

Overview

Large datasets in Science and technology

Practical Uses

- a) Business Transactions
- b) Digging Social Media
- c) AI/ Robotics

Physics / Astro / Bio :

a) HEP : CERN LHC $\rightarrow 25 \text{ GB/s}$ (with all events)
Nominal rate $68 \text{ TB/s} = 6 \times 10^8 \text{ events/s} \rightarrow 100 \text{ events/s}$
 $1 \text{ GB/s} \rightarrow 100 \text{ TB/day}$

b) Astro : Large Survey telescopes
LSST $\rightarrow 20 \text{ TB/night}$

c) Bio : Large Genome centers like
New York Genome Center $\sim 10 \text{ TB/day}$

Need for automated analysis
to flag interesting patterns \rightarrow Machine Learning (ML)

The book we follow emphasizes probabilistic approach. Not all methods are dependent on probabilistic formulation but we will mostly take prob models as these are in familiar territory for physicists.

ML and Statistics overlaps a lot.
We will quickly cover relevant corners of statistics.

Types of machine Learning

a) Supervised Learning:

Given $D = \{(x_i, y_i)\}_{i=1}^N \rightarrow$ training set

feature vectors
 $x_i \in \mathbb{R}^D$

outputs
categorical, real or ordinal

Predict for x . $y = h(x)$

b) Unsupervised Learning

Only $\{x_i\}_{i=1}^N$ are given. Find hidden structure.

Examples: clustering, dimension reduction, ..

c) Reinforcement Learning

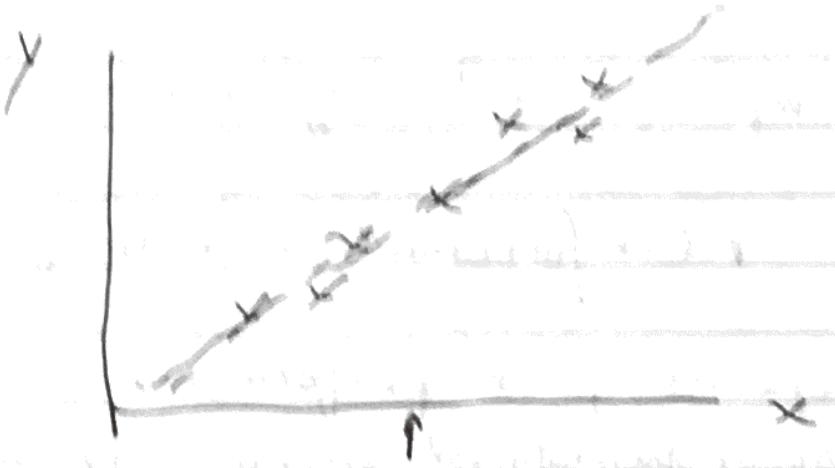
Depending on performance in some task reward or punishment signal is sent to the algorithm. The algorithm then tries itself to perform better.

Dynamic training data.

Examples: Game playing Systems.

Supervised Learning

Regression:



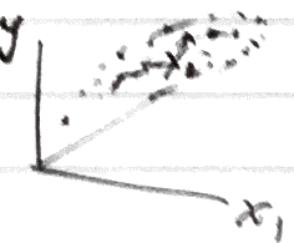
Linear Regression



Non linear Regression.

How complex should our fitting function be?

Multivariate examples:



Classification

$\{(x_i, y_i)\} \quad y_i \in \{0, 1\}$, binary classification



- Given any x , give us a probability
- of corresponding y being 0, 1.

Generalization to y_1, \dots, y_n , multiclass classification

Example: i) Given a galaxy image, classify it as elliptical, spiral, ...

ii) Given some protein sequence, assign it to a structural or functional class.

Probabilistic approach to Supervised learning:

Want a function approx.

$$\hat{y} = f(x) \quad \text{[redacted]}$$

Construct $P(y|x, D)$

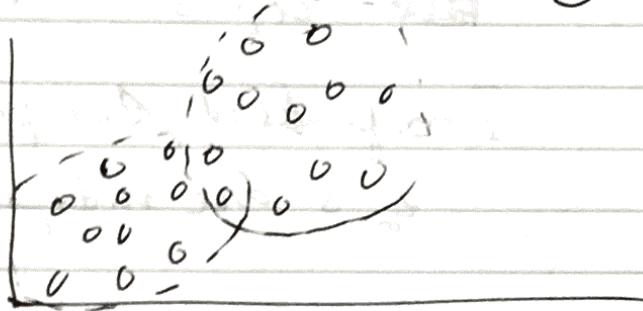
$$\hat{y} = \operatorname{argmax}_y P(y|x, D) = \hat{f}(x)$$

Maximum a posteriori (MAP)
estimate.

