

Machine Learning

Federico Matteoni

A.A. 2021/22

Index

1	Machine Learning	2		
1.1	Machine Learning	2	1.3	Model Selection and Model Assessment 20
1.1.1	Machine Learning	2	1.3.1	Bias-Variance 20
1.1.2	Statistical Learning Theory	5	1.3.2	Motivations 20
1.1.3	Linear Models	6	1.3.3	Validation 20
1.1.4	Gradient Descent	8	1.4	Statistical Learning Theory 23
1.1.5	Extending the linear model	10	1.4.1	VC-dim 23
1.2	Neural Networks	11	1.4.2	Structural Risk Minimization 24
1.2.1	Artificial Neuron	12	1.5	Support Vector Machines 25
			1.5.1	High-Dimensional feature spaces 30
			1.5.2	SVM for non-linear regression 31
			1.5.3	Kernel Methods 33
			1.6	Bias-Variance 33
			1.6.1	Bias-Variance Decomposition 34
			1.7	Ensemble Learning 35
			1.7.1	Bagging 35
			1.7.2	Boosting 35
			1.7.3	Feature Selection 35
			1.8	Applications 35
			1.8.1	Character recognition (classification) 35
			1.8.2	Convolutional Neural Networks 36
			1.8.3	Deep Learning 38

Capitolo 1

Machine Learning

1.1 Machine Learning

What is ML? Area of research combining aims of creating computers that could learn and powerful and adaptive statistical tools with rigorous foundation in computational science. Luxury or necessity? Growing availability and need for analysis of empirical data and difficult to provide intelligence and adaptivity by programming it. Change of paradigm.

Examples: spam classification, written text recognition. . . No or poor prior knowledge and rules for solving the problem, but easier to have a source of training experience.

ML is considered the latest general-purpose technology, capable of drastically affect pre-existing economic and social structures. And already has. The ultimate aim is to bring benefits to the people by solving big and small problems, accelerating human progress and empowering humans to add intelligence in any other science field.

1.1.1 Machine Learning

We restrict to the computational framework: principles, methods and algorithms for learning and prediction, from experience. Building a model to be used for predictions. Common framework: infer a model or a **function** from a set of examples which allows the generalization (accurate response to new data).

When can we use ML? Be aware of the opportunity and awareness. ML is useful when there's no or poor theory surrounding the phenomenon, or uncertain, noisy or incomplete data which hinders formalization of solutions. The requests are: source of training experience (representative data) and a tolerance on the precision of results. The best examples are models to solve real-world problems that are difficult to be treated with traditional techniques: face and voice recognition (knowledge too difficult to formalize in an algorithm), predicting binding strength of molecules to proteins (not enough human knowledge) and personalized behavior, such as recommendation systems, scoring messages according to user preferences. . .

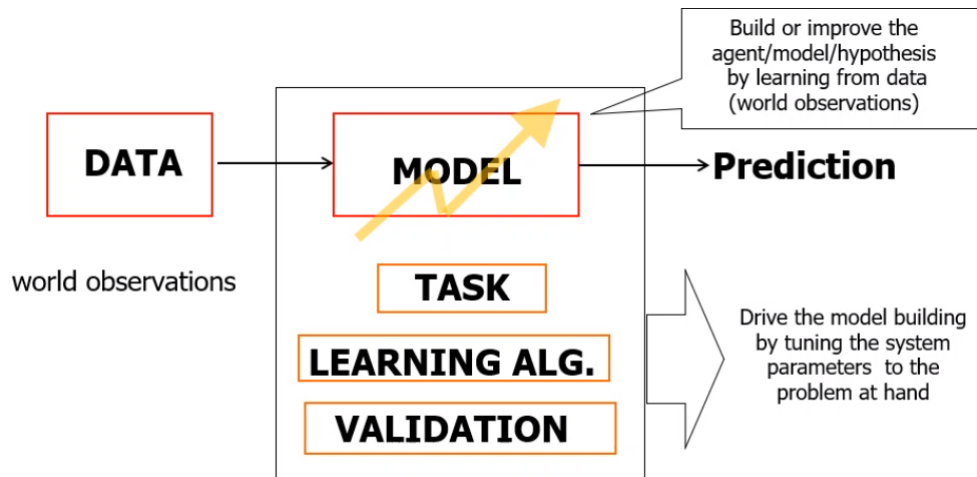
Definition The ML studies and proposes methods to build functions/hypothesis from examples of observed data that fits the known examples and able to generalize, with reasonable accuracy, for new data (according to verifiable results and under statistical and computational conditions and criteria.

Data Data **represents the available experience**. Representation problem: capturing the structure of the analyzed objects. Flat (attribute-value), structured. . . , categorical or continuous, missing data. . . **preprocessing**: variable scaling, encoding, selection. . .

Task The task defines the purpose of the application: knowledge that we want to achieve? which is the helpful nature of the result? what information are available?

Predictive task, classification and regression: function approximation

Descriptive task, cluster analysis and association rules: find subsets or groups of unclassified data.



Also as a guide to the **key design choices**
(ML system “ingredients”)

Supervised learning

Given a set of training examples as $\langle \text{input}, \text{output} \rangle = \langle x, d \rangle$ (**labeled examples**) for an unknown function f , find a *good approximation* of f , an hypothesis h that can be used for making predictions on unseen data x' .

The target d can be:

Discrete value, for **classification tasks**.

$$f(x) \in \{1, 2, \dots, k\}$$

Patterns, feature vectors, are seen as members of a class and the goal is to assign the new patterns observed to the correct class (or label)

If the number of possible classes is two, then f is a *boolean function* and the task is called **binary classification** or **concept learning**: true or false, positive or negative, 0 or 1...

If the number of classes is greater than two then it is a **multi-class classification task**.

Real continuous value, for **regression tasks**.

The patterns are seen as sets of variables (real values), and the task is a curve fitting task. The process aims to estimate a real-value function based on a finite set of noisy samples $\langle x, f(x) + \text{random noise} \rangle$

Unsupervised learning

No teacher. The training set is a set of unlabeled data $\langle x \rangle$. Examples: clustering, finding natural groupings in a set of data.

Learning algorithm

Basing on data, task and model: heuristic search through the hypothesis space H of the **best hypothesis**. I. e. the best approximation of the unknown target function, typically searching for the h with the minimum *error*. H may not coincide with the set of all possible functions and the search cannot be exhaustive, we need to make **assumptions** (**inductive bias**).

Learning Also called:

Inference, in statistics

Adapting, in biology and systems

Optimizing, in mathematics

Training, in neural networks

Function approximations, in mathematics

...

After introducing data, task, model and learning algorithm we will focus on: inductive bias, loss and concepts of generalization and validation.

Inductive bias To set up a model we can make assumptions about the nature of the target function, concerning either:

constraints in the model, **language bias** (in the hypothesis space H , due to the set of hypothesis that we can express or consider

constraints or preferences in learning algorithm/search strategy, **search bias** which is preferred

or both

Such assumptions are needed to obtain an useful model for the ML aims, i.e. a model with generalization capabilities. We can imagine learning a discrete function with discrete inputs assuming **conjunctive rules**, so using a **language bias** to work with a restricted hypothesis space.

Version Space An hypothesis h is consistent with the TR if $h(x) = d(x)$ for each training example $\langle x, d(x) \rangle$. The **version space** $VS_{H,TR}$ is the subset of H of the hypothesis consistent with all the training examples $\langle x, d(x) \rangle$ in the TR.

It's possible to do an exhaustive search in an efficient way, using clever algorithms. This means finding the set of all the hypothesis h consistent with the TR set.

Unbiased Learner The language bias (ex: using only conjunctive rules, may be too restrictive: if the target concept is not in H it cannot be represented in H . We can use an H that expresses every teachable concept (among propositions), that means that H is the set of all possible subsets of X : the power set $P(X)$. If $n = 10$ binary inputs, then $|X| = 2^{10} = 1024$ and $|P(X)| = 2^{1024} = 10^{308}$ possible concepts, which is much more than the number of the atoms in the universe.

An unbiased learner is unable to generalize: the only examples that are unambiguously classified by an unbiased learner represented with the VS are the training examples themselves. Each unobserved instance will be classified positively by exactly half of the hypothesis in the VS and negative by the other half. Indeed: $\forall h$ consistent with x_i , $\exists h'$ identical to h except $h'(x_i) \neq h(x_i)$, $h \in VS \Rightarrow h' \in VS$ (because they are identical on the TR)

Why prefer the search bias? In ML we use flexible approaches (expressive hypothesis spaces with universal capability of the models, for example neural networks or decision trees. We avoid the language bias, so we do not exclude a priori the unknown target function, but we focus on the search bias (ruled by the learning algorithm).

Loss How to measure the quality of an approximation? We want to measure the distance between $h(x)$ and d , using a loss function/measure $L(h(x), d)$ for a pattern x which has high value in cases of bad approximation. The error (or risk or loss) is an expected value of this L , for example $E(w) = \frac{1}{l} \sum_{p=1}^l L(h(x_p), d_p)$. Different L for different tasks. Examples of loss functions:

Regression: $L(h(x_p), d_p) = (d_p - h(x_p))^2$, the squared error. MSE (mean squared error) over the data set

Classification: $L(h(x_p), d_p) = \begin{cases} 0 & h(x_p) = d_p \\ 1 & \text{else} \end{cases}$

Learning and generalization Learning: search for a **good function** in a function space from known data (typically minimizing an error/loss). **Good** with respect to generalization error: it measures how accurately the model predicts over novel samples of data (**measured over new data**).

Generalization is the crucial point of ML. Performance in ML is the generalization accuracy or *predictive accuracy* estimated by the error on the test set.

ML issues Inferring general functions from known data is an ill posed problem, which means that in general the solution is not unique because we can't expect the exact solution with finite data. What can we represent? And so, what can we learn?

Learning phase: building the model including training. The prediction phase is evaluating the learned function over new never-seen-before samples (generalization capability). Inductive learning hypothesis: any h that approximates f well on training examples will also approximate f well on new unseen instances x .

Overfitting: a learner overfits data if it outputs an hypothesis $h \in H$ having true/generalization error (risk) R and empirical (training) error E , but there's another $h' \in H$ with $E' > E$ and $R' < R$, which means that h' is the better one despite having a worse fitting.

1.1.2 Statistical Learning Theory

Under what mathematical conditions is a model able to generalize? We want to investigate the generalization capability of a model, measured as a risk or test error, the role of the model complexity and the role of the number of data.

Formal Setting: approximate a function $f(x)$, with d target ($d = f(x) + \text{noise}$), minimizing the **risk function**

$$R = \int L(d, h(x)) dP(x, d)$$

which is the **true error over all the data**, given:

a value d from the teacher and the probability distribution $P(x, d)$

a loss function $L(h(x), d) = (d - h(x))^2$

We search for $h \in H \mid \min R$, but we only have the finite data set $TR = (x_p, d_p)$ with $p = 1 \dots l$. Looking for h means minimizing the empirical risk (the training error E), finding the best values for the model free parameters

$$R_{emp} = \frac{1}{l} \sum_{p=1}^l (d_p - h(x_p))^2$$

The inductive principle is the **ERM**, Empirical Risk Minimization: can we use R_{emp} to approximate R ?

Vapnik-Chervoneniks dim and SLT

Given the VC dimension (simply VC), a measure of complexity of H and by that we mean its flexibility to fit data. The VC-bound states that it holds with probability $\frac{1}{\delta}$ that

$$R \leq R_{emp} + \epsilon \left(\frac{1}{l}, VC, \frac{1}{\delta} \right)$$

ϵ is a function called VC-confidence, that grows with VC and decreases with higher l and δ

R_{emp} decreases using complex models (with high VC)

δ is the confidence, and it rules the probability that the bound holds.

$\delta = 0.01 \Rightarrow$ the bound holds with probability 0.99

Intuitively:

Higher l (data) \Rightarrow lower VC confidence and bound closer to R

A too simple model, meaning with low VC, can be not sufficient due to high R_{emp} (**underfitting**)

An higher VC with fix $l \Rightarrow$ lower R_{emp} but VC and hence R may increase (**overfitting**)

Structural risk minimization

Minimize the bound! There are different bounds formulations according to different classes of f , of tasks...

In other words, we can make a good approximation of f from examples, provided that we have a good number of data and the complexity of the model is suitable for the task.

Complexity control

The Statistical Learning Theory allows for a formal framing of the problem of generalization and overfitting, providing an analytic upper bound to the risk R for the prediction over all the data, regardless of the type of learning algorithm or the details of the model. So **the machine learning is well founded**, the learning risk can be analytically limited and only a few concepts are fundamental. This leads to new models (such as the Support Vector Machine) and other methods that directly consider the control of the complexity in the construction of the model.

Validation Central role for the applications and the project. Two aims:

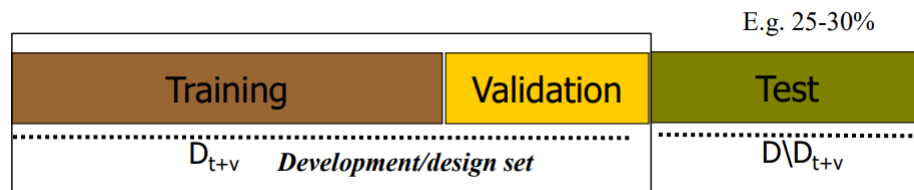
Model Selection: estimating the performance (**generalization error**) of different models in order to choose the best one. This includes searching for the best hyperparameters of the model.
It returns a model.

Model Assessment: with the final model, estimating/evaluating its prediction error/risk (**generalization error**) over new test data.
It returns an estimation.

Golden rule: keep the two goals separated and use different datasets for each one.

In an ideal world, we'd have a large training set, a large validation set for model selection and a very large external unseen data test set. With finite and often small data sets we have just an estimation of the generalization performance. We have to use some techniques: hold-out and k-fold cross validation, for example.

Hold-Out: we partition the dataset D into **training set** TR, **validation/selection set** VL and **test set** TS. All three are disjoint: TR is used to run the training algorithm, VL can be used to select the best model (hyperparameters tuning) and the **TS is only used for model assessment**.



K-Fold: this technique can make use of insufficient data. We split the dataset D into k mutually exclusive subsets D_1, \dots, D_k , we train on $D - D_i$ and test it on D_i .

This can be applied to both VL and TS splitting. Can be computationally very expensive and there's the issue of choosing the number of folds k .



Confusion Matrix

Actual/Predicted	Positive	Negative
Positive	TP	FN
Negative	FP	TN

Specificity = $\frac{TN}{FP+TN}$, and **true negative rate** = $1 - FPR$

Sensitivity = $\frac{TP}{TP+FN}$, also known as **true positive rate** or **recall**

Precision = $\frac{TP}{TP+FP}$

Accuracy: % of correctly classified patterns = $\frac{TP + TN}{total}$. Note that, for example, a 50% accuracy on a binary classifier is equivalent to a random predictor.

ROC Curve We plot **specificity** on **x-axis** and **sensitivity** on the **y-axis**. The diagonal corresponds to the worst classifier, the random guesser. Better curves have greater Area Under the Curve (AUC)

1.1.3 Linear Models

Mainstay of statistics.

Univariate Linear Regression

Simple linear regression: we start with 1 input variable x and 1 output variable y . We assume a model $h_w(x)$ expressed as $out = w_1x + w_0$ where w are real-valued coefficients or **free parameters**, also called **weights**.

Given that the w s are continuous valued, we have an infinite hypothesis space but a nice solution from classical math. We can learn with this basic tool and, although simple, it includes many relevant concepts of modern ML and many methods in the field are based on this.

Least Mean Square: learning means finding w such that it minimizes the error/empirical loss, with best data fitting on the training set with l examples.

So given a set of l training examples (x_p, y_p) with $p = 1, \dots, l$, we have to find $h_w(x)$ in the form $w_1x + w_0$ that minimizes the expected loss on the training data. For the loss, we use the square of errors: **least mean square**, find w to **minimize** the residual sum of squares.

$$Loss(h_w) = E(w) = \sum_{p=1}^l (y_p - h_w(x_p))^2 = \sum_{p=1}^l (y_p - (w_1x_p + w_0))^2$$

To have the mean, divide by l . How to solve? Local minimum as stationary point, so the gradient $\frac{\partial E(w)}{\partial w_i} = 0$ with $i = 1, \dots, \text{dim_input} + 1 = 1, \dots, n + 1$. For the simple linear regression (2 free parameters):

$$\begin{aligned} \frac{\partial E(w)}{\partial w_0} &= 0 & \frac{\partial E(w)}{\partial w_1} &= 0 \\ \frac{\partial E(w)}{\partial w_0} &= -2(y - h_w(x)) & \frac{\partial E(w)}{\partial w_1} &= -2(y - h_w(x)) \cdot x \end{aligned}$$

Classification

The same models used for regression can be used for classification: **categorical targets** y or d , for example 0/1, -1/+1...

