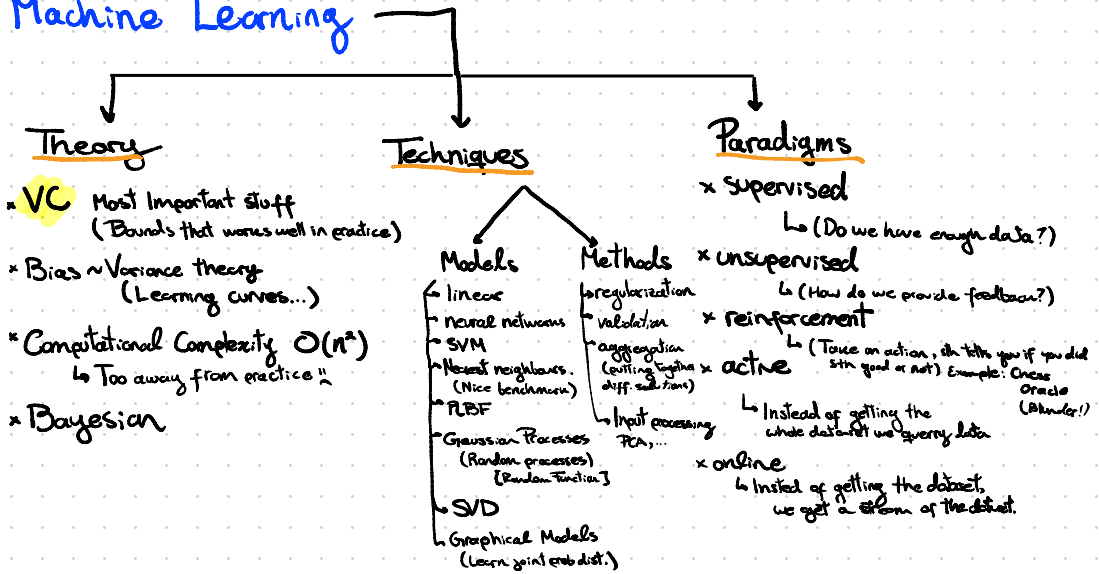


Caltech's
Machine
Learning
CS156

The biggest pitfall in theory is that the assumptions that are taken are divorced from the reality of how we use machine learning in practice

Machine Learning

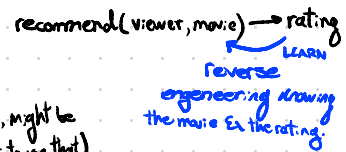


⚡ A *diminute* improvement in a machine learning problem can lead to massive problem.

example:
movie recommendation

⚡ A problem can be taught as a machine learning problem IFF

- We've a pattern *
- We don't know how to pin it down mathematically (If yes, might be better to use that)
- We've data (A MUST)



* How do we know? We don't know! But we can try to apply methods and we can determine if we're learning or not.

In fact we will use machine learning to see if there's a pattern

The Learning Problem.

Components of Learning

// A bank has no magical orb to know if a person is credit worthy or not.

- Input: x (customer data)
- Output: y (good/bad customer)
- Target function: $f: X \rightarrow Y$ (ideal credit approval formula)
- Data: $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$
- Hypothesis: $g: X \rightarrow Y$ (we hope approximate f well). (formula to be used)

UNKNOWN TARGET FUNCTION
 $f: X \rightarrow Y$

Training Examples
 $(x_1, y_1), \dots, (x_n, y_n)$

Final Hypothesis
 $g \approx f$

Learning Algo. A

Hypothesis Set
 H

(Set of candidate formulas)

It can be continuous, discrete... The level of application can be huge.

Ex: $\{x \rightarrow ax + b \mid (a, b) \in \mathbb{R}^2\}$

There is no really loss of generality. H can be the set of all possible functions.

Solution components

Two solution components of the learning problem: 1. The Hypothesis set $H = \{h\}$, $g \in H$
2. The learning algorithm.

Together, they're THE LEARNING MODEL.

Example: The Perceptron

Approve credit if $\sum_{i=1}^d w_i x_i > \text{threshold}$

he H can be written as:

$$h(x) = \text{sign} \left(\left(\sum_{i=1}^d w_i x_i \right) - \text{threshold} \right)$$

→ Add artificial constant $x_0 = 1$

rename: w_0

$$h(x) = \text{sign}(w^T x)$$

Things to learn.



Perceptron Learning Alg.

→ Pick a misclassified point:

$$\text{sign}(w^T x_n) \neq y_n$$

→ Update weight vector:

$$w \leftarrow w + y_n x_n$$

→ Run until no more are misclassified.

(It will end if data is linearly separable.)

How do we know if sth is linearly separable?

↳ Assume is not. We can transform data to be linearly separable. (At what cost?)

The rate of convergence of the perceptron learning algorithm is in fact terrible. There are pathological patterns where it doesn't converge.

Statistics → We make assumptions, we have a math model.

Machine learning → Least assumptions (general as possible). Don't have a math model as in statistics.

Optimization: Tool at disposal for machine learning.

Basic Premise of Learning

"use a set of observations to uncover an underlying process"

- How much is enough data?
↳ In practice is not sth you control. We'll see in theory more about that.
- The battle next of machine learning is its capacity of generalizing.

Is Learning Feasible?

Probabilistic approach:
We pick red ballots from a blue red ballot.

$$P(\text{pick red ballot}) = \mu$$

If we do N indep exp and get a crop of V red ballots, $M \approx V$? No and yes!

They can be totally diff. \uparrow In a big N , V can be close to μ . (Most probably)

In fact,

$$P[|V - M| > \epsilon] \leq 2e^{-2\epsilon^2 N} \quad \forall N, \epsilon > 0.$$

bad event.

