rightarrow 0$ by satisfying:

- 1) $E_{in} \rightarrow 0$ (by selecting a learning method that minimizes E_{in})
- 2) $E_{out} \approx E_{in}$. (by considering probabilistic things + Hoeffding's)

As a consequence: $E_{out} \leq E_{in} + O\left(\sqrt{\frac{C}{N}}\right)$

C : notion of complexity

N : amount of data samples.

* Independently of the data generation process generating we have to match the data complexity and not the model complexity (9)

LEARNING CURVE: curve of training and test errors as the number of training data increases for a given complexity.

(Error rate as: $1 - \text{accuracy over } \begin{matrix} \text{training} \\ \text{test} \end{matrix}$)

- Note that as the # of training data $\uparrow \Rightarrow$ both errors tend to a value (BIAS)
- Note that when \downarrow # of training data \Rightarrow training error is small but test error is large.
- With small complexity \Rightarrow training and test error converge with smaller N
 \Rightarrow the error of converge is larger !

Difference between BIAS and test error is the VARIANCE

OVERFITTING: increment of test error that occurs when a certain complexity level is attained. However, the training error might be reduced as the complexity increases.

How to cure over fitting: $E_{\text{test}} \leq E_{\text{tr}} + O\left(\sqrt{\frac{C}{N}}\right)$

- 1) Cross-validation: simulating the out-of-sample error and checking against unseen data several times.
- 2) Regularization: change learning objective + minimize complexity
(adding to the obj function a term that penalizes complexity)
- 3) Ensemble techniques.