We use an hyperplane (wx) assuming negative or positive values. We exploit such models to decide if a point x belongs to the positive or the negative zone of the hyperplane to classify it. So we want to learn w such that we get a good classification accuracy. The decision boundary is $x^T w = w^T x = w_0 + w_1x_1 + w_2x_2 = 0$ and we can introduce a threshold function which can be written in many ways:

$$h(x) = \begin{cases} 1 & \text{if } wx + w_0 \geq 0 \\ 0 & \text{else} \end{cases}$$

$$h(x) = \text{sign}(wx + w_0)$$

...

w_0 is called **threshold** or **bias**. $h(x) = w^T x + w_0 \geq 0 \Leftrightarrow w^T x \geq -w_0$

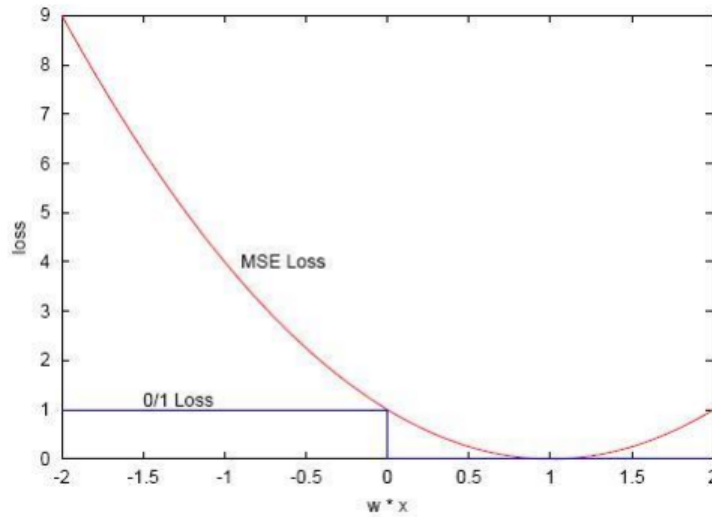
Learning Algorithms

Introducing 2 learning algorithms, both based on LSM and used for the linear model on regression and classification tasks. We start redefining the learning problem and the loss for them (in the case of l data and multidimensional inputs).

Learning problem for classification tasks: given a set of l training examples (x_p, y_p) and a loss function L , with $y_p \in \{0, 1\}$ or $y_p \in \{-1, +1\}$, find the weight vector w that minimizes expected loss on the training data

$$R_{emp} = \frac{1}{l} \sum_{p=1}^l L(h(x_p), y_p)$$

The expected loss can be approximated by a smooth function. We can make the optimization problem easier by replacing the original objective function L (0/1 loss) with a smooth, differentiable function: for example, the MSE loss (mean squared error).



No classification error minimizing either 0/1 loss or MSE loss. Given l training examples (x_p, y_p) , find w that minimizes the residual sum of squares

$$E(w) = \sum_{p=1}^l (y_p - x_p^T w)^2 = \|y - Xw\|^2$$

We can't use $h(x)$ in $E(w)$, as for regression, because $h(x) = \text{sign}(w^T x)$ is non-differentiable. Also, this is a quadratic function so the minimum always exists (but may not be unique). X is a $l \times n$ matrix with a row for each x_p .

Direct Approach with a normal equation Differentiating $E(w)$ with respect to w to get the **normal equation**

$$(X^T X)w = X^T y$$

In the derivation we also find that

$$\frac{\partial E(w)}{\partial w_j} = -2 \sum_{p=1}^l (x_p)_j \cdot (y_p - x_p^T w)$$

If $X^T X$ is not singular, then the unique solution is given by

$$w = (X^T X)^{-1} X^T y = X^+ y$$

with X^+ being the Moor-Penrose pseudoinverse. Else the solutions are infinite, so we can choose the **min norm**(w) solution.

The **Singular Value Decomposition** can be used for computing the pseudoinverse of a matrix (X^+). With $X = U \Sigma V^T \Rightarrow X^+ = V \Sigma^+ U^T$ replacing every non-zero entry by its reciprocal. We can apply SVD directly to compute $w = X^+ y$, obtaining the minimal norm (on w) solution of least squares problem.

$$\frac{\partial E(w)}{\partial w_j} = \frac{\partial \sum_{p=1}^l (y_p - x_p^T w)^2}{\partial w_j} = \dots = -2 \sum_{p=1}^l (y_p - x_p^T w)(x_p)_j$$

1.1.4 Gradient Descent

The derivation suggests an approach based on an iterative algorithm based on

$\frac{\partial E(w)}{\partial w_j} = -2 \sum_{p=1}^l (y_p - x_p^T w)(x_p)_j$. The **gradient** is the **ascent direction**. We can move toward the minimum with a gradient descent $\Delta w = -\text{gradient of } E(w)$. **Local search:** we begin with a initial weight vector and modify it iteratively to minimize the error function. The gradient vector is

$$\Delta w = -\frac{\partial E(w)}{\partial w} = \begin{bmatrix} -\frac{\partial E(w)}{\partial w_1} \\ \vdots \\ -\frac{\partial E(w)}{\partial w_n} \end{bmatrix} = \begin{bmatrix} \Delta w_1 \\ \vdots \\ \Delta w_n \end{bmatrix}$$

Allowing us to work in a multi dimensional space without the need to visualize it. Hence, the iterative approach will move using a learning rule based on a "delta" of w proportional to the opposite of the local gradient. The movements will be made according to

$$w_{new} = w + \eta \cdot \Delta w$$

The simple algorithm is as follows:

1. Start with weight vector $w_{initial}$ and fix $0 < \eta < 1$
2. Compute $\Delta w = -\text{gradient of } E(w) = -\frac{\partial E(w)}{\partial w}$ (or for each w_i)
3. Compute $w_{new} = w + \eta \cdot \Delta w$ (or for each w_i)
 η is the step size or **learning rate**
4. Repeat from 2 until convergence or $E(w)$ sufficiently small

Batch version

The gradient is the sum over all the l patterns. Provides a more precise evaluation of the gradient over l data. We upgrade the weight after the sum

$$\frac{\partial E(w)}{\partial w_j} = -2 \sum_{p=1}^l (y_p - x_p^T w)(x_p)_j$$

Online/Stochastic version

We upgrade the weights with the error that is computed for each pattern. Hence, the second pattern output is based on weights already updated from the first and so on. In makes progress with each example it sees. Can be faster, but needs smaller η

$$\frac{\partial E_p(w)}{\partial w_j} = -2(y_p - x_p^T w)(x_p)_j = -\Delta_p w_j$$

Gradient Descent as error correction delta rule

The error correction rule, also called Widrow-Hoff or delta rule, changes w_j proportionally to the error (target y – output)

$$\Delta w_j = 2 \sum_{p=1}^l (x_p)_j (y_p - x_p^T w)$$

$$w_{new} = w + \eta \cdot \Delta w$$

We improve by learning on previous errors.

Gradient descent is a simple and effective local search approach to a LMS solution. It allows to search through an infinite hypothesis space, can be easily applied for continuous H and differentiable losses and isn't only for linear models (also neural networks and deep learning models).

Many possible improvements (Newton, quasi-Newton methods, conjugate gradients...)

Language bias: H is a set of linear functions.

Search Bias: ordered search guided by the least squares minimization goal. For instance, we could prefer a different method to obtain a restriction on the values of parameters, achieving a different solution with other properties.

Shows that even for a simple model there are many possibilities. We still need a principled approach.

Limitations In geometry, two set of points are linearly separable in an n -dimensional space if they can be separated by a $(n - 1)$ -dimensional hyper-plane. In 2 dimensions, if they can be separated by a line, in 3 dimensions, by a plane...

The linear decision boundary can provide exact solutions only for linearly separable sets of points.

1.1.5 Extending the linear model

We can use transformed inputs, such as $x, x^2, x^3 \dots$ with a non-linear relationship between inputs and output, maintaining the learning machinery used so far.

$$h_w(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Linear basis expansion (LBE)

$$h_w(x) = \sum_{k=0}^K w_k \phi_k(x)$$

Augments the input vector with additional variables which are transformations of x according to a function $\phi_k : R^n \rightarrow R$, so number of parameters $K > n$: linear in the parameters, so we can use the same learning algorithms as before. Which ϕ ? Towards the dictionary approaches. Pro: can model more complicated relationships than linear, so it's more expressive. Cons: with large basis of functions, we easily risk overfitting, hence we require controlling the complexity (as in flexibility of the model to fit the data). How to do that? Many approaches:

Ridge Regression (or Tikhonov Regularization): smoothed model.

Add constraints to the sum of value of $|w_j|$, penalizing models with high values of $|w|$ (so favoring sparse models, using less terms due to weights $w_j = 0$ or close)

$$Loss(w) = \sum_{p=1}^l (y_p - x_p^T w)^2 + \lambda ||w||^2$$

with λ being the regularization hyper-parameter. It implements the control of the model complexity, leading to a model with less VC-dim with a trade-off controlled through a single parameter, λ .

This uses $|| \cdot ||_2$

Lasso uses $|| \cdot ||_1$

Elastic nets uses both $|| \cdot ||_1$ and $|| \cdot ||_2$

Learning Timing

Eager: analyze data and construct an explicit hypothesis

Lazy: store tr data and wait test data point, then construct an ad hoc hypothesis.

k-NN The algorithm is simple: store the training data and given an input x find the k nearest training examples x_i , then output the mean label.

Voronoi Diagram Each cell consists of point closer to x than any other patterns. The segments are all points in plan equidistant to two patterns. It is implicitly used by K-NN.

K-NN vs linear Two extremes of the ML panorama:

Linear	K-NN
Rigid (low variance)	flexible (high variance)
Eager	Lazy
Parametric	Instance-Based

Bayes Error Rate If we know the density $P(x, y)$, we classify the most probable class, using the conditional distribution as: output the class v | is $\max P(v | x)$.

The error rate of this classifier (called the Bayes classifier) is called the Bayes error rate: the minimum achievable error rate given the distribution of the data. K-NN directly approximates this solution (majority vote in a nearest neighborhood) except that conditional probability is relaxed to conditional probability withing a neighborhood and probabilities are estimated by training sample proportions

Inductive bias of K-NN The assumed distance tells us which are the most similar examples. The classification is assumed similar to the classification of the neighbors according to the assumed metric.

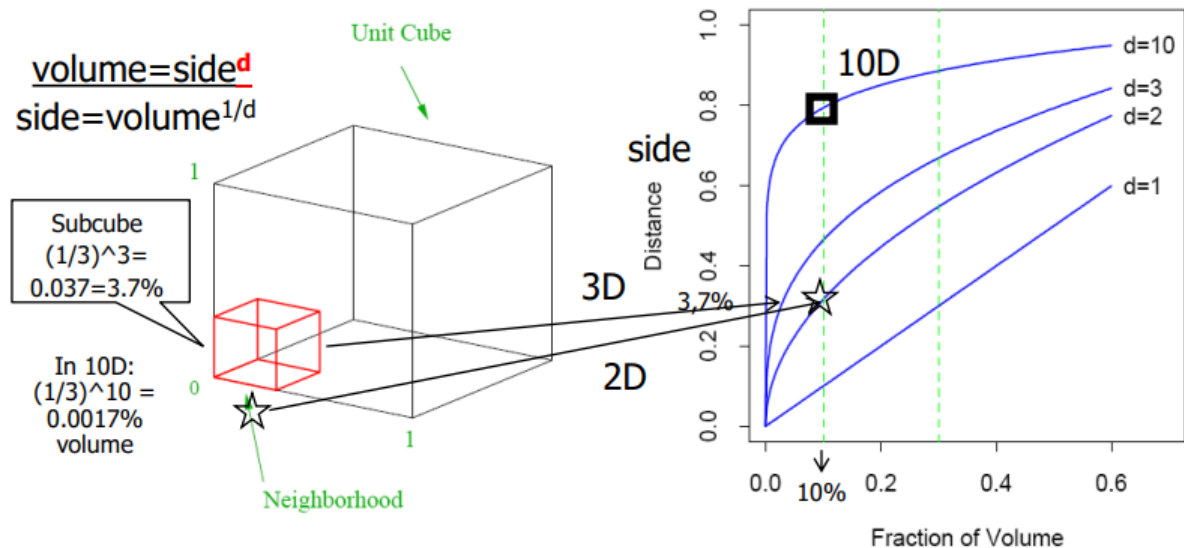
Limitations The computational cost is deferred to the prediction phase: makes the local approximation to the target function for each new example to be predicted.

Moreover, high retrieval cost: computationally intensive, in time, for each new input because computes the distance from the test sample to all stored vectors, so high cost in space too (all training data).

It provides a good approximation if we can find a significant set of data close to any x . Can fail when we have a lot of input variables (high n , high dimensionality) due to the curse of dimensionality:

Hard to find nearby points in high dimensions

K-NN can fail in high dimensions because it becomes difficult to gather K observation close to a target point x_q : near neighborhoods tend to be spatially large and estimates are no longer local



Low sampling density for high-dim data

Sampling density is proportional to $l^{\frac{1}{d}}$ with l data and d data dim. If 100 points are needed to estimate a function in R (1 dim), then 100^{10} are needed in R^{10}

Irrelevant features: the **curse of noisy**

If the target depends only on a few features in x , we could retrieve a similar pattern with the similarity dominated by the large number of irrelevant features.

This grows with dimensionality.

We may weight features according to the relevance, or adopt feature selection approaches (eliminating some variables)

1.2 Neural Networks

Models used to:

Study and model biological systems and learning processes (biological realism is essential)

Introduce effective machine learning systems and algorithms (often losing a strict biological realism, but machine learning, computational and algorithmic properties are essential)

For us: **Artificial Neural Networks** (ANN): a flexible machine learning tool in the sense of approximating functions (builds a mathematical function $h(x)$ with special properties). A neural network:

Can learn from examples

Are universal approximators (**Theorem of Cybenko**): flexible approaches for arbitrary functions (including non-linear)

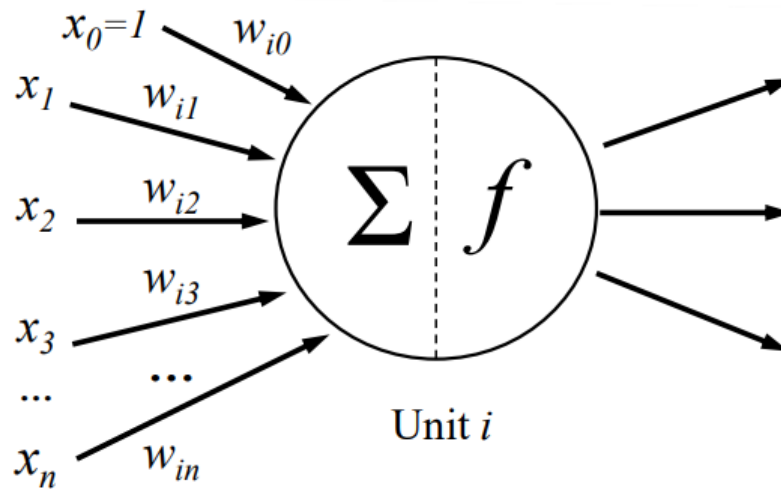
Can deal with noisy and incomplete data, with a graceful degradation of performance

Can handle continuous real and discrete data for both regression and classification tasks

It's a **paradigm**: it encompasses a wide set of models.

1.2.1 Artificial Neuron

Input from external source or other units, with weights w as free parameters: can be modified by the learning process. The unit i computes $f(\sum_j w_{ij}x_j)$ with w_{ij} the weight from input j to unit i . f is called **activation function**: linear, threshold or logistic (sigmoid). The weighted sum ($\sum_j w_{ij}x_j$) is called net input to unit i , or net_i .



