

High level overview of ML:

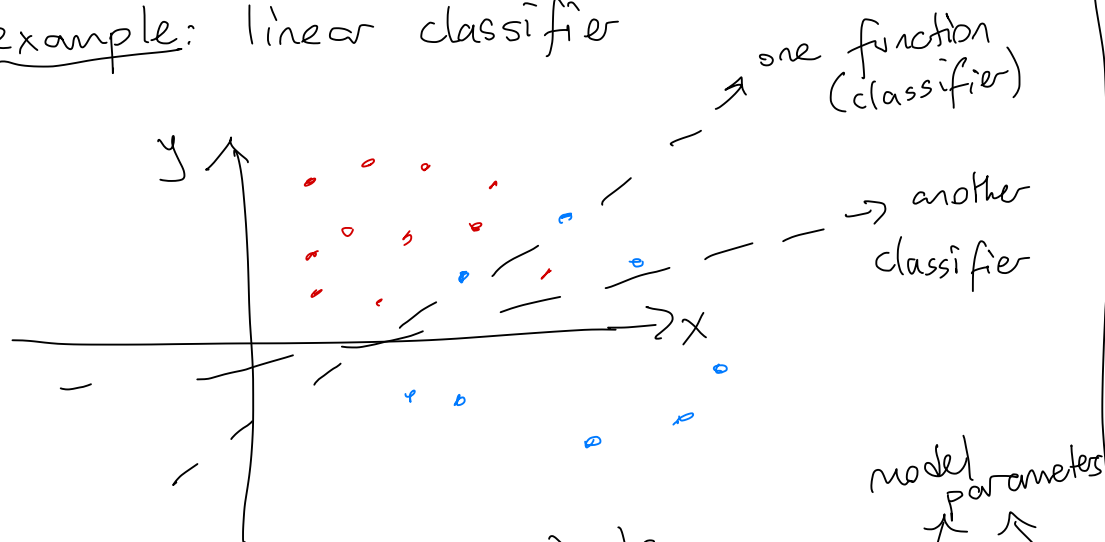
data { if labelled: $(X_{N \times d}, y_{N \times 1})$,
 otherwise $(X_{N \times d})$: N : # of datapoints
 d : # of dimensions of each datapoint

$X_{i,j} \in \mathbb{R}$ or some interval $[a,b]$
 $y_i \in \mathbb{R}$ for regression, or a discrete set for classification

we determine a specific type of function (sometimes called estimator) that we want to use to infer something useful like the y value or class or cluster it belongs to.

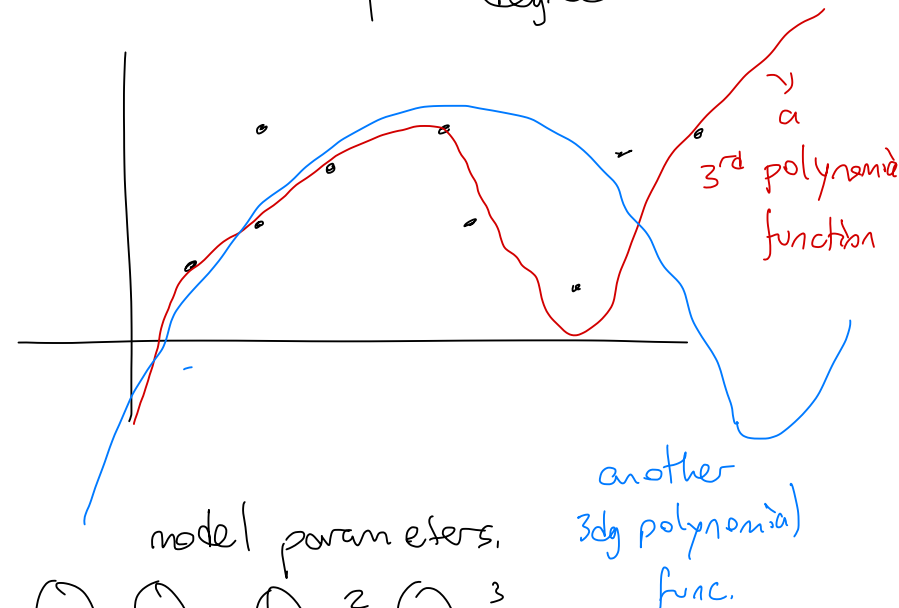
This function is parameterized. These model parameters are what the "machine learns".