Hoeffding's inequality.



A law of large numbers.

So $M \approx V$ is PAC.

(Probably approximately correct.)

- The probability is bounded regardless of μ .
- $V \approx M \rightarrow M \approx V$

- Supervised Learning \rightarrow The output data is explicitly given.
we've (input, correct output) Ex: $\{(x_1, y_1), \dots, (x_n, y_n)\}$
- Unsupervised Learning \rightarrow Classify data without knowing the name of the classes.
we've (input, ?) Example: Clustering.
// A way of getting a high level representation (high patterns) of the input
- Reinforcement Learning
we've (input, partial output, grade for this output)
 \uparrow
Reward for the output.
Ex: game simulators.

Learning: Unknown function. $f: X \rightarrow Y$
Each ballot is a point $x \in X$,

red: right hypothesis $h(x) = f(x)$
blue: wrong hypothesis $h(x) \neq f(x)$
fixed $h!$



Prob dist P on X .
↳ data is generated from P .

connection to learning



Note: This example is verification, not learning.

Solution: Having multiple bins



V is "in sample": E_{in}
 M is out-of sample: E_{out}
 $E_{in}(h)$: Error of approximating my target function using h .
 $E_{out}(h)$

Using Hoeffding inequality

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

\uparrow In sample performance. \uparrow out of sample performance



BUT, Hoeffding's doesn't apply to multiple bins!

Fortunately we can take the worst case. And with that, we get just a M factor which becomes meaningless.

Linear Model.

- × Input rep.
- × Linear classif.
- × Linear Reg.
- × Non linear transf.

→ MNIST Dataset

Input rep: $x = (x_0, x_1, \dots, x_{256})$ $\begin{bmatrix} 1 \\ x \end{bmatrix}^{1 \times 256}$

→ A linear model has 256 dim!

Sol: Features (Extract useful info)

Ex: Intensity and symmetry: $x = (x_0, x_1, x_2)$
(1s are less intense than 8)



Linear reg → real-valued output.

Example: calculate credit line based on client's features.

Regression output → $h(x) = w^T x$

• We use MSE: $(w_0 - f(x))^2 \rightarrow E_{in}(w) = \frac{1}{N} \sum_{i=1}^N (x_i w - y_i)^2$

$$X = \begin{bmatrix} -x_1^T & - \\ \vdots & \vdots \\ -x_n^T & - \end{bmatrix}$$

features

To minimize the error $E_{in}(w) = \frac{1}{N} \sum_{i=1}^N (x_i w - y_i)^2$

$$\nabla E_{in}(w) = \frac{2}{N} X^T (Xw - y) = 0 \Leftrightarrow X^T X w = X^T y$$

$\frac{X^T X}{N}$ is $\frac{X^T X}{N}$ / pseudoinverse.

The linear regression **is** just computing $w = X^+ y$

Linear reg for classification:

Binary valued functions are also real values. We can do
 However, the new problem is that it tries to fit all points into a line.
 So we're just minimizing the distance. So it just works as a good starting value for the weights

Sigmoid($w^T x$)

Non-Linear Transformations

$$[c] = \begin{cases} 0, & \text{if } c \leq 0 \\ 1, & \text{else} \end{cases}$$

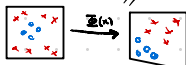
The model weights parameters is fixed $\neq 0$

Same problem!

We just apply non-linear transf. to the data.

$$(x_1, x_2) \xrightarrow{\Phi} (x_1', x_2')$$

• Data that wasn't linearly separable now it can be separated!



$$\vec{g} = \text{sign}(w^T z)$$



$$g(x) = \vec{g}(\Phi(x)) = \text{sign}(w^T \Phi(x))$$

$$x = (x_0, x_1, \dots, x_n) \xrightarrow{\Phi} z = (z_0, z_1, \dots, z_n)$$

// Powerful, but generalization can be terribly poor

$$x_1 \rightarrow x_n \xrightarrow{\Phi} z_1 \rightarrow z_n$$

$$y_1 \rightarrow y_n \xrightarrow{\Phi} y_1 \rightarrow y_n$$

No the weights decay in Z

$$W = (w_0, w_1, \dots, w_n)$$

Error Measures

$E(h, f)$ // error function

Almost pointwise def: $e(h(x), f(x)) \begin{cases} \text{examples} \\ = (h(x) - f(x))^2 & \text{Squared error} \\ = \mathbb{I}[h(x) \neq f(x)] & \text{Binary error} \end{cases}$

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

Outsample error: $E_{out}(h) = E_x[e(h(x), f(x))]$

How to choose error measure → depends on the problem. (Application Domain question)

If we don't know what to do, use **PLAUSIBLE** measures or **FRIENDLY** MEASURES

Squared error (Gaussian noise) → convex optimization

Noisy Targets

// Target function is not a function: $f(x_1) = 0, f(x_2) = 2$, where $x_1 = x_2$

Sol: use "target distribution" $P(y|x)$

(x, y) is now def. by a joint distribution.

Noisy target: $f(x) + \epsilon \Rightarrow f(x) = E(y|x)$

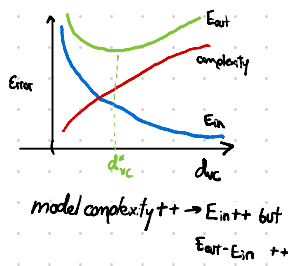
noise

$$P(x) \neq P(y|x)$$

← what are we trying to learn
 quantifies relative importance of X

Merging $P(x)P(y|x)$ as $P(x, y)$ mixes two concepts

We need a metric to measure the sophistication of the model \rightarrow dvc



E_{in} is used as a proxy to E_{out} .