Three types of activation functions:

Linear, or identity: $f(x) = x$

Threshold, for the **perceptrons**: $f(x) = \text{sign}(x)$

Logistic: $f(x) = \frac{1}{1+e^{-\alpha x}}$

Perceptron

A neuron that uses a threshold as activation function. Can be composed and connected to build a network: **MLP**, Multi Layer Perceptron

Xor $x_1 \oplus x_2 = x_1 \cdot \overline{x_2} + \overline{x_1} \cdot x_2$. Let $h_1 = x_1 \cdot x_2, h_2 = x_1 + x_2$ then $x_1 \oplus x_2 = \overline{h_1} \cdot h_2$ with $\wedge = \cdot$ and $\vee = +$. So two layers are sufficient, but single layer cannot model all functions due to limits of single perceptron and the linear separable problems.

Learning for one unit model

1. **Adaline**, adaptive linear neuron: LMS direct solution and gradient descent solution
2. **Perceptron**, non linear: only classification
Minimize number of misclassified patterns, find $w \mid \text{sign}(w^T x) = d$. Online algorithm, a step can be made for each input pattern.
 - (a) Initialize weights
 - (b) Pick learning rate η (between 0 and 1)
 - (c) Until stopping condition (es weights don't change)
For each training pattern (x, d) compute output activation $out = \text{sign}(w^T x)$.
If $out = d$ don't change weights, if $out \neq d$ update weights $w_{new} = w + \eta \cdot d \cdot x$ adding $+\eta x$ if $wx \leq 0$ and $d = +1$ or $-\eta x$ if $wx > 0$ and $d = -1$.
Different form $w_{new} = w + \frac{1}{2} \cdot \eta \cdot (d - out) \cdot x$

Delta Rule: the form $w_{new} = w + \eta \cdot d \cdot x$ is the **Hebbian learning** form, while the form $w_{new} = w + \frac{1}{2} \cdot \eta \cdot (d - out) \cdot x$ is the **error correction learning** form. It's a recall from LMS: a error correction/delta/Widrow-Hoff/Adaline/LMS rule that changes the w proportionally to the error (target d – output).

In terms of neurons, the adjustment made to a synaptic weight is proportional to the product of error signal and the input signal that excite the synapse. Easy to compute when errors signal δ is directly measurable (meaning that we know the desired response for each unit).

Perceptron Convergence theorem A perceptron can represent linear decision boundaries, so it can solve linearly separable problems. Also, it can always learn the solution with the perceptron learning algorithm. The **perceptron convergence theorem** is a milestone: a biologically inspired model with well-defined and proved computational capabilities and proved by a theorem. It states that **the perceptron is guaranteed to converge** (classifying correctly all the input patterns) **in a finite number of steps if the problem is linearly separable**. This independently of the starting point, although the final solution is not unique and it depends on the starting point.

Preliminaries We focus on positive patterns, assuming $(x_i, d_i) \in \text{TR set}$ with $d_i = +1$ or -1 . We also omit T for the dot product $w^T x$ (sometimes).
 Linearly separable $\Rightarrow \exists w^*$ solution $| d_i(w^* x_i) \geq \alpha = \min_i d_i(w^* x_i) > 0$, hence $w^*(d_i x_i) \geq \alpha$.
 With $x'_i = (d_i x_i)$ then w^* solution $\Leftrightarrow w^*$ solution of $(x'_i, +1)$
 This because w^* solves $\Rightarrow d_i(w^* x_i) \geq \alpha \Rightarrow (w^* d_i x_i) \geq \alpha \Rightarrow (w^* x'_i) \geq \alpha \Rightarrow w^*$ solution of $(x'_i, +1)$.
 And if w^* is a solution of $(x'_i, +1) \Rightarrow (w^* d_i x_i) \geq \alpha \Rightarrow d_i(w^* x_i) \geq \alpha \Rightarrow w^*$ solves for x_i

Also assuming $w(0) = 0$ (at step 0), $\eta = 1$ and $\beta = \max_i |x_i|^2$ where $||$ is the euclidean norm.
 After q errors (all false negatives), $w(q) = \sum_{j=1}^q x_{ij}$ with ij denoting the patterns belonging to the subset of misclassified patterns.

Proof The basic idea is that we can find lower and upper bound to $|w|$ as a function of q^2 steps (lower bound) and q steps (upper bound) \Rightarrow we can find the number of steps q | the algorithm converges.
 Lower bound on $|w(q)|$ is $w^* w(q) = w^* \sum_{j=1}^q x_{ij} \geq q\alpha$ recalling that $\alpha = \min_i (w^* x_i)$. With Cauchy-Schwartz we know that $(wv)^2 \leq |w|^2 |v|^2$ where $|w|^2 = ||w||_2^2$.
 $|w^*|^2 |w(q)|^2 \geq (w^* w(q))^2 \geq (q\alpha)^2 \Rightarrow |w(q)|^2 \geq \frac{(q\alpha)^2}{|w^*|^2}$. Also $|w(q)|^2 = |w(q-1) + x_{iq}|^2 = |w(q-1)|^2 + 2w(q-1)x_{iq} + |x_{iq}|^2$ because $|a+b|^2 = |a|^2 + 2ab + |b|^2$
 For the q -th error, $2w(q-1)x_{iq} < 0$, so $|w(q)|^2 \leq |w(q-1)|^2 + |x_{iq}|^2$ and by iteration $w(0) = 0$ we have $|w(q)|^2 \leq \sum_{j=1}^q |x_{ij}|^2 \leq q\beta$ with $\beta = \max_i |x_i|^2$.

So we have an upper bound $q\beta$ and a lower bound $\frac{(q\alpha)^2}{|w^*|^2}$, so

$$q\beta \geq |w(q)|^2 \geq \frac{(q\alpha)^2}{|w^*|^2}$$

$$q\beta \geq q^2 \alpha'$$

$$\beta \geq q\alpha'$$

$$q \leq \frac{\beta}{\alpha'}$$

Differences

$$w_{new} = w + \eta(d - \text{out})x$$

Apparently similar but:

LMS algorithm

LSM rule derived without threshold activation functions

Hence for training $\delta = d - w^T x$

Can still be used for classification using $h(x) = \text{sign}(w^T x)$, LTU

Minimizes $E(w)$ with $\text{out} = w^T x$

Asymptotic convergence also for not linear separable problems

Not always zero classification errors

Can be extended to network of units (NN) using the gradient based approach

Perceptron learning algorithm

Perceptron uses $\text{out} = \text{sign}(w^T x)$

versus $\delta = d - \text{sign}(w^T x)$

Minimizes misclassifications ($\text{out} = \text{sign}(w^T x)$)

Always converges in a finite number of steps for a linear separable problem to a perfect classifier
 Else it doesn't converge

Difficult to be extended to network of units (NN)

Activation functions

Linear, or identity: $f(x) = x$

Threshold, for the **perceptrons**: $f(x) = \text{sign}(x)$

Logistic: $f(x) = \frac{1}{1+e^{-\alpha x}}$

This is a non-linear squashing function like the sigmoidal logistic function: it assumes a continuous range of values in the bounded interval $[0, 1]$. It has the property of being a smoothed differentiable threshold function, with α being the slope parameter of the sigmoid function.

Radial basis: $f(x) = e^{-\alpha x^2}$

Softmax

Stochastic neurons, where the output is +1 with probability $P(\text{net})$ or -1 with $1 - P(\text{net}) \Rightarrow$ Boltzmann machines and other models rooted in statistical mechanics.

For the derivatives, a step function has no derivative, which is exactly why it isn't used. The sigmoids have

$$\frac{df_\sigma(x)}{dx} = f_\sigma(x)(1 - f_\sigma(x)) \text{ and } \frac{df_{\tanh}(x)}{dx} = 1 - f_{\tanh}(x)^2 \text{ for } \alpha = 1.$$

The sigmoid logistic function has the property of being a smoothed *differentiable* threshold function. Hence, we can derive a Least (Mean) Square algorithm, by computing the gradient of the mean square loss function as for the linear units (also for a classifier).

From $\sigma(x) = x^T w$ to $\sigma(x) = f_\sigma(x^T w)$ where f_σ is a logistic function. Find w that minimizes the residual sum of squares

$$\begin{aligned} E(w) &= \sum_p (d_p - \sigma(x_p))^2 = \sum_p (d_p - f_\sigma(x_p^T w))^2 \\ \frac{\partial E(w)}{\partial w_j} &= -2 \sum_p (x_p)_j (d_p - f_\sigma(x_p^T w)) f'_\sigma(x_p^T w) = \\ &= -2 \sum_p (x_p)_j \delta_p f'_\sigma(\text{net}(x_p)) \text{ for 1 unit and patterns } p = 1 \dots l \end{aligned}$$

Gradient descent algorithm The same as for linear units using the new delta rule $w_{\text{new}} = w + \eta \cdot \delta_p \cdot x_p$ where $\delta_p = (d_p - f_\sigma(\text{net}(x_p))) f'_\sigma(\text{net}(x_p)) = (d_p - \text{out}(x_p)) f'_\sigma$

Neural Network In an MLP architecture: units connected by weighted links, organized in layers:

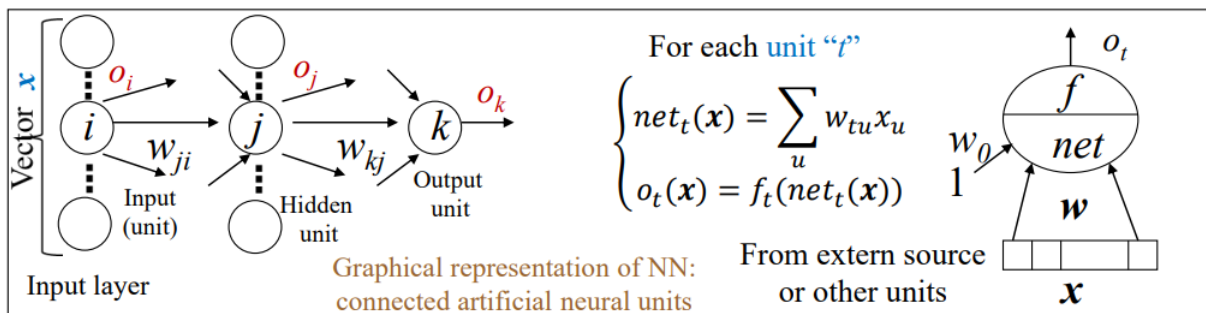
input layer, source of the input x . Copies the source external input patterns x , without computing net and f .

hidden layer, projects onto another hidden or an output layer, computing.

Can be viewed as network of units o flexible function:

$$h(x) = f_k \left(\sum_j w_{kj} f_j \left(\sum_i w_{ji} x_i \right) \right)$$

As for notation, given



We have that the index t denotes a generic unit, can be either k or j . u denotes a generic input component, either i or j .

x is a generic input from an external source (input vector) or from other units according to the position of the unit in the network. If we load the pattern x in the input layer, we can use the notation with o for both the inputs and the hidden units outputs. Hence, inside the network, the input to each unit t from any source u (through the connection w_{tu}) is typically denoted as o_u .

Architectures

Feedforward NN Standard architecture of the NNs. Direction: input \longrightarrow output

Recurrent NN adds feedback loops connections to the network topology. These self-loop connections provide the network with dynamical properties, leaving a "memory" of the past computations inside the model. This allows us to extend the representation capability of the model to the processing of sequences and structured data.

Flexibility The hypothesis space is a **continuous space** of all the functions that can be represented by assigning the weight values of the given architecture. Depending on the class of values produced by the output units (discrete or continuous), the model can deal, respectively, with classification tasks (sigmoidal output f) or regression tasks (linear output f). Also multi-class and multi-regression, with multiple output units.

Universal approximation The flexibility is theoretically grounded (Cybenko 1989, Hornik et al. 1993...). In short, a single hidden-layer network with logistic activation functions can approximate (arbitrary well) every continuous function, provided enough units in the hidden layer.

A MLP network can approximate (arbitrarily well) every input-output mapping, provided enough units in the hidden layers.

Existence theorem: given ϵ , $\exists h(x) \mid |f(x) - h(x)| < \epsilon \ \forall x$ in the hypercube.

With this fundamental result (MLP can represent *any* function), two issues arise: how to learn by neural network and how to decide its architecture.

The **expressive power** of a NN is strongly influenced by the number of units and their configuration (architecture). The number of units can be related to the discussion of the VC-dim, specifically: the network capabilities are influenced by the number of parameters w^* , which is proportional to the number of units. Further studies also report the dependencies on their value sizes.

Es: weights = 0 \Rightarrow minimal VC-dim, small weights \Rightarrow linear part of the activation function, high values weights \Rightarrow more complex model.

The universal approximation theorem is a fundamental contribution, showing that one hidden layer is sufficient in general, but it doesn't assure that a "small number" of units would do the work. For many function families it's possible to find boundaries on that number, but also cases for which a single hidden layer network would require an exponential number of units (in n input dimension).

More layers can help, but is it possible to efficiently train deep networks?

Learning Algorithm The learning algorithm allows adapting the weights w of the model in order to obtain the best approximation of the target function. Often realized in terms of minimization of an error/loss function on the training dataset. Same problem as other models: given a set of l training examples (x_p, d_p) and a loss measure L (ex. the MSE $L(h(x_p), d_p) = (d_p - h_w(x_p))^2$), find the weight vector w that minimizes the expected error on the training data

$$E(w) = R_{emp} = \frac{1}{l} \sum_{p=1}^l L(h(x_p), d_p)$$

Credit assignment problem: which credit to the hidden units? Not easy when the error signal is not directly measurable: we don't know the error (delta) or the desired response for the hidden units, which is useful for changing their weights. Supposed too difficult, but the backpropagation algorithm brought a renaissance of the NN field.

Loading problem, NP-complete: given a network and a set of examples, is there a set of weights so that the network will be consistent with the examples?

In practice, networks can be trained in a reasonable amount of time, although optimal solutions are not guaranteed. How to solve? Key steps:

Credit assignment problem: how to change the hidden layer weights?

Gradient descent approach can be extended to MLP (provided that loss and activations are differentiable functions, to find the delta for every unit in the network)

Backpropagation algorithm We need a differentiable loss, differentiable activation functions and a network to follow the information flow.

Find w by computing the gradient of the error function

$$E(w) = R_{emp} = \frac{1}{l} \sum_{p=1}^l (h(x_p) - d_p)^2$$

It has nice properties: easy because of the **compositional** form of the model, and keeps track only of the quantities local of each unit (local variables), so the **modularity** of the units is preserved.

```

1 def backprop():
2     # 1. Initialize all the weights w in the network and eta
3     # 2. Compute out and e_tot
4     while e_tot < epsilon: # with epsilon desired value or other criteria
5         for w_i in w:
6             d_w_i = - (gradient of e_tot respect to w_i) # step (1)
7             w_new = w + eta*d_w_i + ... # step (2)
8         # Compute out and e_tot
9     end

```

Step (1)

$$\Delta w = -\frac{\partial E_{tot}}{\partial w} = -\sum_p \frac{\partial E_p}{\partial w} = \sum_p \Delta_p w$$

Issues in training neural networks Heuristic guidelines in setting the backward propagation (backprop) algorithm. Generally, the models are over-parametrized, the optimization problem is not convex and is potentially unstable. We will discuss few of the issues. A good interpretation is to see backprop as a path through the loss/weight space. The path depends on: data, neural network, starting point (initial weight values), rate of convergence, final point (stopping rule). This defines a control for the search over the hypothesis space. The basic algorithm, once again, is:

1. Start with weight vector $w_{initial}$ and fix $0 < \eta < 1$
2. Compute $\Delta w = -\text{gradient of } E(w) = -\frac{\partial E(w)}{\partial w}$ (or for each w_i)
3. Compute $w_{new} = w + \eta \cdot \Delta w$ (or for each w_i)
 η is the step size or **learning rate**
4. Repeat from 2 until convergence or $E(w)$ sufficiently small

But now the $\Delta w = -\frac{\partial E(w)}{\partial w}$ were obtained through backprop derivation/algorithm for any weight in the network. At step 2, to compute the error, we first apply inputs to the network computing an output (**forward phase**), then we retro-propagate the deltas for the gradient (**backward phase**). So how to choose the $w_{initial}$, η and the convergence?

Starting values In the basic algorithm, $w_{initial}$. The weights are initialized with random values close to zero. To be avoided: all zero, high values or all equal (symmetry), because this hampers the training. For standardized data, values in $[-0.7, +0.7]$. There are other heuristics: $range \cdot \frac{2}{fanin}$ with $fanin$ being the number of inputs to a hidden unit, or orthogonal matrices...

Multiple minima The loss is not convex, has local minima. This affects the results, which depends on the starting weight values, hence: try a number of random starting configurations (5-10 or more **training runs** or **trials**). Useful taking the mean results (mean of errors) and looking at the variance to evaluate the model and then, if only one response is needed: we can choose the solution giving the lowest (penalized) validation error or we can take advantage of different end points and outputting a mean of the outputs (**committee approach**).