Practical applications:

- a) document classification, e.g. spam filtering
- b) Image classification, e.g. handwriting recognition,
- c) Object recognition, e.g. face detection.

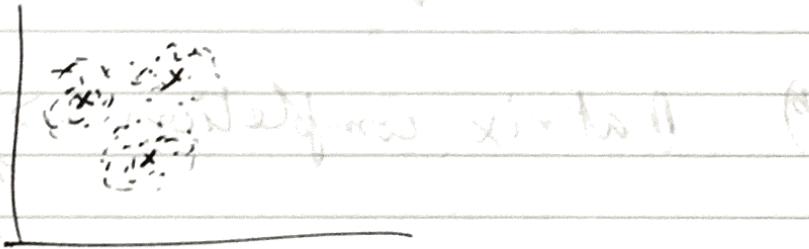
Unsupervised Learning



Knowledge discovery, density estimation, ..

Want. $P(x_i | \theta)$ as opposed to $P(y_i | x_i; \theta)$

Discovering clusters



Discovering latent factors

$$x_i = A z_i = \sum_{\alpha=1}^k a_{\alpha} z_{\alpha i}$$

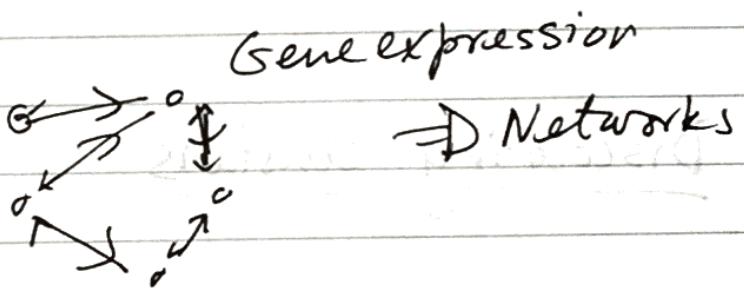
$$A = [a_1 \cdots a_k]$$

Examples: Principal component analysis
~~a_i~~'s orthogonal (PCA)

Independent Component Analysis
 z_i 's independent. (ICA)

Additional examples:

a) Discovering underlying graph structures



b) Matrix completion \rightarrow ~~decomposing~~
imputing missing data,

recommendation system.

Some basic Concepts

Parametric vs non-parametric ~~models~~ models

(The number of parameters fixed)