example: linear classifier



treated as a parameterized function because these all look like $(a)x + (b)$
 $A.X + b$

Example: polynomial regression:
 of 3rd degree



$$(a) + (a_1)x + (a_2)x^2 + (a_3)x^3 = y$$

How does the learning work? Usually there is some metric (i.e. loss, cost) that is being minimized by changing the model parameters. In the case of the example 2 it can be squared error or absolute error or others (meaning that your learning changes the model parameters to minimize the squared error between the guess \hat{y} and true value y). In example 1 it can be cross-entropy loss or hinge loss or squared loss between one-hot label vectors.

$$\begin{matrix} 2^{\text{nd}} \\ \text{class} \end{matrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad \begin{matrix} 1^{\text{st}} \\ \text{class} \end{matrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

* One popular way of learning (or optimizing your loss) is called gradient descent. (question 2c) In this method, you take numeric gradient of the loss with respect to all model parameters: $\frac{\partial \text{loss}}{\partial a}, \frac{\partial \text{loss}}{\partial b}$ and then update your parameters: $a^{(i+1)} = a^{(i)} - \eta \cdot \frac{\partial \text{loss}}{\partial a^{(i)}}$. η is called learning rate.

then repeat until your loss converges (or maxiter)

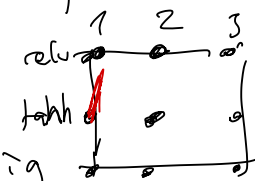
Hyperparameters and Grid Search Hyperparameters are different from parameters they are not learnt from data.* It is you, the ML practitioner that tweaks these hyperparameters (or searches over them automatically). The set of hyperparameters determine a family of functions.

Example 3: polynomial regression: hyperparameter: degree of the polynomial
parameters: coefficients of the polynomial

Example 4: neural network: hyperparameters: depth of the network
of neurons in each layer
activation function
neuron connectivity (convolutional, dense etc)
optimizer arguments.
parameters: weights & biases

Grid search: Divide the space of possible hyperparameters and for each point in this space, train the model. See which parameters work the best.

grid: dg 1 2 3 4 Ex3

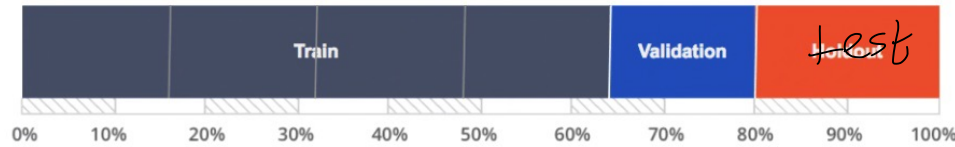
Ex4: grid,  → depth of net

* Technically you, as the programmer fit the validation data to the hyperparameters.

Validation vs test data

* Separate validation and test data to not fool yourself. "You are the easiest person to fool."

* By changing hyperparameters you implicitly fit the model to validation data. If validation = test it's like solving the exam question in class.



Cross validation

5-fold:

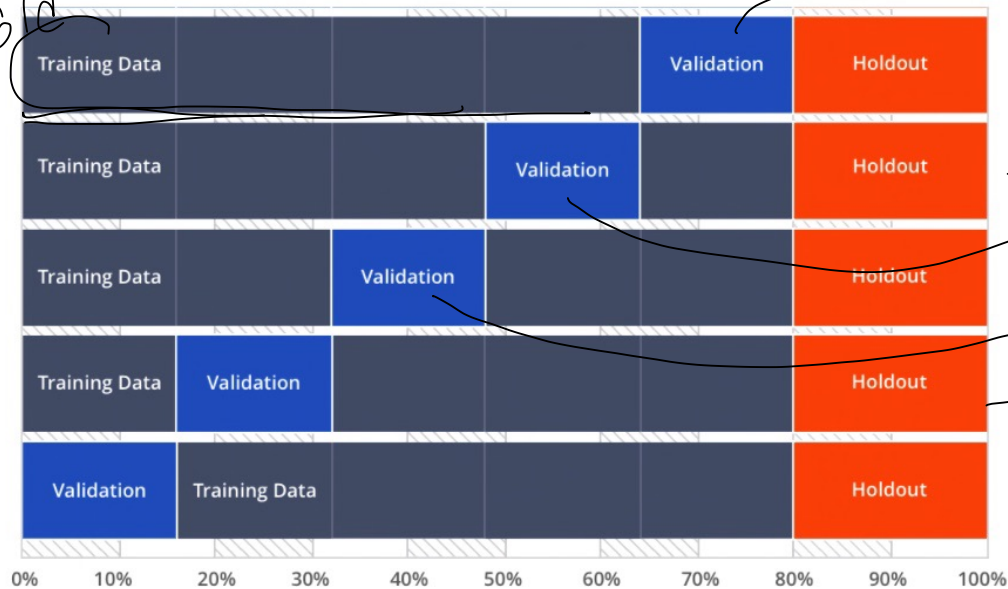
10

K-fold

K times

train

1st fold



→ 95%

→ 94%

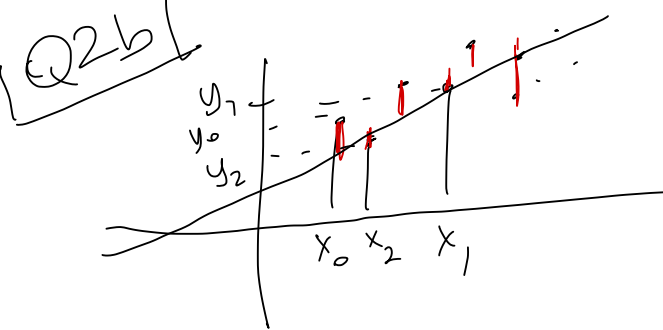
→ 93%

→ 96%

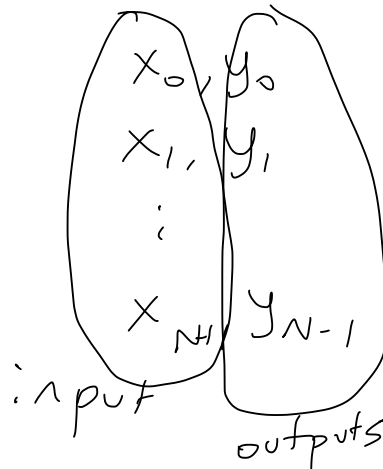
avg 94.5%

- Divide training data into 5 "folds", keep 1 as validation set. Train model on the remaining 4/5's.
- Repeat using other folds.

Importance: Decrease the bias in accuracy.



N : number of data points



find $a, b \in \mathbb{R}$

$$\underbrace{a \underline{x}_i + b = \underline{y}_i}_{\text{eqn}} \}_{N}$$

$$\underline{X} = \begin{pmatrix} x_0 & 1 \\ x_1 & 1 \\ x_2 & 1 \\ \vdots & \vdots \end{pmatrix}_{N \times 2}$$

$$\underline{y} = \begin{pmatrix} y_0 \\ \vdots \\ y_{N-1} \end{pmatrix}$$

$$\underline{X}_{N \times 2} \underbrace{\begin{bmatrix} a \\ b \end{bmatrix}}_{\underline{\delta}} = \underline{y}_{N \times 1}$$

can be rewritten

$$\underline{X} \cdot \underline{\delta} = \underline{y}$$

Least squares solution to $\begin{bmatrix} a \\ b \end{bmatrix}$:

$$\text{"squared error"} = \sum_{i=0}^{N-1} \left(\underbrace{\underline{X} \cdot \underline{\delta}}_{\hat{y}} - \underline{y} \right)^2_{N \times 1} = (\underline{X} \cdot \underline{\delta} - \underline{y})^T_{1 \times N} (\underline{X} \cdot \underline{\delta} - \underline{y})_{N \times 1}$$

Convex in $\underline{\delta}$ so $\frac{\partial \text{"squared error"}}{\partial \underline{\delta}_i} = 0$ gives min or max.

$$\frac{\partial}{\partial \underline{\delta}} \underline{\delta}^T \underline{X} \underline{X} \underline{\delta} - \underbrace{\underline{\delta}^T \underline{X}^T \underline{y}}_{\text{same}} - \underline{y}^T \underline{X} \cdot \underline{\delta} = 0 \Rightarrow \underbrace{-2 \underline{X}^T \underline{y} - 2 \underline{X}^T \underline{X} \underline{\delta}}_{\text{same}} = 0 \Rightarrow$$

$$\underline{X}^T \underline{X} \underline{\delta} = \underline{X}^T \underline{y} \quad \left. \begin{matrix} \underline{X}^T \underline{X} \\ \text{invertible} \end{matrix} \right\} \underline{\delta}^* = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

solution