A "good" local minima is often sufficient: in ML we don't need the global or local minimum on R_{emp} , as we are searching the minimum of R (which we can't compute). Often we stop early, in a point of non-zero gradient so being neither a local nor a global minima for the training error. The neural network builds a variable size hypothesis space, so VC-dim increases during training and the training error decreases toward zero (or global minimum) while the neural network becomes too complex. We **stop before this condition of overtraining**, avoiding overfitting.

Online/batch Batch version: sum all the gradients of each pattern over an epoch and then update the weights (Δw after each epoch of l patterns). Online/stochastic: upgrade w for each pattern p , making progress with each example it sees. Faster but need smaller η .

Batch

1. Start with weight vector $w_{initial}$ and fix $0 < \eta < 1$
 2. Compute $\Delta w = -$ gradient of $E(w)$ on the entire training set (**epoch**)
 3. Compute $w_{new} = w + \eta \cdot \Delta w$ (or for each w_i)
 4. Repeat from 2 until convergence
1. Start with weight vector $w_{initial}$ and fix $0 < \eta < 1$
 2. For each pattern p
 - (a) Compute $\Delta_p w = -$ gradient of $E(w)$
 - (b) Compute $w_{new} = w + \eta \cdot \Delta_p w$ η is the step size or **learning rate**
 3. Repeat from 2 until convergence

More accurate estimation of gradient **Online**

Since the gradient of a single data point can be considered a noisy approximation to the overall gradient, this is also called stochastic gradient descent.

Many variations exists, for example stochastic gradient descent with minibatch.

Learning rate With batch training we have more accurate estimation of gradient, higher η . With online we have training faster but potentially unstable, lower η . So high vs low $\eta \Leftarrow$ fast but unstable vs slow but stable. Typically $\eta \in [0.01, 0.5]$.

The learning curve, plotting the errors during training, allows to check the behavior in the early phases of the model design. Of course the absolute value depends also on model capability and other hyperparameters, but η plays a big role in the curve quality. It's useful to have the mean of the gradients over the epoch: uniform approach (Least **Mean** Square). Some improvements: momentum (Nesterov), variable and adaptive learning rates, varying depending on the layers (in deep networks)...

With momentum it becomes $\Delta w_{new} = -\eta \frac{\partial E(w)}{\partial w} + \alpha \Delta w_{old}$, saving Δw_{new} in Δw_{old} for the next step. Becomes faster in plateaus but damps in oscillations (inertia effect, allows higher η)

Can be used in online by considering the previous example, Δw_{p-1} as Δw_{old} .

It smooths the gradient over different examples. A variant is to evaluate the gradient after the momentum is applied (so using $\bar{w} = w + \alpha \Delta w_{old}$), improves the rate of convergence for the batch mode (not online!).

The variable learning rate starts high and decays linearly for each step until iteration τ , using $\alpha = \frac{s}{\tau}$ with s current step so that $\eta_s = (1 - \alpha)\eta_0 + \alpha\eta_\tau$, then stops and uses a fixed small η_τ . Set up as $\eta_\tau = 1\%$ of η_0 and τ as few hundred steps. η_0 same no instability/no stuck trade off.

With adaptive learning rate, it's automatically adapted during training possibly avoiding/reducing the fine tuning phase via hyperparameters selection.

Stopping criteria When to stop training? Basic: error. The best metric if we know the tolerance of data. Other: max tolerance instead of mean, number of misclassified, no more relevant weight changes or no more significant error decreasing.

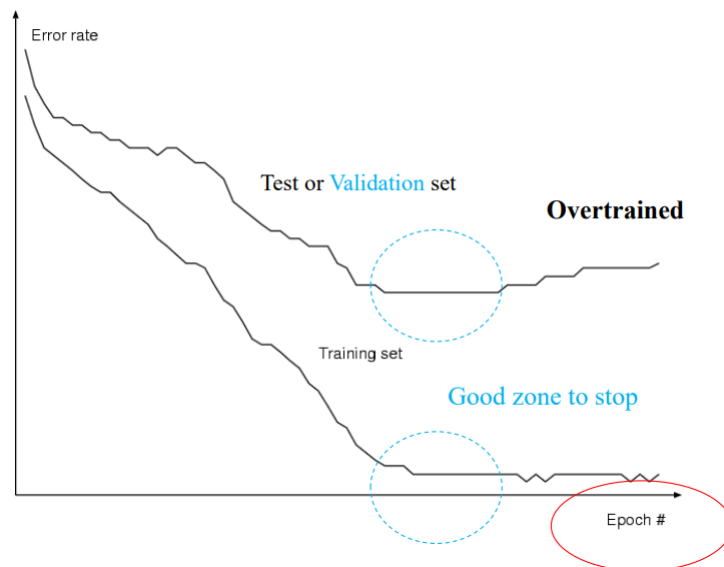
In any case stop after too many epochs, but avoid stopping at an arbitrary fixed number of epochs, and not necessarily stop with very low training error.

Control of complexity is the main aim to achieve best generalization capability.

Overfitting and regularization Typically, stopping at the global minimum of $R_{emp}(w)$ is likely to be an overfitting solution. The control of complexity is our main aim to achieve the best generalization capability. For instance,

we need to add some **regularization**: can be achieved directly (**penalty term**) or indirectly (**early stopping**). Model selection with cross validation on empirical data to find the trade-off.

In neural networks, we start learning with small random weights (breaking the symmetry!). As optimization proceeds, hidden units tend to saturate, increasing the effective number of free parameters (hence increasing the VC-dim). As we discussed, this is a variable-sized hypothesis space (changes during training).



How to act on the overtraining?

1. Early stopping: using a validation set to determine when to stop (vague indication: when the validation error increases, so use more than one epoch before estimating (**patience**))

Since the effective number of parameters grows during the training, halting the process effectively limits the complexity.

2. Regularization on the loss: we can optimize the loss considering the weights values.

Related to Tikhonov, so well principled approach: add a penalty term to the error function: $Loss(w) = \sum_p (d_p - f(x_p))^2 + \lambda ||w||^2$ with $||w||^2 = \sum_i w_i^2$

The effect is a weight decay, basically $w_{new} = w + \eta \cdot \Delta w - 2\lambda w$.

λ is the **regularization parameter**, generally very low (0.01) and selected in the model selection phase. Applied on the linear model it's the ridge regression.

More sophisticated penalty terms have been developed (ex: weight elimination, Haykin). Misunderstandings:

Regularization is not a technique to control the stability of the training convergence but controls the **complexity of the model**, measured by VC-dim and related to the number of weights and values of the weights in the neural networks

Early stopping needs a validation set to decide when to stop, which sacrifices some data. The regularization is a principled approach, as it allows the VL curve to follow the TR curve so that early stopping is not needed.

But you can use both!

Typically the bias (w_0) is omitted because its inclusion causes the results to be not independent from target shift/scaling. May be included with its own regularization coefficient.

Typically it's applied in the batch version. So for online/mini batch we need to take into account possible effects over many steps, so it's better to use $\lambda \cdot \frac{mb}{l}$ with l being the number of total patterns.

Other techniques: **dropout**.

3. Pruning methods

Number of units This is related to the control of complexity but also to the input dimension and the size of the TR set. In general, this is a model selection issue. The number of units, along with the regularization parameters, can be selected with the model selection phase (cross-validation).

Too few units \Rightarrow underfitting, viceversa too many units \Rightarrow overfitting. The number can be high with proper regularization.

Constructive approaches: the learning algorithm decides the number of hidden units, starting with small networks and adding new units.

Incremental approach: algorithms that build a network starting with a minimal configuration and add new units and connections during training. Examples: Tower, Tiling Upstart for classification and Cascade Correlation for both regression and classification.

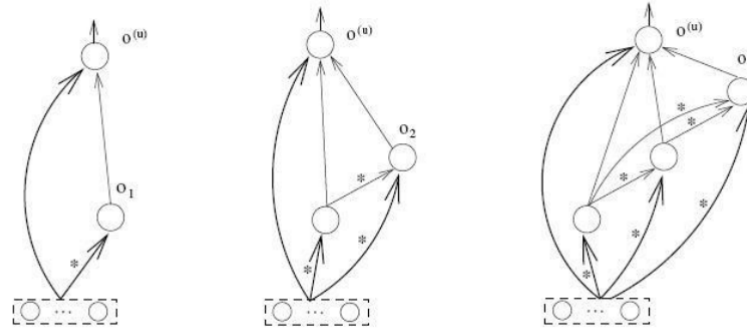
The **Cascade Correlation** algorithm starts with N0, a network without hidden units, which is trained and evaluated. If N0 cannot solve the problem, go to N1: a hidden unit is added such that the correlation between the output of the unit and the residual error of N0 is maximized (by training its weights).

After training, the weights of the new unit are frozen and the remaining weights are restrained. If the obtained network N1 cannot solve the problem, new hidden units are progressively added, which are connected with all the inputs and previously installed units. The process continues until the residual errors of the output layer satisfy a specified stopping criteria.

This method dynamically builds up a neural network and terminates once a sufficient amount of hidden units has been found to solve the given problem. Specifically, it works by interleaving the minimization of the total error function (LMS) by simple backpropagation training of the output layer and the maximization of the (non-normalized) correlation (the covariance) of the new inserted hidden (candidate) units with the residual error. With $E_{p,k} = (o_{p,k} - d_{p,k})$ residual error, with $o_{p,k}$ output, p pattern and k output unit

$$S = \sum_k \left| \sum_p (o_p - \text{mean}_p(o_{p,k})) (E_{p,k} - \text{mean}_p(E_{p,k})) \right|$$

$$\frac{\partial S}{\partial w_j} = \sum_k \text{sign}(S_k) \sum_p (E_{p,k} - \text{mean}_p(E_{p,k})) f'(net_{p,h}) I_{p,j} \text{ with } h \text{ candidate index}$$



© A. Micheli 2003

* = weights frozen after candidate training

The role of hidden units is to reduce the residual output error (solves a specific sub-problem and becomes a permanent "feature detector").

Typically, since the maximization of the correlation is obtained using a gradient descent technique on a surface with several maxima, a pool of hidden units is trained and the best one selected to avoid local maxima. It's also greedy: easy to converge may also find a minimal number of units but may lead to overfitting.

Pruning methods: start with a large network and progressively delete weights or units.

Input scaling and output representation Preprocessing can have large effects: normalization (via standardization and rescaling), categorical inputs, handling missing data. . .

For the output, one or more linear units for regression. For classification, one unit (binary classification) or one of k (multioutput):

sigmoid (choose the threshold to assign to class)

rejection zone

one-of- k encoding (winner class chosen by taking the highest value among the outputs)

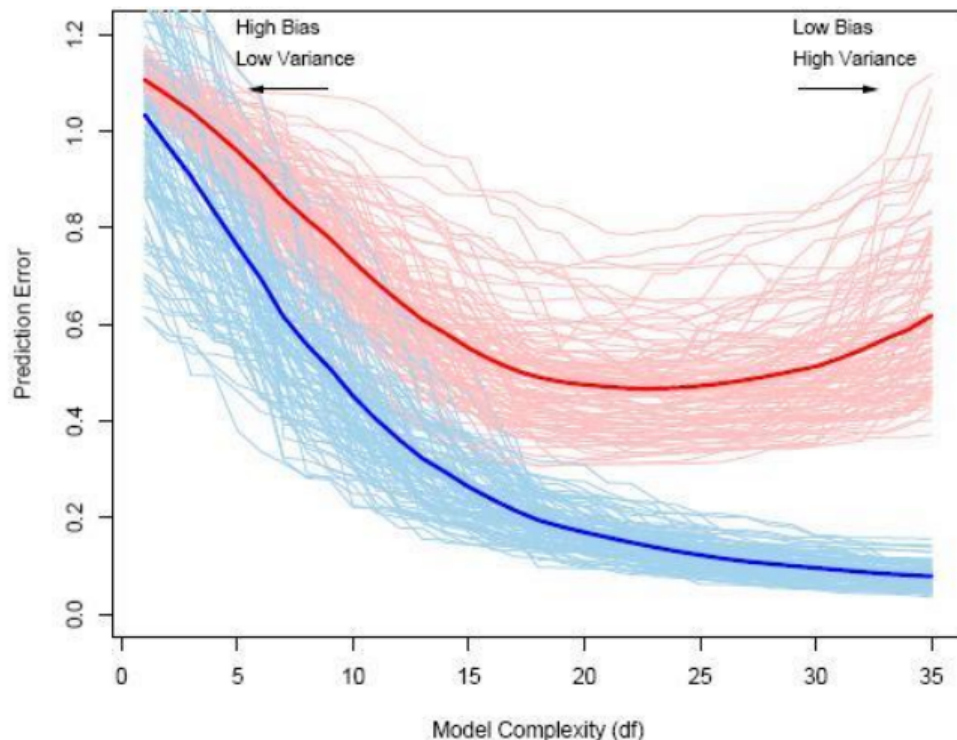
Often symmetric logistic (tanh) learn faster

Softmax for 0/1 targets, sum to 1 can be interpreted as probability of belonging to class i among k classes.

1.3 Model Selection and Model Assessment

1.3.1 Bias-Variance

A bias-variance decomposition provides an useful framework to understand the validation issue, showing how difficult is the estimation of a model's performance, again showing the need of a trade-off between fitting capabilities (**bias**) and model flexibility (**variance**).



1.3.2 Motivations

We are looking for the best solution, with minimal test error, looking for a balance between fitting (accuracy on training data) and model complexity. The **training set is not a good estimate of test error**.

Assuming that we have a set of tuning parameters Θ , implicit or explicit, that varies the model complexity, we wish to find the value of Θ that minimizes test error: **methods for estimating the expected error** for a model (or each model of a class of models, or event a set of models...)

1.3.3 Validation

We can approximate the estimation analytically:

AIC, BIC, Bayesian Information Criterion (limited to linear models)

MDL

Structural Risk Minimization and VC-Dimension

In practice, we can approximate the estimation on the data set by resampling, a direct estimation of error via cross-validation (hold-out, K-fold...), bootstrap...

Two aims

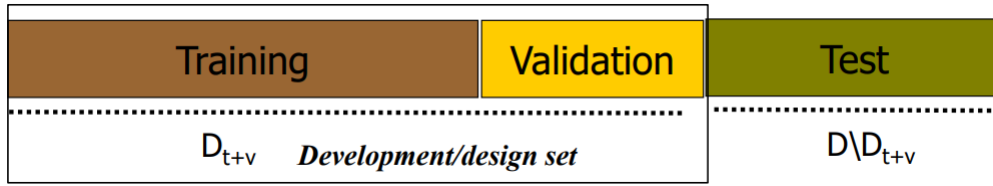
Model Selection: estimating the performance of different learning models in order to choose the best one (to generalize), which includes looking for the best hyperparameters of the model.

It **returns a model**.

Model Assessment: after choosing the final model (or class of models), estimating/evaluating the generalization error on **new test data**.

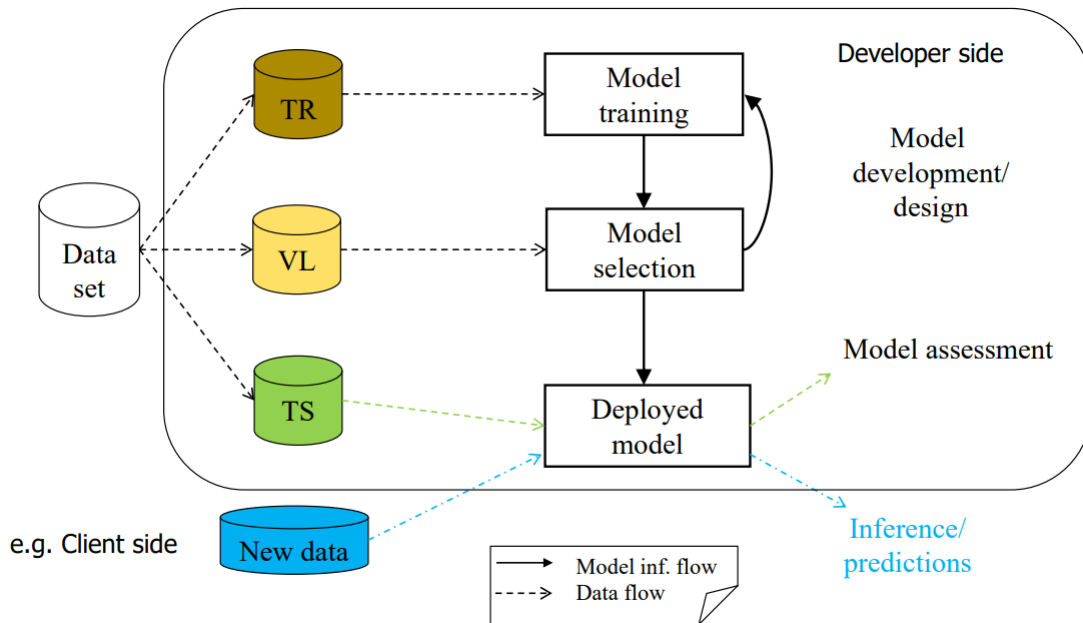
It **returns an estimation value**.