$E_{out}(g) \approx 0$ is achieved through:

$$\uparrow E_{out}(g) \approx E_{in}(g)$$

$$\uparrow E_{in}(g) \approx 0$$

Learning is reduced to two questions

- Can we assure that $E_{out}(g)$ refers to $E_{in}(g)$?
- Can we make $E_{in}(g)$ small enough?

Obs: Sometimes it's impossible to have $E_{out}(g) \approx 0$. Ex: stock market. Having an error of 45% would mean you're rich already.

Training vs Testing.

Testing: $IP[|E_{in} - E_{out}| > \epsilon] \leq 2 \exp(-2\epsilon^2 N)$

Training: $IP[|E_{in} - E_{out}| > \epsilon] \leq 2M \exp(-2\epsilon^2 N)$

\hookrightarrow Objective: Find a better bound than M

M comes from

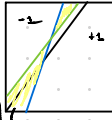
Bad events: $P[B_1 \text{ or } \dots \text{ or } B_M] \leq P[B_1] \dots P[B_M]$

no overlaps: M

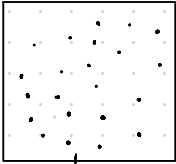
In practice, Bad Events overlaps!

$\hookrightarrow \Delta E_{out}$: change in ± 1 across
 ΔE_{in} : change of labels in data part

$$|E_{in}(h_1) - E_{out}(h_1)| \approx |E_{in}(h_2) - E_{out}(h_2)|$$



If there is data here, the E_{in} changes. If it's over almost the same, there is almost no change in ΔE_{in} !



From this constellation of points, how many patterns of red & blue can I get?

\hookrightarrow Number of dichotomies

$h: X \rightarrow \{\pm 1\}$ // hypothesis

$h: \{x_1, \dots, x_N\} \rightarrow \{\pm 1\}$ // dichotomy

The number of hypotheses |H| can be infinite, but the number of dichotomies is at most 2^N

The growth function

≥ 1 give you a budget N, choose where to place N points so the dichotomies are maximized.

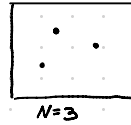
The growth function counts the most dichotomies on any N points

$$m_H(N) = \max_{x_1, \dots, x_N \in X} |H(x_1, \dots, x_N)|$$

set of dichotomies.

It satisfies: $m_H(N) \leq 2^N$

Apply $m_H(N)$ to perceptrons



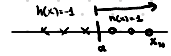
We have 2^3 ways to separate it with a line.
 $m_H(3) = 2^3$

And for $m_H(4)$?



We can't make this dichotomy with a line!
 In fact, $m_H(4) = 14$
 \hookrightarrow 14 perceptrons

Illustration of the growth function:



Example (Positive Rays) $f: \mathbb{R} \rightarrow \{\pm 1\}$
 $h(x) = \text{sgn}(x-a)$ $m_H(N) = N+1$

(Positive Intervals) $m_H(N) = \binom{N+1}{2} + 1$

// ways of assigning N+1 points to two separate sets.

(Convex sets) $h: \mathbb{R}^2 \rightarrow \{\pm 1\}$
 $h(x) = +1$ is convex.
 $h(x) = 0$ not convex.



Sol: put the N points in a circle, we get 2^N dichotomies (max bound!)

$$m_H(N) = 2^N$$

\hookrightarrow When this happens, we say f shatters the points

$$\mathbb{P}[|E_{in} - E_{out}| > \epsilon] \leq 2M e^{-2\epsilon^2 N}$$

Breakpoint

→ The point where you fail to get all possible subsets (we can't shatter the data anymore).

If no data set of size k can be shattered by H , then k is a break point of H .

Breakpoint examples

- positive rays: $k=2$
- positive intervals: $k=3$
- convex sets: $k=+\infty$

$$\text{lowest } k / m_{H^k}(k) < 2^k$$

- Let's replace M with $m_{H^k}(N)$.
 - If $m_{H^k}(N)$ has polynomial order. I've won.
- Just proving that $m_{H^k}(N)$ is polynomial is enough to prove that learning is possible.

Le main point

$$\nexists \text{ break point} \rightarrow m_{H^k}(N) = 2^N$$

$$\exists \text{ break point} \rightarrow m_{H^k}(N) \in P(N)$$

\hookrightarrow polynomials in N

How to replace M in Hoeffding's inequality?

$$\mathbb{P}[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 4m_H(2N)e^{-\frac{1}{8}\epsilon^2 N}$$

• This bound can be improved. But this is enough to prove that learning is possible.

C'est Vapnik-Chervonensis inequality

(We want to bound $m_{H^k}(N)$)

tactic: $m_{H^k}(N) \leq \dots \leq$ a polynomial.

Key quantity: $B(N, k)$: Max num of distances on N points, with break point k .

\hookrightarrow BIN

Derived from this.

$$m_{H^k}(N) \leq \sum_{i=0}^{k-1} \binom{N}{i}$$

then power is $N^k \rightarrow$ constant! (A polynomial)

Upper bound for anything with break point k .

$\rightarrow m_{H^k}(N)$ is a polynomial

The VC dimension

Def: The VC dimension of a hypothesis set H denoted by $d_{VC}(H)$ is the largest value of N for which $m_H(N) = 2^N$

"the most points H can shatter"

$N \leq d_{VC}(H) \rightarrow H$ can shatter N points.