Parametric models for regression and classification

Regression

$$y(x) = w^T x + \epsilon$$

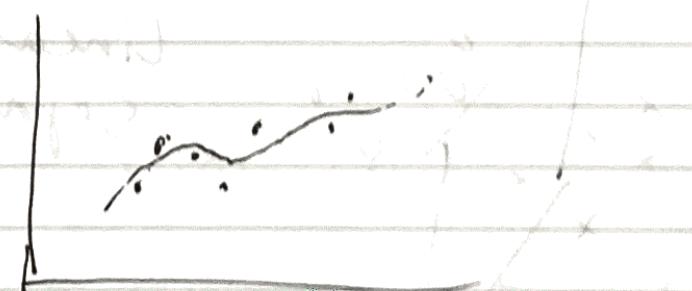
↑ residual error

$$P(y|x,\theta) = N(y|\mu(x), \sigma^2(x))$$

$$\mu(x) = w^T x$$

~~More generally~~ More generally $\mu(x) = w^T \phi(x)$

For example $\phi(x) = [1, x, x^2, \dots, x^n]$



$$y_i = w_0 + w_1 x_i + w_2 x_i^2 + \dots$$

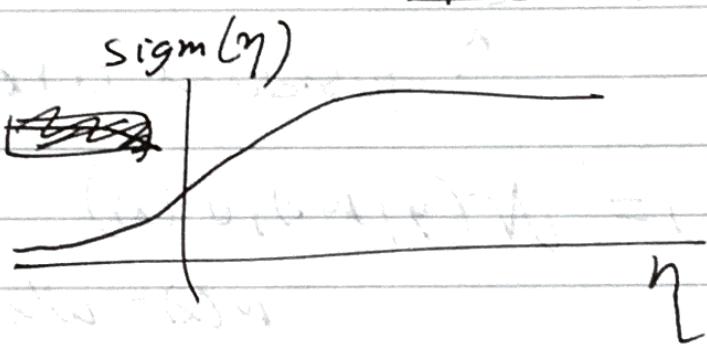
Classification

Logistic regression

$$P(y|x, w) = \text{Ber}(y|\mu(x))$$

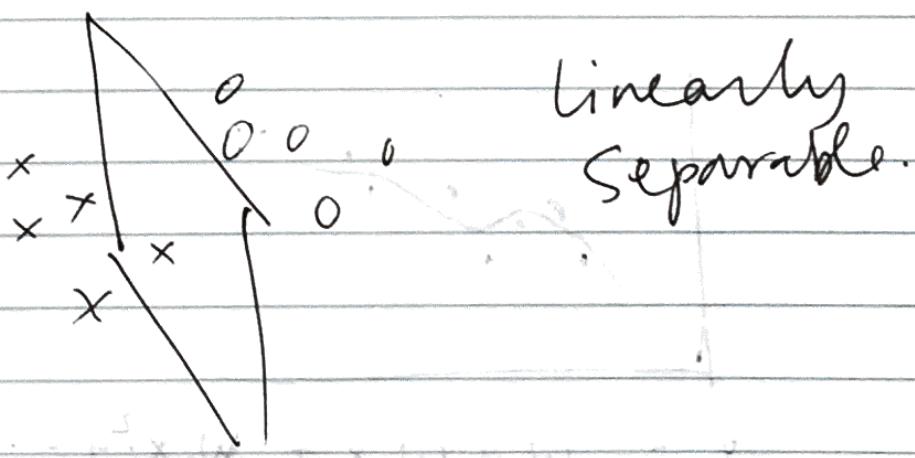
$$\mu(x) = \text{sigm}(w^T x)$$

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



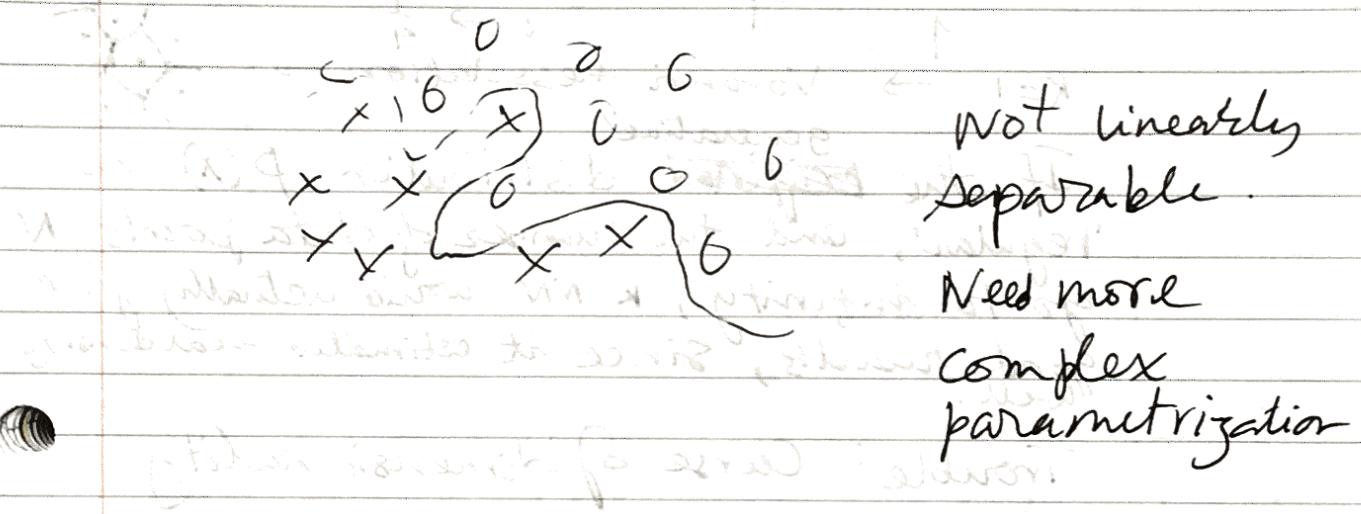
Decision rule $\hat{y}(x) = \begin{cases} 1 & P(y|x) > \frac{1}{2} \\ 0 & \text{otherwise} \end{cases}$