The gold rule is to keep a separation between goals and to use separate data sets: **hold out**, if data set is sufficient, for example 50% TR, 25% VL and 25% TS (**disjoint sets!**)



The TR is the **training set**, VL is the **validation/selection set** (used to perform the **model selection**) and the TS is the **test set** (used for the model assessment).

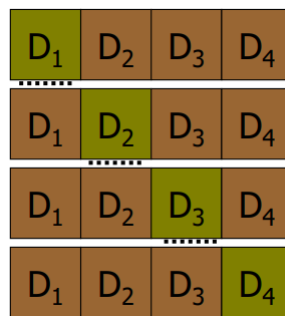
If the test set is used repeatedly in the design cycle we would be doing model selection, and not reliable assessment. **Blind test set** concept: *if you see the solution it's not a test*. In that case the test error would be an overoptimistic evaluation of the true test error. It's very easy to obtain very high classification accuracy over random tasks even when using the test set only implicitly.



Grid Search For the model selection: hyperparameters can be set by searching in a hyperparameters space. Try every hyperparameter-value combination and record the result. The best result will corresponds to the best hyperparameter values.

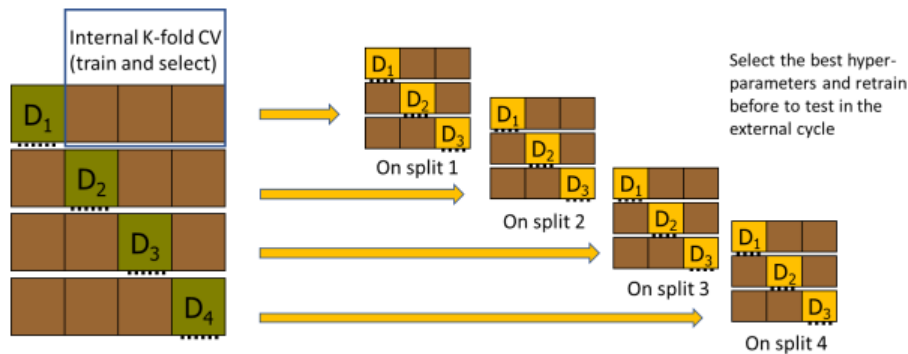
The cost can be high (cartesian product between sets of values per each hyperparameter: what if 6 hyperparameters each with tens of possible different values?), so can be useful to fix some hyperparameter value before: more levels of grid search, where we do a first coarse grid search to find good intervals of values, then a finer grid-search can be performed over smaller intervals.

K-fold Cross Validation Split the data set D into k mutually exclusive subsets D_1, \dots, D_k , we train on $D - D_i$ and test it on D_i . No unique model, but we get a variance (std. dev.) over different folds for the class of models or algorithm.



The issue is to find k and that it can be very computationally expensive but in can be combined with validation set, double- k -fold cross-validation. . .

Double cross-validation After dividing the data set into k mutually exclusive subsets D_1, \dots, D_k , for each D_i : do another **internal** cross-validation using the other $k - 1$ subsets.



It doesn't provide a model but an estimation of the risk, because can provide a different model per external step cycle.

Example with k -fold CV Split data in TR and TS.

For the **model selection**, use k -fold CV over TR, obtaining new TR and VL in each fold, to find the best hyperparameters of the model.

Train the final model on the whole, original, TR.

Perform the **model assessment** by evaluating it on the external TS.

There are more combination and more ways of doing this.

Particular cases

Lucky/unlucky sampling Can we avoid to be sensible to the particular partitioning of the examples?

Stratification: the process of grouping examples into relatively homogeneous subgroups. For example, partitioning in a way that every class (in a classification task) is represented in approximately the same proportions as in the full data set.

Repeated hold-out or k -fold CV: repeating the splitting with different random sampling, average the results to yield an overall estimation.

Very few data In this case it's difficult to say if the sampling is representative or not. We can use stratification. To avoid (or to be considered during evaluation): missing classes or features in training data, special classes, known outliers that can affect the results. . .

Also the blind test set can be misleading if: is from a different distribution, measured with a different scale/tolerance, uncleaned, unprocessed. . .

When to stop? Best to avoid stopping the NN training by fixing an arbitrary number of epochs: if it's too small it may be too early (underfitting), too high may be too late (overfitting), may not old for all the configuration in a cross-validation. . .

Also selecting the number of epochs by model selection it's not the best practice: better than a fixed number, but still not accurate.

Early Stopping Should be part of the model selection as well. Tricky: uses part of the data just to decide when to stop, also complex with CV.

How to select the best model with early stopping for the final retraining (ex: different stop point for each validation set in a CV)? You could take the average number of epochs, but varying the data for retraining after model selection could lead to a different (best) stop point. So it's better to consider the average of the TR error at the best point and use it to gain the same level of fitting on the retraining. Else you need a VL every time you train.

Random initialization Different $w_{initial} \Rightarrow$ different models: how to choose one for the final model? You can compute the mean and variance of error/accuracy across the different trials, so that a random case would be in this range.

1.4 Statistical Learning Theory

1.4.1 VC-dim

The VC-dim is used to provide the analytical bound we need to evaluate H . It's a measure of complexity (that is to say, capacity/expressive power/flexibility) of a class of hypothesis.

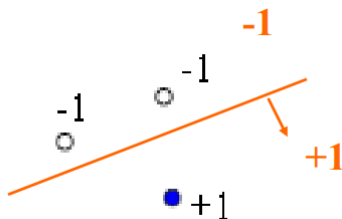
The bound is expressed not in terms of $|H|$ but in terms of distinct instances that can be completely discriminated using H . Intuitively (for classification): how much can H discriminate points? The maximum number of points that can be correctly classified/learned without error for all the possible labelings.

Shattering Given X instance/input space, n size of the instance space (number of instances) and l for the number of data available. H is the hypothesis space.

In the case of binary classification, there are 2^n possible **dichotomies** (partitions or labeling of the n points in $\{-1, +1\}$): a particular dichotomy is represented in H if there exists a hypothesis $h \in H$ that realizes the dichotomy.

Definition: H shatters $X \Leftrightarrow H$ can represent all the possible dichotomies on X (0 errors). the points in X can be separated by an $h \in H$ in all the possible ways, and for every possible dichotomy of X there exists a consistent hypothesis $h \in H$

Example



3 points in R^2 , H as a set of lines: $h(x) = \text{sign}(wx + w_0)$ with $x \in R^2$

A dichotomy is a particular labeling in $\{-1, +1\}$ of the points. This specific dichotomy can be represented in H : there exist a line that correctly separates the points (in pic)

VC-Dimension Definition: the VC dimension of a class of functions H is the maximum cardinality of a set (configuration) of points in X that can be shattered by H .

$VC(H) = p \Rightarrow H$ shatters **at least one** set of p points $\wedge H$ **cannot shatter any** set of $p + 1$ points.

If arbitrarily large but finite sets of X can be shattered by H then $VC(H) = \infty$

$VC(H) \geq 3$ shown before, note that not all the possible configurations of 3 points can be shattered (example follows), but it's sufficient to find one configuration of three points which is separable for every labeling.



VC-dim is related to the number of parameters, but it's not the same thing: we may add redundant free parameters, for example, and there exists models with one parameter and infinite VC-dim. For example, k -NN has infinite VC-dim.

Analytical Bound With N number of data

$$R[h] \leq R_{emp}[h] + \epsilon(VC\text{-dim}, N, \delta)$$

with:

Guaranteed risk $R_{emp}[h] + \epsilon(VC\text{-dim}, N, \delta)$

VC Confidence $\epsilon(VC\text{-dim}, N, \delta)$

For example, for 0/1 loss

$$\epsilon(\text{VC-dim}, N, \delta) = \sqrt{\frac{VC\left(\ln \frac{2N}{\text{VC-dim}}\right) - \ln \frac{\delta}{4}}{N}}$$

with probability at least $1 - \delta$ for every $\text{VC-dim} < N$

There are different bound formulations for different classes of functions, of tasks. . .

This gives us a way to estimate the error on future data based only on the training error and the VC-dim of H . The resulting bounds are the worst case scenario, because they hold for all but $1 - \delta$ of the possible approximation function/training sets.

Remarks For many reasonable hypothesis classes (ex: linear approximators), the VC-dim is linear in the number of free parameters of the hypothesis. This shows that to learn "well" we need a number of examples that is linear in the VC-dim.

1.4.2 Structural Risk Minimization

SRM uses VC-dim as a controlling parameter for minimizing the generalization bound on R . Assuming finite VC-dim, we can define a nested structure of models-hypothesis spaces according to the VC-dim in the following way:

$$H_1 \subseteq H_2 \subseteq \dots \subseteq H_n$$

$$VC(H_1) \leq VC(H_2) \leq \dots \leq VC(H_n)$$

Some examples:

Neural Networks with increasing number of hidden units, but also the number of epochs

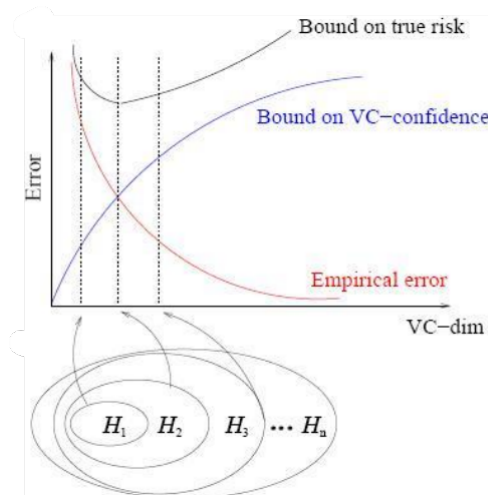
Polynomial of increasing degree

Increasing values for c , in $\|w\| < c$ for regularization

Increasing number of nodes in a decision tree

Model selection Growing VC-dim: empirical (training) error decreasing, VC confidence increasing.

SRM: find a trade-off in the bound and choose the model (h) with the best bound on the true risk.



Use of the bound Provide a fundamental theoretical ground for principled ML: independently from specific model details or learning algorithms, highlighting the role of complexity control. The optimal choice of the model complexity (structure) provides the minimum expected risk (**inductive principle of SRM**)

Also to provide a direction for new model development guided by SRM. As estimation of predictive errors is rarely used: the upper bound is overly pessimistic and may not be adequate for reliable evaluation of the generalization error (model assessment) and tighter bounds are under development, also is difficult to compute the VC-dim for specific classes of H .

Towards **principled approaches** less based on trial and error. For example, two practical approaches:

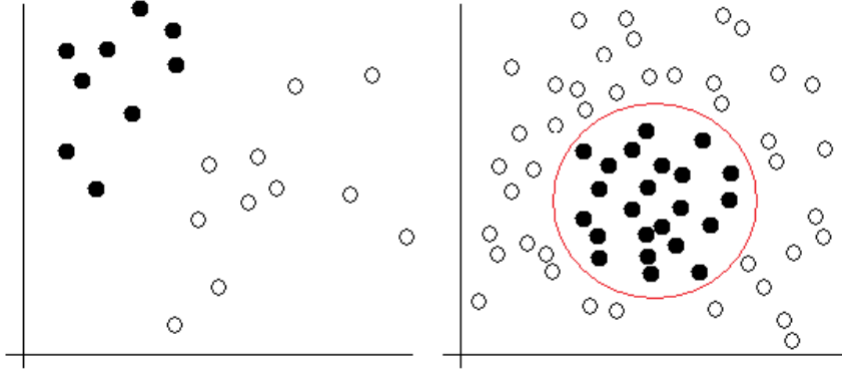
Choose appropriate structure/complexity, fix the model (hence the VC-dim) and minimize the TR error.
 Can be used in neural networks, however regularization by training heuristic can further introduce implicit SRM (early stopping... Or SRM with Tikhonov regularization where we have minimum loss with R_{emp} + complexity term, considering both the terms.

Fix the TR error, automatically minimize the VC confidence (SVM)

1.5 Support Vector Machines

Linear machine, with maximization of the separation margin and structural risk minimization. Initially **hard margin SVM**, we assume to deal with linearly separable problems without errors in the data.

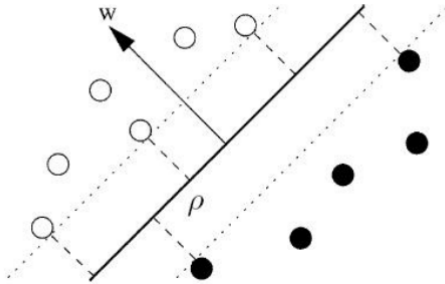
Example of linearly and non linearly separable problems:



Separating hyperplane Given the training set $T = \{(x_i, d_i)\}_{i=1}^N$ we want to find an hyperplane of equation $w^T x + b = 0$ to separate the examples and get $w^T x_i + b \geq 0$ for $d_i = +1$ and $w^T x_i + b < 0$ for $d_i = -1$

$g(x) = w^T x + b$ is the discriminant function and $h(x) = \text{sign}(g(x))$ is the hypothesis.

An example of separation margin:



In this case the hyperplane has equal distance to both the closest negative and positive examples. The separation margin (p) is evaluated as the double of the distance between the linear hyperplane and the closest data point (a "safe zone"). Not all hyperplanes solving the task are equal: the margin changes, bigger or smaller. The **optimal hyperplane is the hyperplane which maximizes p** $w_O^T x + b_O = 0$ with O being "optimal".

We will find $p = \frac{2}{\|w\|}$, so maximize $p \Leftrightarrow$ minimize $\|w\|$.

Support Vector We can rescale w and b so that the closest points to the separating hyperplane satisfy $|g(x_i)| = |w^T x_i + b| = 1$, so we can write

$$w^T x_i + b \geq 1 \text{ if } d_i = +1$$

$$w^T x_i + b < -1 \text{ if } d_i = -1$$

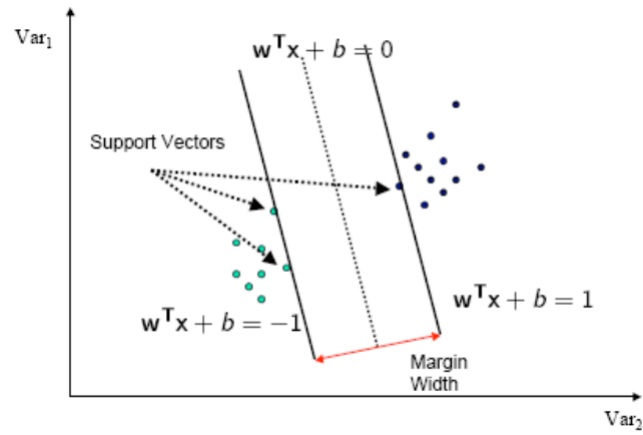
which is compact form is

$$d_i(w^T x_i + b) \geq 1 \quad \forall i = 1, \dots, N$$

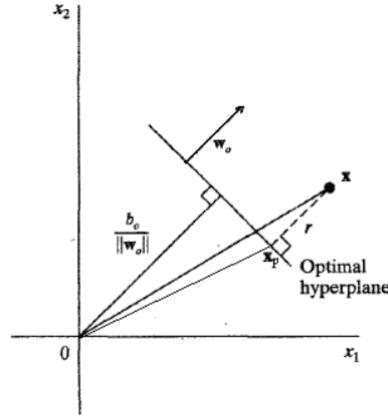
A **support vector** $x^{(s)}$ satisfies the previous equation exactly

$$d^{(s)}(w^T x^{(s)} + b) = 1$$

so the **support vectors** are the closest data points to the hyperplane and "lay on the margin".