$k > d_{VC}(H) \rightarrow k$ is a breakpoint of H

$$\rightarrow m_H(N) \leq \sum_{i=0}^{d_{VC}(H)} \binom{N}{i}$$

poly of order N

Ex: Positive rays: $d_{VC} = 1$

2D perceptrons: $d_{VC} = 3$

Convex sets: $d_{VC} = \infty$

VC dimension and Learning

$d_{VC}(H) < \infty \rightarrow g \in H$ will generalize.

> This is independently from the Learning alg, Input distribution & target function.

> You will generalize with high probability with the input data.

Example: For d -dim perceptrons $d_{VC} = d + 1$

→ Proof: Let $N = d + 1$, $X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{d+1} \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 1 \end{bmatrix}$. For arg $y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{d+1} \end{bmatrix} = \begin{bmatrix} \pm 1 \\ \pm 1 \\ \vdots \\ \pm 1 \end{bmatrix}$

$d + 1$ is also the number of parameters in the perceptron.

Let's find $w / \text{sign}(Xw) = y$. We can do $Xw = y$. Because X inv, $w = X^{-1}y$!
 Thus, we can shatter these $d + 1$ points → $d_{VC} \geq d + 1$
 For $d \geq 2$, more points than dim → $x_j = \sum a_i x_i$ (lin. dep), not all the $a_i = 0$.
 The non-zero a_i get $y_i = \text{sign}(a_i)$ and x_j gets $y_j = -1$. No perceptron can implement this dichotomy!
 $x_j = \sum a_i x_i \rightarrow w^T x_j = \sum a_i w^T x_i$. If $y_i = \text{sign}(w^T x_i) = \text{sign}(a_i)$, thus $a_i w^T x_i > 0$
 So $w^T x_j = \sum a_i w^T x_i > 0 \rightarrow y_j = \text{sign}(w^T x_j) = +1$

Interpreting the VC dimension.

of params in model \sim degrees of freedom.
 $d_{VC} \sim$ binary degrees of freedom (if you can separate data or not)

Obs: Parameters may not contribute degrees of freedom.

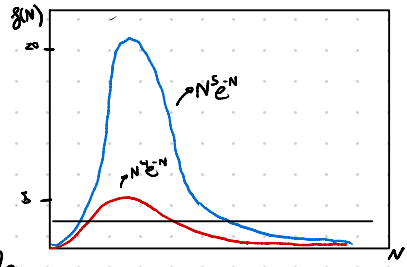
$x \rightarrow \textcircled{5} \xrightarrow{\text{perceptions}} \textcircled{5} \rightarrow \dots \rightarrow \textcircled{5} \rightarrow y$
 Here there are ∞ params, but 1 degree of freedom.
 → d_{VC} measures the EFFECTIVE number of parameters.

2- How many data points we need

VCineq: $P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq \underbrace{4m_H(2N)}_S e^{-\frac{1}{2} \epsilon^2 N}$

If we want certain ϵ and S , how does N depend on d_{VC} ?

simplification of S
 $f(N) = N^d e^{-N}$



$P[|E_{out} - E_{in}| > \epsilon] \leq 4m_H(2N) e^{-\frac{1}{2} \epsilon^2 N}$

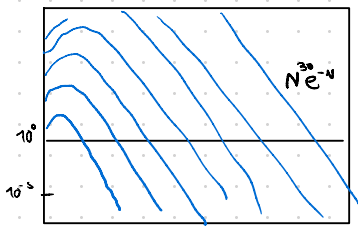
Get $\epsilon(S)$: $\epsilon = \sqrt{\frac{S}{N} \ln(4m_H(2N))} = \Omega$

What prob $1 - \delta$, $|E_{out} - E_{in}| \leq \Omega(N, H, \delta)$

↔ $E_{out} - E_{in} \leq \Omega(N, H, \delta) \rightarrow E_{out} \leq E_{in} + \Omega$ (with prob $1 - \delta$)

goes down with bigger H
 goes up with bigger H

Rescale:



> Bigger VC dimension → Need for most samples
 $d_{VC} \propto$ number of samples needed to obtain a certain performance.

Rule of thumb: $N \geq 10 d_{VC} \rightarrow N \propto d_{VC}$

↳ Number of samples needed is proportional to d_{VC}

Bias Variance Tradeoff

Approximation - Generalization tradeoff.

Small E_{out} : good approx of f out of sample

More complex $H \rightarrow$ Better chance of approximating f
 Less complex $H \rightarrow$ Better chance of generalizing out of sample

$f \in H$ never happens...

The biasvariance is a **new approach**

VC analysis: $E_{out} \leq E_{in} + \Omega$

Biasvariance: Decompose E_{out} into $\left\{ \begin{array}{l} \text{How well it approximates } f \\ \text{How well we can zoom in on a good } h \in H \end{array} \right.$
 Applies to real-valued targets and uses square error

$$E_{out}(g^{(0)}) = \mathbb{E}_x [(g^{(0)}(x) - f(x))^2]$$

I want to get rid of (0) , my dataset.