$$w^T x = 0$$



Advantage : Efficient

Disadvantage: ⚡ Less flexible

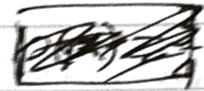
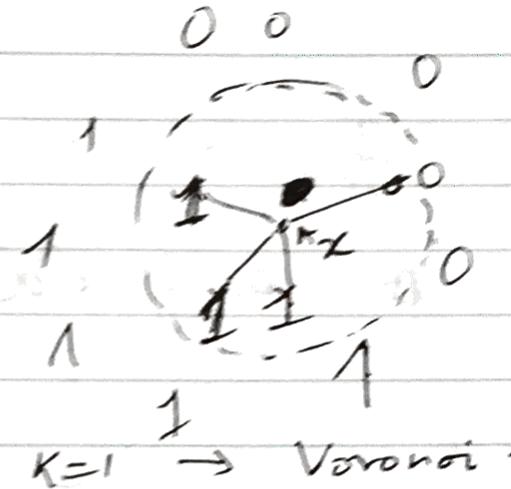


Non-parametric Classifier, Example

K-nearest neighbor model.

$$p(y=c|x, D, K)$$

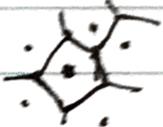
$$= \frac{1}{K} \sum_{i \in N_k(x|D)} I(y_i = c)$$



$$1 \rightarrow \frac{3}{4}$$

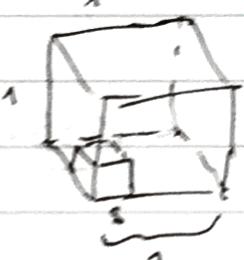
$$0 \rightarrow \frac{1}{4}$$

$K=1 \rightarrow$ Voronoi tessellation \rightarrow



If the ~~regular~~ distribution $p(x)$ is 'regular', and the number of data points N go to infinity, k -NN would actually give good results, since it estimates local density well.

Trouble: Curse of dimensionality



Example: data distributed uniformly in a ^{unit} hypercube

A smaller cube of size s has fractional volume s^D

The ~~#~~ number of points within the smaller volume $\propto N s^D$, implying $s = (\frac{K}{N})^{1/D} = f^{1/D}$

For large D and practical values of K, N s is often close to 1. For example, if

$$K \sim 10, D \sim 10, N \sim 10^4$$

$$\Rightarrow s \sim .5$$

So points involved are not closeby.

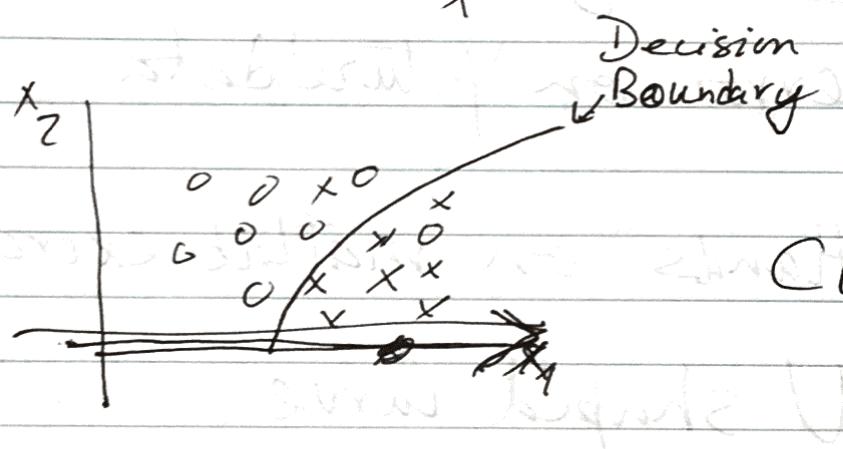
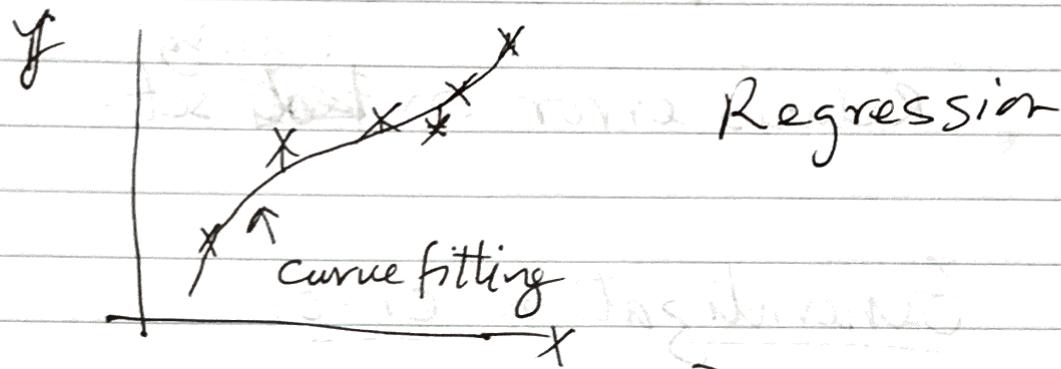