Computing the distance With the discriminant function $g(x) = w^T x + b$, recalling that w_O is a vector orthogonal to the hyperplane. Let's denote with r the distance between x and the optimal hyperplane.



$$x = x_p + r \frac{w_O}{\|w_O\|}$$

$$g(x) = g\left(x_p + r \frac{w_O}{\|w_O\|}\right) = w_O^T x_p + b_O + w_O^T r \frac{w_O}{\|w_O\|} = g(x_p) + r w_O^T \frac{w_O}{\|w_O\|} = r \frac{\|w_O\|^2}{\|w_O\|} = r \|w_O\|$$

$$\text{thus } r = \frac{g(x)}{\|w_O\|}$$

Computing the margin Consider the distance between the hyperplane and a positive support vector $x^{(s)}$, for example

$$r \text{ for } x^{(s)} = \frac{g(x^{(s)})}{\|w_O\|} = \frac{1}{\|w_O\|} = \frac{p}{2}$$

$$\Rightarrow p = \frac{2}{\|w_O\|}$$

So the optimum hyperplane maximizes $p \Leftrightarrow$ minimizes $\|w\|$

Quadratic Optimization Problem The formal derivation of the SVM solution requires techniques in the constrained optimization framework, and we assume that quadratic programming is solved elsewhere.

For the hard margin SVM, the quadratic optimization problem asks to find the optimum values of w and b in order to maximize the margin.

Primal Form Given the training examples $T = \{(x_i, d_i)\}_{i=1}^N$, find the optimum values of w and b which minimizes

$$\psi(w) = \frac{1}{2} w^T w$$

satisfying the constraints

$$d_i(w^T x_i + b) \geq 1 \quad \forall i = 1, \dots, N$$

The objective function $\psi(w)$ is quadratic and convex in w , while the constraints are linear in w : solving this problem scales (the computational cost) with the size of the input space m .

The Lagrangian multipliers method is used: constructing the Lagrangian function corresponding to the quadratic optimization problem

$$J(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i (d_i (w^T x_i + b) - 1)$$

with $\alpha_i \geq 0$ the N Lagrangian multipliers. Each term in the sum correspond to a constraint in the primal problem. J must be minimized with respect to w and b , and maximized with respect to α , and the solution correspond to a saddle point in J .

We will find

$$w_O = \sum_{i=1}^N \alpha_{O,i} d_i x_i$$

thus the optimal hyperplane is expressed as

$$w_O^T x + b_O = 0 \Leftrightarrow \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x + b_O = 0$$

The optimal conditions will be:

$$\text{Minimize } J \text{ with respect to } w \Rightarrow \frac{\partial J}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^N \alpha_i d_i x_i$$

$$\text{Minimize } J \text{ with respect to } b \Rightarrow \frac{\partial J}{\partial b} = 0 \Rightarrow w = \sum_{i=1}^N \alpha_i d_i$$

And we may substitute these in J to study the dual form.

Kuhn-Tucker Conditions From the KT conditions follows that

$$\alpha_i (d_i (w^T x_i + b) - 1) = 0 \quad \forall i = 1, \dots, N$$

in the saddle point of J :

$$\alpha_i > 0 \Rightarrow d_i (w^T x_i + b) = 1 \text{ and } x_i \text{ is a support vector}$$

$$x_i \text{ isn't a support vector} \Rightarrow \alpha_i = 0$$

Hence we can restrict the computation to N_s so $w_O = \sum_{i=1}^{N_s} \alpha_{O,i} d_i x_i$: **the hyperplane depends solely on the support vectors!**

To obtain the Lagrangian multipliers $\{\alpha_i\}_{i=1}^N$ we must solve the dual form: given the training examples $T = \{(x_i, d_i)\}_{i=1}^N$, find the optimum values of $\{\alpha_i\}_{i=1}^N$ that maximizes

$$Q(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j d_i d_j x_i^T x_j$$

satisfying the constraints

$$\alpha_i \geq 0 \quad \forall i = 1, \dots, N$$

$$\sum_{i=1}^N \alpha_i d_i = 0$$

Then we compute $w_O = \sum_{i=1}^N \alpha_{O,i} d_i x_i$ and $b_O = 1 - w_O^T x^{(s)}$ corresponding to a positive support vector $x^{(s)}$, so

$$b_O = 1 - \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x^{(s)}$$

So we don't need to explicitly compute w_O , but we only need the Lagrangian multipliers α_i (by solving the dual problem), then we compute the optimal bias b_O . The decision surface is given by

$$w_O^T x + b_O = 0 \Leftrightarrow \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x + b_O = 0$$

So given the input pattern x :

$$\text{Compute } g(x) = \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x + b_O$$

$$\text{Classify } x \text{ as the sign of } g(x)$$

Note that it's not necessary to compute w_O , also the sum can be restricted to the number of support vector N_s

How does this improve the generalization? Minimizing the norm of w is equivalent to minimizing the VC-dim and thus to minimizing the VC confidence ϵ in

$$R[h] \leq R_{emp}[h] + \epsilon(\text{VC-dim}, N, \delta)$$

Theorem (Vapnik) Let D be the diameter of the smallest ball around the data points x_1, \dots, x_N . For the class of separating hyperplanes described by the equation $w^T x + b = 0$, the upper bound to the VC-dim is

$$\text{VC-dim} \leq \min(\lceil \frac{D^2}{p^2} \rceil, m_O) + 1$$

An elegant approach For linearly separable data there are many solutions. Vapnik proposed an "optimal separating hyperplane" maximizing the margin providing:

An unique solution with zero errors for the binary classifier

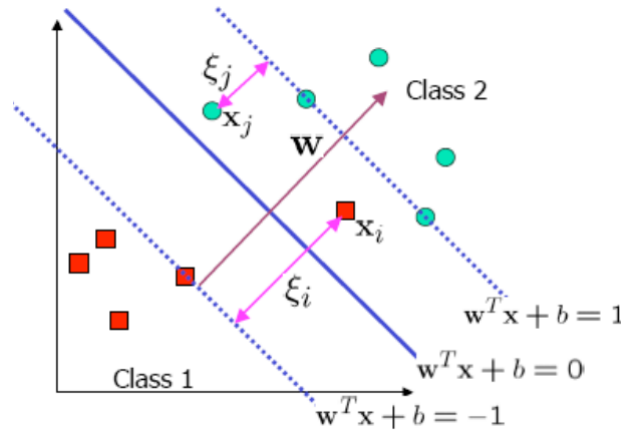
An automatized approach to SRM that minimizes VC confidence (by maximizing the margin) as part of the training process, without hyper parameters in the linear separable case.

The use of a solver in the class of constrained quadratic programming (instead of gradient descent) with a dual form (showing support vectors and dot product among the patterns)

A solution focused on "selected" training data (support vectors)

For noisy or not linearly separable data, you can have a soft margin (support vector still on the border, but some data may fall closer to the hyperplane: at least one point violate $d_i(w^T x_i + b) \geq 1$)

Soft margin SVM We introduce $\xi_i \geq 0 \ \forall i = 1, \dots, N$ called **slack variables**: $d_i(w^T x_i + b) \geq 1 - \xi_i$, so a support vector satisfies that exactly $d_i(w^T x_i + b) \geq 1 - \xi_i$



Note that Vapnik does not hold anymore. The primal problems becomes: given the training examples $T = \{(x_i, d_i)\}_{i=1}^N$, find the optimum values of w and b which minimizes

$$\psi(w, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$

satisfying the constraints

$$\begin{aligned} d_i(w^T x_i + b) &\geq 1 - \xi_i \quad \forall i = 1, \dots, N \\ \xi_i &\geq 0 \quad \forall i = 1, \dots, N \end{aligned}$$

We introduced C as a regularization hyperparameter, losing the fully automated SRM. Low C means many TR errors allowed (possible underfitting), while high C means less or no TR errors allowed (smaller margin, possible overfitting)

Kuhn-Tucker Conditions The KT conditions are now defined as

$$\forall i = 1, \dots, N \quad \alpha_i (d_i(w^T x_i + b) + \xi_i - 1) = 0$$

$\forall i = 1, \dots, N \quad \mu_i \xi_i = 0$ with μ_i Lagrangian multipliers introduced to enforce the non negativity of the slack variables.

We have $0 < \alpha_i < C \Rightarrow \xi_i = 0$ on the margin and $\alpha_i = C \Rightarrow \xi_i > 0$ inside the margin.

Solving the problem We solve the dual problem to compute the $\{\alpha_i\}_{i=1}^N$, then we find w_O and b_O from $\{\alpha_{O,i}\}_{i=1}^N$ with $w_O = \sum_{i=1}^N \alpha_{O,i} d_i x_i$. Then, for a pattern j such that $0 < \alpha_j < C$, we have

$$b_O = d_j - \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x_j$$

or an average of all the solutions for numerical stability. The non-zero Lagrangian multipliers correspond to the support vectors. We use it as before: $g(x) = \sum_{i=1}^N \alpha_{O,i} d_i x_i^T x + b_O$ and $h(x) = \text{sign}(g(x))$

1.5.1 High-Dimensional feature spaces

We can use:

Non-linear mapping of input patterns to a high-dimensional feature space

Cover's Theorem: the patterns are linearly separable with high probability in the feature space under such conditions.

Finding the optimal hyperplane to separate the patterns in the feature space.

However we know that using high dimensional feature spaces (large basis function expansion) can be computationally unfeasible and can lead to overfitting.

We will propose the kernel approach to implicitly manage the feature space while regularizing.

Kernel Non linear function mapping

$$\begin{aligned} \phi : R^{m_0} &\rightarrow R^{m_1} \\ x &\mapsto \phi(x) \end{aligned}$$

The problem is formulated as before, but the training set is now $\{(\phi(x_i), d_i)\}_{i=1}^N$ and the hyper plane is now $w^T \phi(x) + b = 0$

The weight vector is a linear combination of the feature vectors

$$w = \sum_{i=1}^N \alpha_i d_i \phi(x_i)$$

so the hyperplane equation is

$$\sum_{i=1}^N \alpha_i d_i \phi(x_i)^T \phi(x) = 0$$

Evaluating $\phi(x)$ could be intractable, but with certain conditions we don't need to evaluate it directly. We do not even need to know the feature space itself! This is possible using a function to compute directly the dot products in the feature spaces $k : R^{m_0} \times R^{m_0} \rightarrow R$, with k called inner product kernel function: $k(x_i, x) = \phi(x_i)^T \phi(x)$, and it's symmetric meaning $k(x_i, x) = k(x, x_i)$

Kernel Matrix We can arrange the dot products in the feature space between the images of the input training patterns in a $N \times N$ matrix called kernel matrix $K = \{k(x_i, x_j)\}_{i,j=1}^N$, symmetrical.

Mercer's Theorem This property holds only for kernels with positive semi-definite kernel matrices. It's related to having non-negative eigenvalues in the kernel matrix.

Properties With k_1, k_2 kernels over $R^{m_0} \rightarrow R^{m_0}$, the following are also kernel functions:

$$k_1(x, y) + k_2(x, y)$$

$$\alpha k_1(x, y) \text{ with } \alpha \in R^+$$

$$k_1(x, y) \cdot k_2(x, y)$$

Wrapping up

The training set is $T = \{(x_i, d_i)\}_{i=1}^N$

We can choose the trade-off parameter C and the kernel function k

We find $\{\alpha_i\}_{i=1}^N$ by solving the optimization problem via quadratic programming algorithms
Remember that the solution is sparse, as every α_i corresponding to non-support vector is 0

The bias b_O is computed knowing the Lagrangian multipliers and the kernel matrix

Given an input pattern x we compute $w^T \phi(x) = \sum_{i=1}^N \alpha_i d_i k(x, x_i)$
Fundamental: we don't need to compute w

x is classified as $\text{sign}(g(x)) = \text{sign}\left(\sum_{i=1}^N \alpha_i d_i k(x, x_i)\right)$ with $g(x)$ called discriminant function

At **test phase** we have the Lagrangian multipliers and the kernel matrix. To classify an unseen input pattern x , we compute $\sum_{i=1}^N \alpha_i d_i k(x, x_i)$ and classify x and the sign of that computation, so $h(x) = \text{sign}\left(\sum_{i=1}^N \alpha_i d_i k(x, x_i)\right)$

Examples of kernels

Polynomial Learning Machine $k(x, x_i) = (x^T x_i + 1)^p$ with p user-specified parameter

Radial Basis Function or **Gaussian Kernel** $k(x, x_i) = e^{-\frac{1}{2\sigma^2} \|x - x_i\|^2}$ with σ^2 user-specified parameter.

Narrow peaked kernels with small σ , imply that the reply for x_i is only d_i

Feature space with an infinite number of dimensions

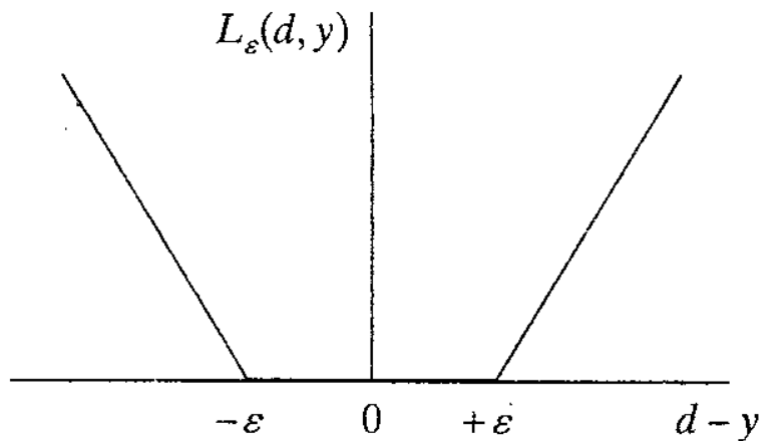
Two-layer perceptron $k(x, x_i) = \tanh(\beta_0 x^T x_i + \beta_1)$ where $\beta_0 > 0$ and $\beta_1 < 0$ are user-specified parameters
Here the Mercer's Theorem holds only for some choices of β_0, β_1

1.5.2 SVM for non-linear regression

A regression problem requires to find f such that $d = f(x) + v$ with a training set $T = \{(x_i, d_i)\}_{i=1}^N$ and with v statistically independent from x

We estimate d using a linear expansion of non-linear functions $\{\phi_j(x)\}_{j=0}^{m_1}$ so $y = h(x) = w^T \phi(x)$ where $w = (w_0 = b, w_1, \dots, w_{m_1})^T$ and $\phi(x) = (\phi(x)_0 = 1, \phi(x)_1, \dots, \phi(x)_{m_1})^T$

ϵ -insensitive loss function $L_\epsilon(d, y) = \begin{cases} |d - y| - \epsilon & \text{if } |d - y| \geq \epsilon \\ 0 & \text{otherwise} \end{cases}$



Optimization problem Introducing the non-negative slack variables ξ'_i and $\xi_i \forall i = 1, \dots, N$

$$-\xi'_i - \epsilon \leq d_i - w^T \phi(x_i) \leq \epsilon + \xi_i$$

leading to the following constraints, all $\forall i = 1, \dots, N$

$$d_i - w^T \phi(x_i) \leq \epsilon + \xi_i$$

$$w^T \phi(x_i) - d_i \leq \epsilon + \xi'_i$$

$$\xi_i \geq 0$$

$$\xi'_i \geq 0$$

The primal problem then is: given the training set $\{(x_i, d_i)\}_{i=1}^N$ find the optimal values of w such that the following objective function is minimized

$$\psi(w, \xi, \xi') = \frac{1}{2} w^T w + C \sum_{i=1}^N (\xi_i + \xi'_i)$$

under the constraints, all $\forall i = 1, \dots, N$

$$d_i - w^T \phi(x_i) \leq \epsilon + \xi_i$$

$$w^T \phi(x_i) - d_i \leq \epsilon + \xi'_i$$

$$\xi_i \geq 0$$

$$\xi'_i \geq 0$$

The dual problem yields the optimal $\{\alpha_i\}_{i=1}^N$ and $\{\alpha'_i\}_{i=1}^N$. With those, we can compute the optimal w

$$w = \sum_{i=1}^N (\alpha_i - \alpha'_i) \phi(x_i) = \sum_{i=1}^N \gamma_i \phi(x_i)$$

with $\gamma_i = \alpha_i - \alpha'_i$: in this case support vectors correspond to non-zero values of γ_i

The estimate is defined as $h(x) = y = w^T \phi(x)$, and using the linear expansion for w we get

$$h(x) = \sum_{i=1}^N \gamma_i \phi(x_i)^T \phi(x) = \sum_{i=1}^N \gamma_i k(x_i, x)$$

Wrapping up

Important: select values for the user-specified parameters C and ϵ

Important: choose an inner product kernel function k

Compute the kernel matrix K

Solve the dual form and get the optimal values of the Lagrangian multipliers ($\{\gamma_i\}_{i=1}^N$)

Compute the optimal value for the bias (b)

Obtain the estimate function as a linear combination of dot products in a feature space we can ignore ($h(x) = \sum_{i=1}^N \gamma_i k(x_i, x)$)

The test: given an input pattern x , we estimate the value of the unknown function f in that point using the estimate computed before

$$h(x) = \sum_{i=1}^N \gamma_i k(x_i, x)$$

Summary of the main characteristics

Pros

The regularization is embedded in the optimization problem (margin)

Approximation of the theoretical structure risk minimization

Convex problem (**training always finds a global minimum**)

Implicit feature transformation using kernels

Cons

Must choose the kernel and its parameters

It's a batch algorithm

Very large problems are computationally intractable

Problems with more than 20000 examples are very difficult to solve with standard approaches. However many solutions are proposed (including gradient descent approaches)

In practice There is no theory which guarantees that a given family of SVMs will have high accuracy on a given problem.