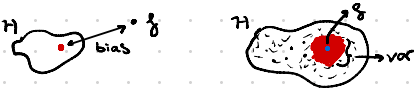
$$\mathbb{E}_D [E_{out}(g^{(0)})] = \mathbb{E}_D [\mathbb{E}_x [(g^{(0)}(x) - f(x))^2]] \rightarrow \text{just integration! (w/ } d(x)) \\ = \mathbb{E}_x [\mathbb{E}_D [(g^{(0)}(x) - f(x))^2]]$$

Lets define $\bar{g}(x) := \mathbb{E}_D [g^{(0)}(x)] \rightarrow \mathbb{E}_D [(g^{(0)}(x) - \bar{g}(x) + \bar{g}(x) - f(x))^2] = \mathbb{E}_D [(g^{(0)}(x) - \bar{g}(x))^2] + (\bar{g}(x) - f(x))^2 + 2(g^{(0)}(x) - \bar{g}(x))(\bar{g}(x) - f(x))$

$$= \underbrace{\mathbb{E}_D [(g^{(0)}(x) - \bar{g}(x))^2]}_{\text{var}(x)} + \underbrace{(\bar{g}(x) - f(x))^2}_{\text{bias}(x)}$$

Variance of my model
This is my best candidate from my hypothesis set. We differ from f because

$$\rightarrow \mathbb{E}_D [E_{out}(g^{(0)})] = \mathbb{E}_x [\text{bias}(x) + \text{var}(x)]$$



Example: $f: [1, 1] \rightarrow \mathbb{R}, f(x) = \sin(\pi x)$

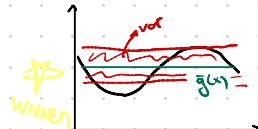
Our training dataset: $N=2$ (Lol!)

Hyp. set: $H_0: h(x) = b$

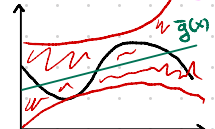
$H_1: h(x) = ax + b$

who is better?

FOR WHAT?



bias = $1/2$
var = $1/4$ + = $3/4$

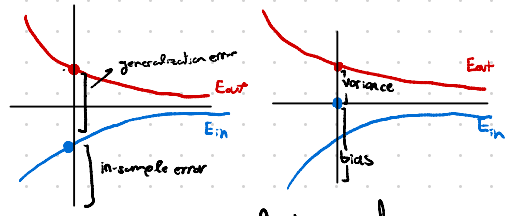
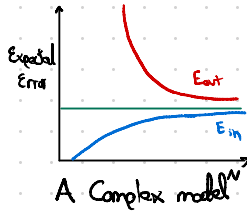
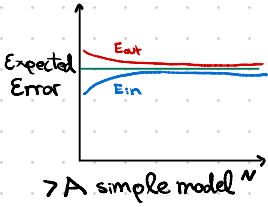


bias = 0.21
var = 1.69

Learning Curves

VC Analysis / Bias Variance

Dataset D of size N . How does E_{in} and E_{out} vary with N ?



careful when you do the transformations!

Looking at the data BEFORE choosing the model can be hazardous to your E_{out}


Do just "snoop" at the data.
 ↳ The dataset with the dvc guarantees is the one you had before.
DATA SNOOPING

Non-linear transformations

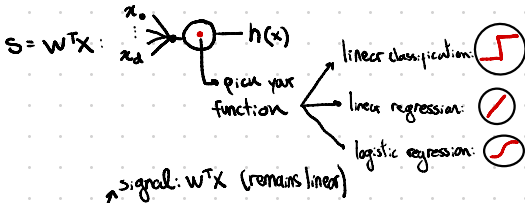
$$Z = \Phi(x)$$

ex: $Z = (1, x_1, x_2, x_1^2, x_2^2) = (z_0, \dots, z_d)$

Obs: Finite hyp on \mathcal{X} .

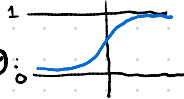
with this we can do: 
 $dvc \leq d+1$
 ↳ we're limited to Φ .

Logistic Regression



Logistic function θ :

$$\theta(s) = \frac{e^s}{1 + e^s}$$



↳ soft threshold.
 ↳ models: uncertainty
 This is also called **sigmoid**

We pick this because is great for optimization

Example: (x, y) . y is noisy

$$P(y|x) = \begin{cases} \beta(x), & \text{for } y=+1 \\ 1-\beta(x), & \text{for } y=-1 \end{cases}$$

My target is $f: \mathbb{R}^d \rightarrow \{0, 1\}$ probability

$$\text{Learn } g(x) = \theta(w^T x) \approx f(x)$$

↳ How to train weights.

How we build an error measure in this context?
 ↳ Plausible error measure based on likelihood.

$h(x) = \theta(s)$ is interpreted as a probability.

$$\theta(-s) = 1 - \theta(s)$$

$$P(y|x) = \theta(y w^T x)$$

$$\prod_{n=1}^N P(y_n | x_n)$$

$$\prod_{n=1}^N \theta(y_n w^T x_n) \text{ (Likelihood)}$$

↳ we want to maximize this w.r.t w .

We max the log likelihood instead (Equivalent). And we minimize with -1 instead.

$$\min -\frac{1}{N} \ln \left(\prod_{n=1}^N \theta(y_n w^T x_n) \right) = \frac{1}{N} \sum \ln \left(\frac{1}{\theta(y_n w^T x_n)} \right)$$

$$\rightarrow E_{in}(w) = \frac{1}{N} \sum \ln \left(\frac{1}{\theta(y_n w^T x_n)} \right) = \frac{1}{N} \sum \ln \left(\frac{1}{e^{h(x_n, y_n)}} \right)$$

To minimize this, there is no close form solution

"cross-entropy" error

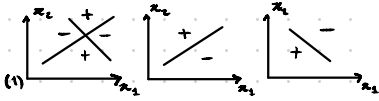
Iterative Method!
 ↳ Gradient descent.