Overfitting:

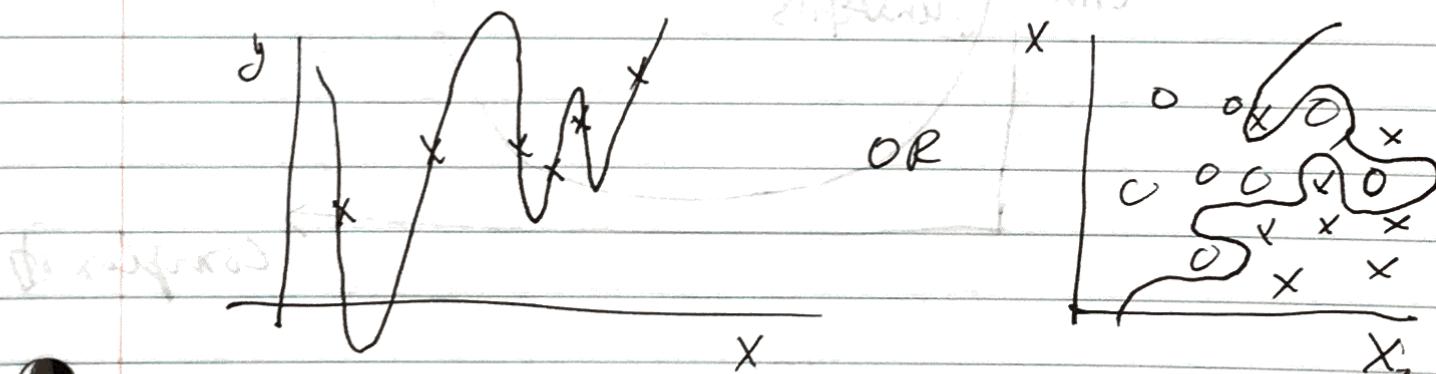
We usually minimize

$$\text{err}(f, D) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i), y_i)$$

least squares
misclassification count.



We do not want



even if it does a perfect job
on training data.

Model Selection

Too simple a model \Rightarrow increased error

Too complex a model \Rightarrow ?