The nice property of hard-margin SVMs cannot be directly extended to soft-margin and kernels: the C parameter and kernels can lead to infinite VC-dim of the SVM classifier.

Gaussian RBF SVMs of sufficiently small width can classify an arbitrarily large number of training points correctly (only 1 SV point, the closest one, will contribute to the solution), thus have infinite VC-dim. On the opposite, with large width of the gaussian all SV points are considered and you get a sort of "global average", low VC-dim.

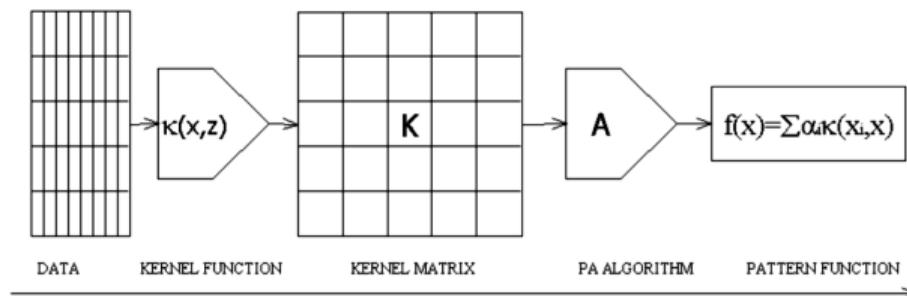
⇒ controlling the width is controlling the VC-dim.

Rigorous selection of the kernel and the C requires an estimation of the complexity (VC-dim). Hyperparameters affecting model complexity are used for model selection. In practice a careful empirical evaluation.

Nowadays, deep neural networks have largely outperformed previous records of SVMs on image and speech recognition tasks.

1.5.3 Kernel Methods

The kernel trick introduces the kernel methods, also without SVM. **Kernelization** of previous approaches, whenever you had a dot product inside your model or a similar measure enabling them to operate in a new implicit high-dimensional space just by specifying a kernel (and changing it: **the K can change without changing the learning machine.**



An SVM is largely characterized by the kernel: the best choice of kernel for a given problem is still a research issue.

1.6 Bias-Variance

A training set is only one possible realization from the universe of the data: different TR sets can provide different estimate. The expected error (on various TR) at a point x is decomposed as:

Bias: quantify the discrepancy between true function and $h(x)$ (averaged on data)

Variance: quantify the variability of the response of model h for different realizations of the TR data

Noise: error in the label

Let's assume the regression scenario with target y and L_2 squared error loss. Suppose we have examples (x, y) where the true function is $y = f(x) + \epsilon$ with ϵ being Gaussian noise with zero mean and std. dev. σ . In linear regression, given a set of examples (x_i, y_i) with $i = 1, \dots, l$, we fit a linear hypothesis $h(x) = wx + w_0$ to minimize sum-squared error over the training data $\sum_{i=1}^l (y_i - h(x_i))^2$

Because of the hypothesis class that we chose for some function f (linear hypothesis), we have a **systematic prediction error**. Depending on the dataset that we have, the parameters w that we find will be different.

Given a new data point x , what is the expected prediction error? Assume that the data points are drawn independent and identical distributed from a unique underlying probability distribution P . The goal is to compute, for an arbitrary new point x , $E_P[(y - h(x))^2]$ noting that there's a different h and y for each different "extracted" training set. y is the value for x that could be present in a data set, and the expectation is over **all training set** that are drawn according to P . We will decompose this expectation into three components: bias, variance and noise.

Recall of statistics With Z random variable of possible values z_i with $i = 1, \dots, l$ and probability distribution $P(Z)$

Expected value or **mean** of Z is

$$\bar{Z} = E_P[Z] = \sum_{i=1}^l z_i \cdot P(z_i)$$

with the sum replaced by an integral and the distribution by a density function if Z is continuous.

Variance of Z is

$$Var[Z] = E[(Z - E[Z])^2] = E[Z^2] - E[Z]^2$$

with

$$E[Z^2] = E[Z]^2 + Var(Z)$$

1.6.1 Bias-Variance Decomposition

$$E_P[(y - h(x))^2] = E_P[h(x)^2 - 2yh(x) + y^2] = E_P[h(x)^2] + E_P[y^2] - 2E_P[y]E_P[h(x)]$$

Let $\bar{h}(x) = E_P[h(x)]$ denote the **mean prediction** of the hypothesis at x when h is trained with data drawn from P .

$$E_P[(y - h(x))^2] = E_P[(h(x) - E_P[h(x)])^2] + \text{variance} \\ \text{By doing some calculations, we obtain} \quad (E_P[h(x)] - f(x))^2 + \text{bias}^2 \\ E_P[(y - f(x))^2] \quad \text{noise}^2$$

So expected error is Variance + Bias² + Noise²

Bias: quantify the discrepancy between true function and $h(x)$, with $h(x)$ averaged over different TR data (**systematic error**)

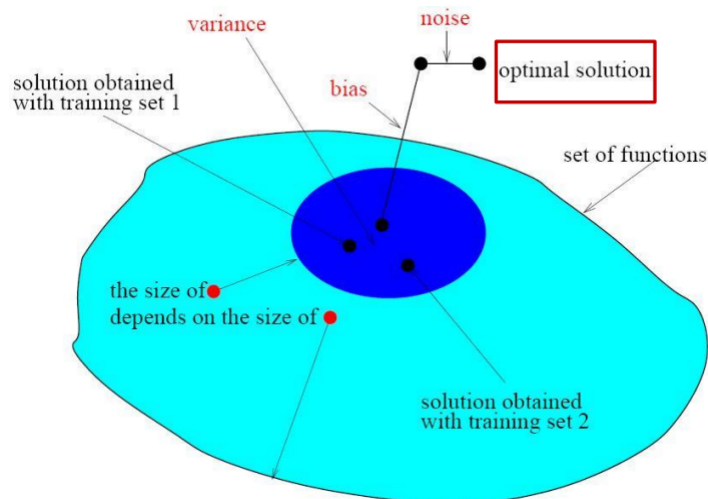
E.g. due to too small H , too rigid model.

Variance: quantify the variability of the response of model h for different realizations of the training data. Due to too high flexibility

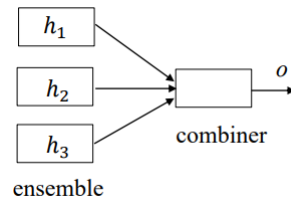
Noise: event the optimal solution can be wrong.

E.g. if for a give nx there are more than one possible d .

It's irreducible, doesn't depend on the model.



1.7 Ensemble Learning



Take advantage of multiple models. In regression: simple average committee, mean square error of a committee... In classification, take a vote over many classifier.

1.7.1 Bagging

Bootstrap Aggregating Combine many classifiers. Train n classifiers on different subsets of TR and differentiate each training using bootstrap (resampling with replacement). Thing to bias-variance: high variance model can perform well on average. In other terms: the average of h reduce the variance. In regression use the mean, in classification use the vote of the n classifiers.

1.7.2 Boosting

If the models have the same errors we have no advantage on esembling.

Differentiate each training concentrating on errors (more weight to difficult instances, e.g. more likely to be included in the TR of the next classifier, for example: train the 1st classifier, then train the 2nd with more weight on the instances misclassified before...)

Combine the results by output weights (weighted vote for classifiers, with more weight to low error classifiers).

If not stopped, boosting will learn to classify correctly all instances from TR, providing an incremental approach to construct complex models. Can suffer from noisy data (which are weighted more)

1.7.3 Feature Selection

Find a selection of features that are more informative for the problem at hands. Benefits for: dimensionality reduction and filtering of irrelevant information and noise, interpretability...

Computationally hard (many possible subsets of features, retraining...) typically an heuristic search of the best subset by greedy or other optimization techniques.

1.8 Applications

A huge number of successful applications in all the fields of ML.

1.8.1 Character recognition (classification)

First approaches Standard NN with 256 inputs (16×16 pixels): misclassification rate of 5/20% (mostly due to lack of invariance as rotations ecc.)

Basic idea Exploiting the architecture design to include prior information into NN:

Restricting the network architecture by extracting local features

Constraining weights by sharing them among different units, as to reduce the number of free parameters while still allowing more complex connectivity (same operation on different part of the input, because the features of the characters should appear in different areas of the image)

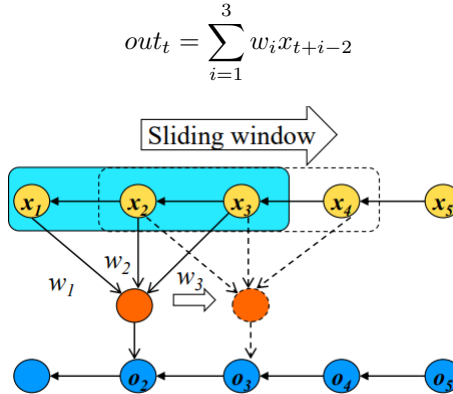
⇒ **Convolutional neural networks**

1.8.2 Convolutional Neural Networks

The name Takes inspiration from the convolutional operator

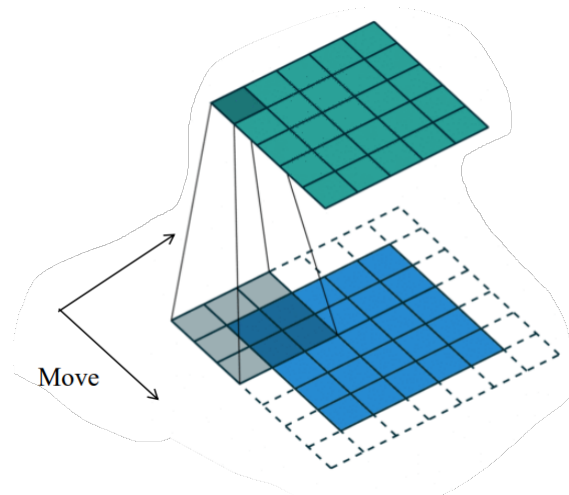
$$(f * g)(t) = \int_{-\infty}^{\infty} f(\tau)g(t - \tau) d\tau = \int_{-\infty}^{\infty} f(t - \tau)g(\tau) d\tau$$

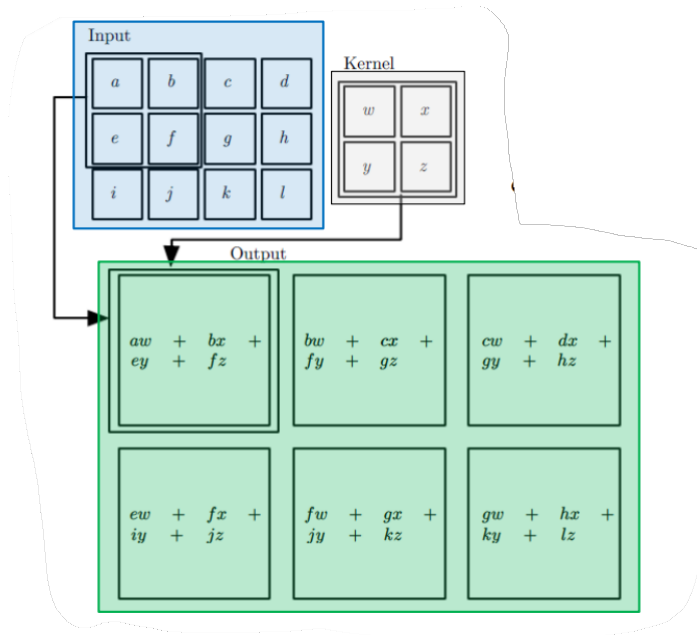
Concept: weighted average of a function f weighted by another function g moved over time (**sliding**) \Rightarrow a simple neural unit example over a stream



We will call this a "time-delay NN": weights are tuned as usual by learning and there's weight sharing.

2D convolution Using a $n \times n$ kernel (**local receptive field**), we traverse the image with 1px shifts (**stride**), with padding used to treat the border.





With a 2D image I and a 2D kernel K

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$$

or alternatively

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(i + m, j + n) K(m, n)$$

The units weights are a **filter**, trained to detect some features in the image:

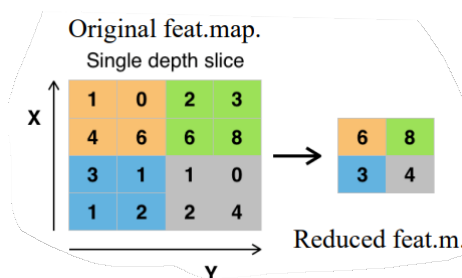
Local and small receptive field (kernel): by this architecture constraints, learnt "filters" produce the strongest response to a spatially local input pattern

The same unit/filter is **applied across the entire image**: this allows for features to be detected regardless of their position in the visual field, thus constituting the property of **translation invariance**

Learnable filters: we can learn the filters that in traditional algorithms for image processing were hand-engineered. This independence from prior knowledge and human effort in feature design is a major advantage.

Stacking many layers: feature maps unit represent larger and larger areas of the original image (assembling areas of the previous feature maps).

Pooling So we reduce the feature map by subsampling (see before), or by operating a mean (average or weighted average), or by a max pool operation over the neighbors of each input. In practice, instead of producing a value for each input pixel, we produce a value for a rectangular set of pixels taking the mean or the max of the kernel/neuron outputs.



Max pooling with a 2×2 filter and stride of 2

Pooling also further helps to make the representation become approximately invariant to small translations of the input.

Overview CNNs exploits weight sharing to setup a shifting windows or **local field** of the units over a segment of the input signal, also extended to 2D images and re-apply for many layers (**feature maps**).

Advantages

Local connections and weight sharing: detect local patterns and invariant to translations, helping reducing the number of free-parameters.

Pooling: helps reducing the dimension of the representation (in each layer, so a pyramid of layers), also helps to small shifts and distortions

Historical instance of deep neural network, progressive abstracting features from images.

How to use? Training is typically made by backpropagation, with the many heuristics that we studied and some specializations for weight sharing, pooling ecc. . .

Given the typical usage of huge networks and large amount of data, many hyperparameters are fixed by experience and expert suggestion, as it would be very expensive to run cross-validation over a large set of hyperparameters.

Modern CNNs There are many variants and specialized architectures and models fro processing of images, and efficient implementations through tensor representation of matrices, exploiting GPUs. Even using the features of models already trained with large benchmarks: many pre-trained CNNs for image classification, segmentation, face recognition and text detection available.

Parallelize linear operations on GPU Matrix multiplication (i.e. dot product) is a typical linear operation performed efficiently on GPU. It can be parallelized by parallelizing the sum $(AB)_{ij} = \sum_{k=1}^m A_{ik}B_{kj}$ or a set of independent products $A_{ik}B_{kj}$ for each i, j

In computer programming, a **tensor** is a multi-dimensional matrix with a linear operation called **tensordot**. In many libraries such as TensorFlow we can use tensors and perform linear operations as **tensordot0s** between them.

1.8.3 Deep Learning

Framework *"Deep learning allows computational models that are composed of multiple processing layers to learn representations of data with multiple levels of abstraction. These methods have dramatically improved the state-of-the-art in speech recognition, visual object recognition, object detection and many other domains".*

Starting from 2010, deep learning architectures are no longer unfeasible, thanks to the availability of large datasets, powerful computing systems. . .

Deep learning is a **general framework** which includes different models:

Deep Neural Networks

Convolutional Neural Networks

Deep Belief Networks and Generative Approaches

Recurrent and Recursive Neural Networks

...