Neural Networks

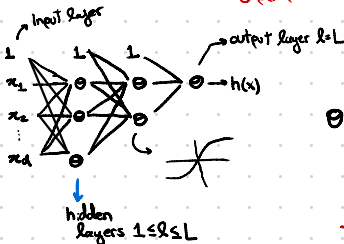
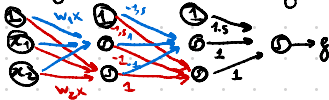
// features from features

biological function \rightsquigarrow biological structure

Key: combining perceptrons



A multilayered perceptron implementing this (1)



$$\theta(s) = \tanh(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$

$w_{ij}^{(k)}$ $\begin{cases} 2, 1 \leq k \leq L \text{ layers} \\ i, 0 \leq i \leq d^{(k-1)} \text{ inputs / weights like } x_0 \text{ constant!} \\ j, 1 \leq j \leq d^{(k)} \text{ outputs} \end{cases}$

$$x_j^{(k)} = \theta(s_j^{(k)}) = \theta\left(\sum_{i=0}^{d^{(k-1)}} w_{ij}^{(k)} x_i^{(k-1)}\right)$$

Apply x to $x_1^{(0)} \dots x_{d^{(0)}}^{(0)} \rightarrow x_1^{(1)} = h(x)$

Algo: Backpropagation.
 1. Init weights $w_{ij}^{(k)}$ RANDOMLY.
 2. for $t=0, 1, 2, \dots$ do
 Pick $n \in \{1, \dots, N\}$
 Forward: compute ALL $x_j^{(k)}$
 Backward: compute ALL $\delta_j^{(k)}$
 Update $w_{ij}^{(k)} \leftarrow w_{ij}^{(k)} - \eta \pi_{ij}^{(k-1)} \delta_j^{(k)}$
 Iterate until criteria stop
 return $w_{ij}^{(k)}$

SGD

$$E_{in}(w) = \frac{1}{N} \sum \ln(1 + e^{-y_n w^T x_n}) \quad // \text{logistic regression}$$

We pick a small subset (batchsize) of the dataset and do gradient descent.

// Theog: $E_n(\nabla g(x_n, y_n)) = E_{in}$!

Advantages of SGD $\begin{cases} 1. \text{ less order of comp.} \\ 2. \text{ randomization (prevents you escape of silly local minima)} \\ 3. \text{ simple} \end{cases}$

Rule of thumb: $\eta = 0.1$ is ok (learning rate)

SGD @ NN

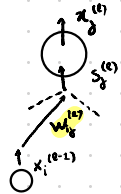
Error on sample: $e(h(x_n), y_n) = e(w)$

We need the gradient:

$$\nabla e(w) = \frac{\partial e(w)}{\partial w_{ij}^{(k)}}$$

We can evaluate $\frac{\partial e(w)}{\partial w_{ij}^{(k)}}$ analytically or numerically.

$$\rightarrow \frac{\partial e(w)}{\partial w_{ij}^{(k)}} = \frac{\partial e(w)}{\partial s_j^{(k)}} \frac{\partial s_j^{(k)}}{\partial w_{ij}^{(k)}}$$



$$x_i^{(k-1)} = \frac{\partial s_j^{(k)}}{\partial w_{ij}^{(k)}}$$

What we need:
 $\hookrightarrow \delta_j^{(k)} = \frac{\partial e(w)}{\partial s_j^{(k)}}$

δ for the final layer

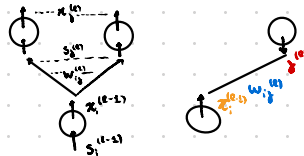
$$\delta_1^{(L)} = \frac{\partial e(w)}{\partial s_1^{(L)}}$$

$$e(w) = e(h(x_n), y_n)$$

$$\pi_1^{(L)} = \theta'(s_1^{(L)})$$

The value of the output!

$$\rightarrow \theta'(s) = 1 - \theta^2(s) \text{ for tanh.}$$



$$\delta_i^{(k-1)} = \frac{\partial e(w)}{\partial s_i^{(k-1)}} = \sum_{j=1}^{d^{(k)}} \frac{\partial e(w)}{\partial s_j^{(k)}} \times \frac{\partial s_j^{(k)}}{\partial x_i^{(k-1)}} \times \frac{\partial x_i^{(k-1)}}{\partial s_i^{(k-1)}}$$

$$\delta_i^{(k-1)} = \sum_{j=1}^{d^{(k)}} \delta_j^{(k)} \times w_{ij}^{(k)} \times \theta'(s_i^{(k-1)})$$

$$\Rightarrow \delta_i^{(k-1)} = (1 - \pi_i^{(k-1)})^2 \sum_{j=1}^{d^{(k)}} w_{ij}^{(k)} \delta_j^{(k)} \quad // \text{For tanh}$$

Overfitting

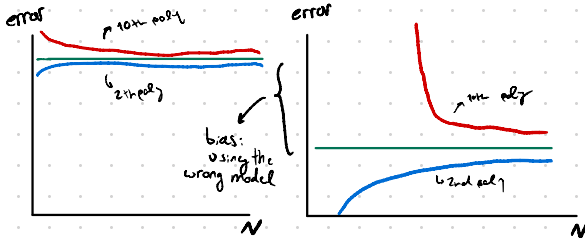
> Overfitting is what separate ML students from professionals.

Overfitting ≠ Bad generalization

↳ You tried too hard

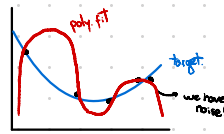
Def: Overfitting: fitting the data more than is warranted

→ Criterion: fitting the noise (harmful)

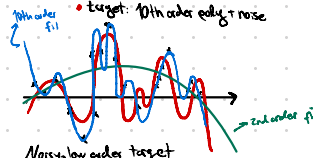


> With enough N , the H_{10} model is a better fit. But before, we had smaller error out of sample with H_2 .