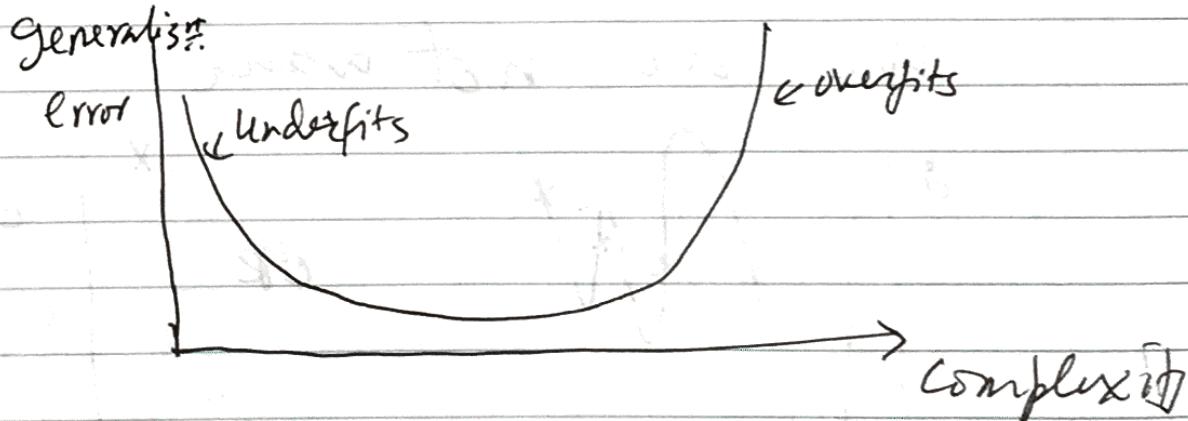
Reduces error on ~~test~~ set.
training

Generalization error

error on 'future' data

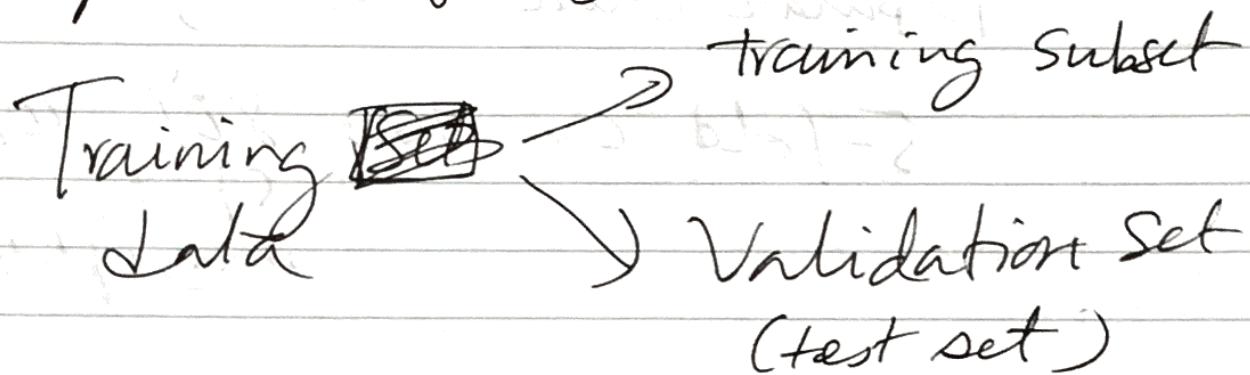
Hands on matLab exercise.

U shaped curve



Common strategy to guard

against overfitting!

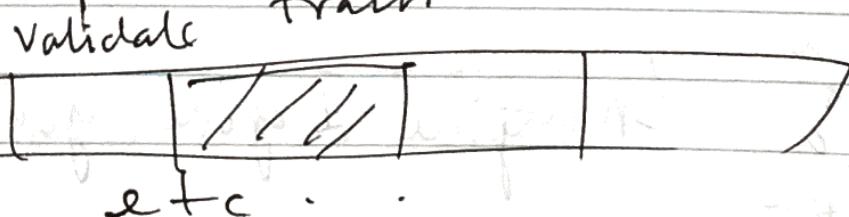


Often training set 80%
test set 20%.

If the dataset is small, often not enough statistics:

Popular solution: Cross validation (cv)

Split data & fold



Popular choice $K=5$

5-fold CV has 80%, 20%.
(split, but
many splits)

For size N training set

$K=N \Rightarrow$ Leave one out CV
LOOCV

Choosing number of nearest neighbors
 K_{NN} .

parameters in
regression.

No free lunch theorem:

(Wolpert 1996)

Training

Imagine I give you classification

task: $D = \{(x_i, y_i)\}_{i=1}^N$ is the training set

$D_V = \{(x_{j+N}, y_{j+N})\}_{j=1}^M$ is the validation set

For simplicity assume no x 's are the same

Let $f: X \rightarrow Y$

On $\{x_1, \dots, x_{N+M}\}$ there are 2^{N+M} choices for f . Consider

a prob. distr. $P(f)$ that uniformly samples from these 2^{N+M} functions.

If D^+ and D_V^+ are generated by f , and algorithm training on D produces a function \hat{f}
 Generalization error $\text{err}(\hat{f}, D_V) = \frac{1}{M} \sum_{j=1}^M \mathbb{I}(\hat{f}(x_{j+N}) \neq y_{j+N})$

$$E[\text{err}(\hat{f}, D_V)] \sum_f \text{err}(\hat{f}, D_V) P(f) = ?$$

$$E[\text{err}(\hat{f}, D_V)] = \frac{1}{2}$$

If algorithm $A \rightarrow f_{A, D^+}$ $B \rightarrow f_{B, D^+}$
 Their ^{average} performance over the validation set is the same

So it cannot be that one algorithm always outperforms another under all circumstances.

⇒ Importance of tailoring algorithms to ~~match~~ the structure of data to be processed.

Role of bias in learning?