This contrasts with "shallow models", e.g. neural networks with few layers.

Implement Many approaches, basically by building MLPs with many layers and, for example, pre-training techniques. Common aspects among approaches:

Multiple layers of nonlinear processing units

Supervised or unsupervised **learning of feature representations in each layer**, with the layers forming a **hierarchy from low-level to high-level features**, giving different levels of abstraction

Hence, hierarchical **sparse/distributed representation** of the input data

Hierarchical Abstraction Unsupervised or semi-supervised feature learning and hierarchical feature extraction **instead of features engineering**: moving towards the concept of "learning representations of data".

The abstract features can be "reused" at the higher levels: **combine to generalize** to unseen during training.

Techniques In general, deeper networks are able to use less units in each layer, so less parameters and less training data to achieve a good generalization. But, many layers can be harder to train, hence emphasis on methods to improve

- gradient descent (also gradient vanishes/explodes through layers)
- regularization (due to large networks)
- better exploitation of data

Do we need many layers? The general idea is that when a function can be compactly represented by a deep architecture, it might need a very large architecture to be represented by an insufficiently deep one.

An example from logic circuits: a two-layer circuit of logic gates can be represented by any boolean function, and any boolean function can be written as a sum of products (**disjunctive normal form**, \wedge gates on the first layers with optional negation of inputs and \vee gates on the second layer). With depth-two logical circuits, most boolean functions require an exponential (in input size) number of logic gates to be represented.

Examples: the parity function (return 1 \Leftrightarrow there are an odd number of 1 over N binary inputs) with 2 inputs can be done with 3 gates, with 3 inputs 5 gates, with N inputs $2^{N-1} + 1$ gates, exponential. If we use $\log N$ layers, we have fewer gates. However this doesn't hold for all classes of functions.

The **universal approximation theorem** is a fundamental contribution, it shows that 1 hidden layer is sufficient in general but doesn't assure that a "small number" of units could be sufficient and it also doesn't provide a limit on such number. For many function families it's possible to find boundaries on such number and "**no flattening**" results (on efficiency): cases for which the implementation by a single hidden layer would require an exponential number of units (in terms of n input dimension) or exponential non-zero weights, while more layers can help for the number of units/weights and also for learning.

Examples Cases with an exponential number of units are for example classes of problems requiring one hidden unit corresponding to each input configuration that needs to be distinguished. For example, the number of possible binary functions on vector $v \in \{0, 1\}^n$ is 2^{2^n} and selecting one such function requires 2^n bits which in general requires $O(2^n)$ degrees of freedom. It's a practical issue for the dimension of the network, but also implies that it's difficult to learn complex tasks with few examples.

On the other side there are families of functions that can be approximated efficiently by an architecture with depth greater than some value d , but which would require a much larger model if the depth is restricted to be $\leq d$. There's in general **no guarantee** that your task shares this property.

Theoretical Analysis Deep models can exploit the **compositionality** of internal representation: an exponential gain in representation power, due to the fact that simpler concepts represented in a layer of the network can be exploited as primitives by the next layer to represent more complex concepts, avoiding explicit combinatorial representation and learning of features.

More in general, but it's an open theoretical research topic:

No flattening results: no-flattening theorems, compositional functions that can be well implemented by a deep neural network cannot be implemented retaining the same efficiency while flattening the neural network

A complete list of no-flattening theorems would show exactly when deep networks are more efficient than shallow networks

So, deep networks can be seen as compact models with respect to potentially larger shallow models for the task (even exponentially more compact). Hence more efficient, not just from the computational point of view but from the learning point of view too: less units and less weights help learning on complex tasks to achieve good generalization with less examples.

Inductive Bias Choosing a deep model encodes a very general belief that the function we want to learn should involve composition of several simpler functions. If our task matches this bias, then of course the deep shape of the learner is suitable, and it happens that generalization is better due to the use of many layers. Typical examples are: the structure of images (composition of sub-graphical parts), the structure of language (text and speech), music... new fields are under discovering.

This doesn't apply to all data and tasks, of course. Also note that it's still a quite general bias, much more than other ad hoc constraints.

Curse of Dimensionality Many learners rely on local approximation (K-nn, local kernels, decision trees...), but often the smoothness assumption (or local consistency) is not enough. They need training examples to generalize the surrounding and they may need many examples: $O(k)$ examples to distinguish $O(k)$ regions, or exponential number of regions (in k) with additional assumptions.

To approach this, one can make strong, task-specific assumptions, losing the generality. Instead, the deep learning framework chooses a quite general inductive bias related to composition of functions: we assume that the data was generated by the composition of factors or features, potentially at multiple levels in hierarchy. This allows to:

- achieve a potential exponential gain between number of examples and number of regions that can be distinguished
- generalize non-locally
- learn with less examples

Practical issues How many layers? How many units? Model selection!

In general, for deep learning: deeper network often with less units, so less parameters, so less training data. But many layers harder to optimize during learning (exploding/vanishing gradients...)

Representation Learning "*Representation learning is a set of methods that allows a machine to be fed with raw data and to automatically discover the representations needed for detection or classification.*"

Deep learning methods are representation-learning methods with multiple levels of representation. But this concept is more abstract and can be applied to many models, and neural networks in general.

This is in most cases referred to raw data such as images or texts.

Basic ideas Many information processing tasks can be very easy or very difficult, depending on the representation of the information (this applies to every field of computer science).

In machine learning, a good representation is one that makes a subsequent learning task easier. The **manual design of features** is difficult, even decades for specific communities (languages, images,...). Supervised learning in a deep neural network is an example that leads to an **automatic representation** at every hidden layer, taking on properties that make the output layer task easier.

Obtaining or exploiting hidden representation The historical case in obtaining the hidden representation is the semi-supervised learning: we can learn a representation for the unlabeled data and then use it to solve supervised tasks (**pretraining approaches**).

To exploit the hidden representation, we can use the learned representation for other tasks (**transfer learning approaches**)

Implementing Deep Learning

Many approaches. Let's start with a simple one: building a multi-layer perceptron with many layers and pretraining techniques.

Pretraining The first approach to make possible to train a deep supervised network was the **Greedy Layer-Wise Unsupervised Pretraining**: unsupervised learning as pretraining to capture the shape of the input distribution. It makes training the whole network easier, used by autoencoders.

Each layer is optimized independently (*greedy layer-wise*) in an *unsupervised* way, constituting a *pretraining* for the final fine-tuning of the network. Works in terms of: good initialization strategy, regularization (in terms of discovering features that simplify the unsupervised process, which can be a more appropriate regularization technique when the underlying functions are "complicated" and shaped by regularities of the input distribution) and reducing the variance of the estimation.

Autoencoders A neural network that is trained to attempt copy its input to its output. Internally, it has a hidden layer h that describes a code used to represent the input. The network may be viewed as consisting of two parts: an encoder function $h = f(x)$ and a decoder that produces a reconstruction $r = g(h)$.

There are many forms of autoencoders, for example:

Undercomplete: hidden layer smaller than the input.

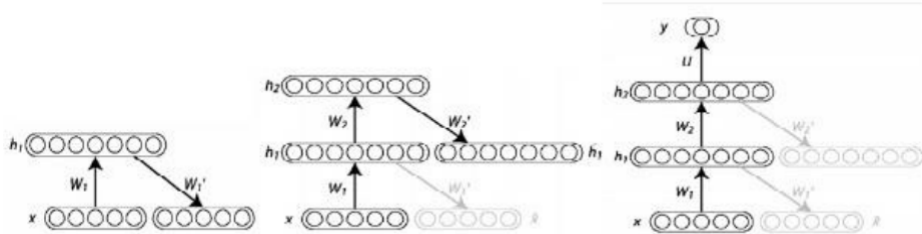
Forced to capture the most salient features of the training data by architectural constraints.

Overcomplete: hidden layer bigger than the input.

Regularization make the constraint for sparsity of the representation, robustness to noise and other properties of interest beside the trivial input-output copy capability.

Unsupervised!

Layer-wise pretraining Unsupervised, it can exploit also unlabeled data for the layer by layer training.



1. Train the first layer as an autoassociator to minimize the reconstruction error of the raw input (unsupervised because we only need unlabeled examples)
2. The hidden units outputs in the autoassociator are now used as input for another layer, also trained to be an autoassociator (again, unsupervised)
3. Iterate 2. to add the desired number of layers
4. Take the last hidden layer output as input to a supervised layer and initialize its parameters (either randomly or by supervised training, keep the rest of the network fixed)
5. Fine-tune all the parameters of this deep architecture with respect to the supervised criterion

The hope is that the pretraining has put the parameters of all the layers in a region from which a good local optimum can be reached by local descent. Despite many initially successful cases, the general role of pretraining is **nowadays under critical revision by researchers**.

Needed? Pretraining can yield improvements for some tasks, but not in other. It allowed to start deep learning, but difficult to be managed.

Transfer Learning Using the representation discovered in a model to improve another model. We assume there exist features that are useful for the different settings or tasks, corresponding to underlying factors that appear in more than one setting. But also we can use a trained model for another task, with same inputs but different targets (**multi-task learning**), or changing the input domain but sharing features (**domain adaption**), or take advantage of pre-training models from a larger dataset... (note that the bold terms are subfields of machine learning)

Example of pretrained networks

AlexNet, a CNN that has been trained on $\simeq 1.2$ million images, classifying them on 1000 categories: as a result, this model has learned rich feature representations for a wide range of images.

It has 5 convolutional layers and 3 fully connected layers. Its learned features can be transferred to new classification tasks, e.g. by replacing the last 3 layers for your task and train those new layers.

There are many pretrained CNNs for image classification, segmentation, face recognition and text detection.

AlexNet in matlab, also MatConvNet, VGG, ResNet, GoofLeNet... and similarly for other DL libraries such as Keras...

Distributed Representation "In a distributed representation, their elements (the features) are not mutually exclusive and their many configurations correspond to the variations seen in the observed data."

Deep learning methods exploit distributed representation with multiple levels of representation, but the concept is more abstract and can be applied to many models.

An example of **symbolic representation** is the **one-hot representation**: one element is 1, the others are 0 and the distance between two elements is always $\sqrt{2}$. With a **distributed representation** we have a dense vector of real numbers, and each symbol or concept is represented as the set of unit activations: this allows to share similarities and some learned features among concepts. The distance between two concepts reflects the meaning, and their similarity.

Input or internal representation? Symbolic and distributed refers to the input or the internal representation? We are talking in general, but **learning acts on the internal representation**. In general, distributed representation is used:

In input if you have a background knowledge to build it (if it's done improperly then it could hamper the task solving)

Automatically, by learning, if distributed representation can be used internally in the model (e.g. hidden neural network layers)

Typically, one-hot is used for the input and the model is free to develop internally the distribution representation needed for the task at hand

Count the difference Distributed representation can represent n features with k values to describe k^n different concepts.

Sharing attributes By sharing a learned feature we can share its concept, disentangling it from the task.

Example 4 concepts: blue car, blue bike, red car, red bike. There's no need for 4 neurons, just 2 are enough: a neuron for blue/red and a neuron for car/bike. The neuron describing redness is able to learn about the concept of redness from images of cars and bikes, not just one.

If we share that neuron, we share the concept of redness.

Beyond Neural Networks The debate on distributed representation extend to the debate between **logic-inspired paradigms** and **neural network-inspired paradigms** for cognition.

Interpretability Less easy for distributed representation.

Deep Distributed Representation Deep learning exploits distributed representation through many layers, obtained by composing different levels of abstraction or by a hierarchy of reused features. This compositionality, as already discussed, can lead to further exponential boos to efficiency. Globally, deep learning learn a distributed representation of the data by disentangling shared causal factors that generates the data through different levels of abstraction.

Deep Learning Techniques Key techniques for training a layered neural network:

Originally: pretraining approaches

Nowadays:

SGD with momentum, with decay on η or minibatch, or Adam...

ReLU activation functions on hidden units

Cross-Entropy (max log-likelihood) loss with softmax for output layers, also to avoid saturation and small gradient effects

Regularization: early stopping and weight decay, also drop-out and batch normalization (for CNN and sigmoids)

Gradient Issues The magnitude of the gradients, as its backpropagated through the many layers, suffers from two main issues:

If the weights are small, the gradient shrinks exponentially (**vanishing gradient**)

If the weights are big, the gradient grows exponentially (**exploding gradient**)

Some approaches include the **gradient clipping**: if $\|g\| > v$ then $g = \frac{vg}{\|g\|}$ with g being the gradient and v the norm threshold, as to move in the gradient direction but bounding the weight update.

Most of the heuristics try to deal with the problem of vanishing gradient tat hampers the training in early deep neural network in the low layers. Traditional activation functions, such as the hyperbolic tangent, have gradients in the range of $(-1, 1)$ and the backpropagation computes gradients by the chain rule repeated through the layers. This has the effect of multiplying n of these small numbers to compute gradients of the "front" layers (closer to the inputs), in a n -layer network: the gradient decreases exponentially with n , with the result of a very slow training in the front layers. Same with big weights (exploding gradients), but less frequent.

To deal with this, many approaches: Rprop, short-cut connections, randomized neural networks... let's see the new ones.

ReLU Rectified Linear Unit: $f(x) = \max(0, x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$

Sparse activation: in a randomly initialized network, only about 50% of hidden units are activated (non-zero output)

Efficient gradient propagation with respect to vanishing and exploding effects

Efficient computation: only comparison, addition and multiplication

Faster and effective training of deep neural architectures

But it's non differentiable in 0!

Used in deep learning since 2009 due to the simplified and better gradient propagation through many layers and the avoiding of saturation effects of sigmoidal functions. Non differentiable in 0, but often assumed to be 0 for the left derivative and 1 for the right derivative: the approximation is acceptable and safe because there's already an approximation for x input, so unlikely to be $x = 0$ effectively.

Beyond ReLU there are:

ELU (Exponential Linear Unit), $f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{if } x \leq 0 \\ x & \text{if } x > 0 \end{cases}$

Leaky ReLU $f(x) = \begin{cases} 0.01x & \text{if } x < 0 \\ x & \text{if } x \geq 0 \end{cases}$

Batch Normalization Method that normalizes each batch by calculating each individual batch statistic such as mean and variance for each layer (**reparametrization**): normalize each matrix (databatch \times activation of units) with mean and variance, by shifting values to zero-mean and unit variance, and include it in backpropagation.

Normalizing inputs is a standard approach, batch normalization helps by making the data flowing between intermediate layers to stay normalize. It has a regularization effect, achieving faster learning and higher accuracy for deep learning.

Dropout Method that randomly selects a subset of the network during its training. Can be explained in terms of:

Ensembling: implicit bagging, different sub-networks

Regularization

Bagging is the combination of many classifiers/models: train n classifiers/models on different subsets of the TR and differentiate each training using bootstrap (resampling with replacement). Dropout aims to approximate this process, with an exponentially large number of sub neural networks (while still having a single network at test time), while also aiming at maximizing the diversity of the ensemble preventing complex co-adaptions on training data.

A basic algorithm:

Each time an example is loaded into a minibatch, we randomly sample a different binary mask to apply to all of the input and hidden units in the network.

The mask is sampled independently for each unit, and the probability of a mask value of 1 (causing the unit to be included) is an hyperparameter, typically 0.8 for input units and 0.5 for hidden units.

The forward propagation, back propagation and the learning algorithm are done as usual (only on a subset of units at the time)

The removed nodes are then reinserted into the network with their original weights.

There's parameter sharing among sub networks: in each single step only a small fraction of the possible sub networks are trained, but the parameter sharing causes the remaining sub networks to arrive at good settings of parameters.

It also has regularization effect:

Avoids training all units on all training data, reducing units interactions

Variance reduction as for bagging

Insert structured noise

For example, recognize a face with the hidden unit for the nose, compare with standard regularization by adding deformed inputs with random noise)

Moreover, it regularizes each unit to be not merely a good feature but a feature that is good in many contexts (different sub-networks). It can be shown to be equivalent to L2 weight decays using a different λ for each input, can be used for any model that uses distributed representation and SGD training and can be extended to any kind of random modification.

L1 Regularization Using the L1 norm (sum of absolute values) instead of the L2 norm in the penalty term. Favorite features elimination (weight zero) with respect to L2, which often only shrinks the weights without setting any of them to zero. This produces simpler final models.

Adversarial Training **Generative Adversarial Network:** two neural networks contesting with each other, one network generates candidates and the other evaluates them.