Example:

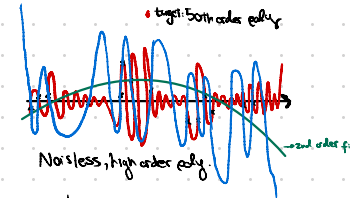
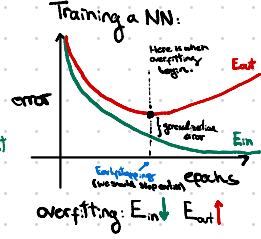


> Using a 4th order polynomial is an overfit.



Noisy, low order target

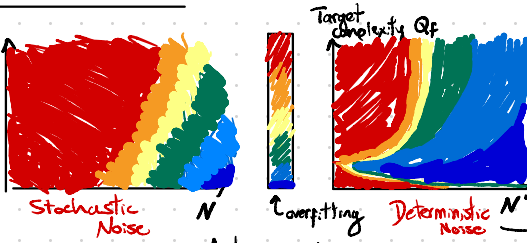
	2nd order	10th order
E_{in}	0.015	0.034
E_{out}	0.127	9.00



	2nd order	10th order
E_{in}	0.029	10^{-3}
E_{out}	0.120	7.680

Why? If there's no noise!
↳ There is another type of noise...

Overfitting Heatmap



number of data points $N \uparrow$: Overfitting \downarrow
 stochastic noise \uparrow : Overfitting \uparrow
 deterministic noise \uparrow : Overfitting \uparrow

Deterministic noise: the part of H cannot capture.
 • It depends on H and is fixed for a given x .

Bias variance decomposition with noise:

$$\underbrace{\mathbb{E}_{\mathcal{D}, x} [(g(x) - \bar{g}(x))^2]}_{\text{var}} + \underbrace{\mathbb{E}_x [\bar{g}(x) - f(x)]^2}_{\text{bias}} + \underbrace{\mathbb{E}_{\mathcal{E}, x} [(\epsilon(x))^2]}_{\sigma^2}$$

Fitting the noise is like fitting to a pattern that doesn't exist.

The deterministic noise comes from the limitations of our hypothesis set H .

Regularization

Instead of having wild lines, we want mild lines.

Two approaches

Math: Ill-posed problems in function approximations. Problem: The assumptions don't hold in practice sometimes !!

A tool to put a break on overfitting.

We greatly reduce the variance at the cost of mildly increasing the bias.

EXAMPLE

Def: Let L_n the n th Legendre polynomial of order n . They are orthogonal to each other. $L_1 = x, L_2 = \frac{1}{2}(3x^2 - 1), L_3 = (5x^3 - 3x)$

$$Z = \begin{bmatrix} 1 \\ \vdots \\ L_q(x) \end{bmatrix} \quad \mathcal{H}_q = \left\{ \sum_{q=0}^q w_q L_q(x) \right\}$$

> Legendre polynomials for an hypothesis set.

Given $(x_1, y_1), \dots, (x_n, y_n) = (z_1, y_1), \dots, (z_n, y_n)$

Minimize $E_{in}(w) = \frac{1}{N} \sum_{n=1}^N (w^T z_n - y_n)^2 = \frac{1}{N} (Zw - y)^T (Zw - y)$
 $\rightarrow w_{in} = (Z^T Z)^{-1} Z^T y \quad // \hat{G}_{in}$

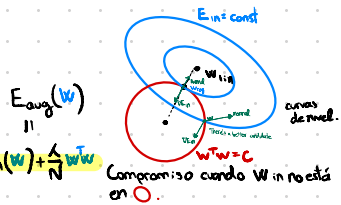
However, if we constrain the weights:

Obs: \mathcal{H}_2 is a constrained version of \mathcal{H}_3 // Hard constraint: $w_2 = 0 \forall q > 2$.
 a softer constraint: $\sum_{q=0}^q w_q^2 \leq C \leftrightarrow \min \frac{1}{N} (Zw - y)^T (Zw - y)$
 st $w^T w \leq C$

- We call this sol w_{reg}
- This is easily solvable using KKT multipliers.

$$\nabla E_{in}(w_{reg}) \propto -w_{reg} = -\frac{2\lambda}{N} w_{reg}$$

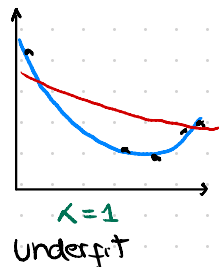
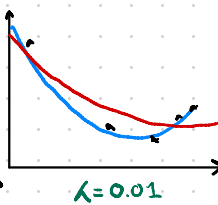
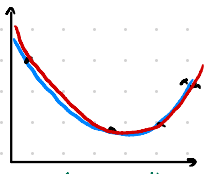
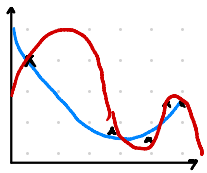
$$\nabla E_{in}(w_{reg}) + \frac{2\lambda}{N} w_{reg} = 0 \rightarrow \text{minimize } E_{in}(w) + \frac{\lambda}{N} w^T w$$



> Minimizing $E_{avg}(w)$ unconditionally is equiv to $\begin{cases} \min E_{in}(w) \\ \text{st } w^T w \leq C \end{cases}$

The final solution is $Z^T(Zw - y) + \lambda w = 0$
 $\leftrightarrow w_{reg} = (Z^T Z + \lambda I)^{-1} Z^T y$ // with reg.
 $w_{in} = (Z^T Z)^{-1} Z^T y$ // without reg.

One shot learning with regularization!



This regularizer is called **WEIGHT DECAY**.

Gradient descent $w(t+1) = w(t) - \eta \nabla E_{in}(w(t)) - 2\eta \frac{\lambda}{N} w(t)$
 $= w(t) (1 - 2\eta \frac{\lambda}{N}) - \eta \nabla E_{in}(w(t))$

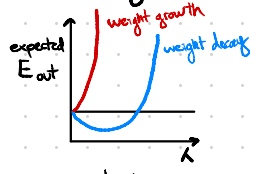
Variations of weight decay: $\rightarrow 1$ - Emphasize certain weights: $\sum_{q=1}^Q \gamma_q w_q^2$

Examples: $\gamma_1 = 2^2$ // low order fit.
 $\gamma_2 = 2^{-2}$ // high order fit

> In Neural Networks: Different γ for each layer.

\rightarrow Tikhonov regularizer: $W^T \Gamma W$

2- We play the inverse: Weight growth!



X Terrible idea.
 Best $\lambda = 0$.

Rule of thumb: stochastic noise is "high frequency"
 deterministic noise is also non-smooth

\Rightarrow constrain learning towards smoother hypothesis

In general smaller weights \rightarrow smoother hypothesis

> We want to penalize the noise, not the signal.

General form of Augmented Error

Call reg $\Omega(h) = \Omega$, we minimize $E_{out}(h) = E_{in}(h) + \frac{\lambda}{2} \Omega(h)$

E_{out} is better than E_{in} as a proxy for E_{out} .

$$E_{out}(h) \leq E_{in}(h) + \Omega(h)$$

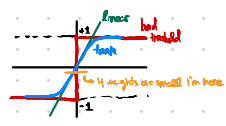
Choosing A Regularizer

> We move to smoother postures because the noise is not smooth.

> If Ω is bad, we have λ to cheer. If we have it wrong.

Neural Network regularizers

Weight decay: If all weights are small we end up with a linear function convoluted (if linear, if linear \rightarrow sig error). As weights increase we get a full sigmoid function.



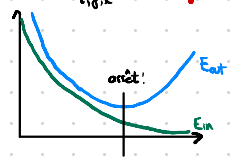
Weight elimination: fewer weights \rightarrow smaller VC dimension

Soft weight elimination: $\Omega(w) = \sum_{i,j=1}^n \frac{(w_{ij})^2}{\beta^2 + (w_{ij})^2}$

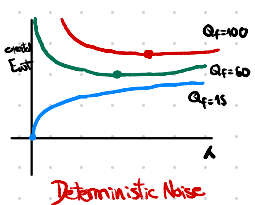
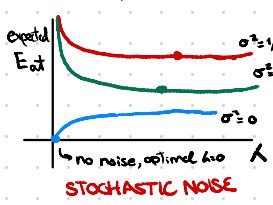
// Big weights are left alone, small weights are pushed toward zero.

Early stopping as a Regularizer

> Regularization through the optimizer!
 > When to stop?



The optimal λ



> Deterministic noise behaves almost exactly as if it were stochastic noise.

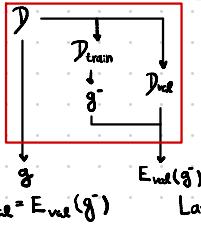
Validation

On act of sample point (x, y) the error is $e(h(x, y))$

$$\rightarrow \mathbb{E}[e(h(x, y))] = \mathbb{E}_{out}(h) \\ \text{var}[e(h(x, y))] = \sigma^2$$

Is K is too big you are giving a reliable estimate of the error for a very poor model.

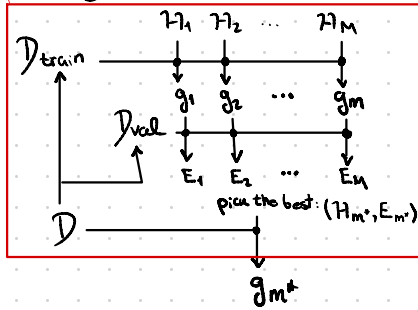
$$\mathcal{D} \rightarrow \mathcal{D}_{train} \cup \mathcal{D}_{val} \\ \mathcal{D} \rightarrow \mathcal{g}, \mathcal{D}_{train} \rightarrow \mathcal{g}^-$$



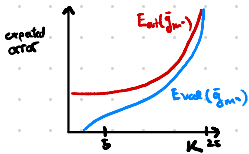
$$E_{val} = E_{val}(\mathcal{g}^-) \quad \text{Large } K \rightarrow \text{bad estimate.}$$

> \mathcal{D}_{val} is used to make learning choices
If an estimate of E_{out} affects learning:
The set is no longer a TEST set.
It becomes a VALIDATION set!

Using \mathcal{D}_{val} many times



We selected the model H_{m^*} using \mathcal{D}_{val}
 $E_{val}(\mathcal{g}_{m^*})$ is a biased estimate of $E_{out}(\mathcal{g}_{m^*})$



On a validation set $(x_1, y_1), \dots, (x_n, y_n)$ the error is $E_{val}(h) = \frac{1}{n} \sum e(h(x_i), y_i)$

$$\rightarrow \mathbb{E}[E_{val}(h)] = E_{out}(h)$$

$$\text{var}[E_{val}(h)] = \text{var}\left[\frac{1}{n} \sum e(h(x_i), y_i)\right] \\ = \frac{1}{n} \sum \text{var}[e(h(x_i), y_i)] = \frac{\sigma^2}{n} \rightarrow \text{Big n word!}$$

$$\Rightarrow E_{val}(h) = E_{out}(h) \pm \underbrace{\mathcal{O}\left(\frac{1}{\sqrt{n}}\right)}_{\text{error}}$$

Rule of thumb: $K = \frac{N}{5}$

Validation \neq Regularization

$$E_{out}(h) = E_{in}(h) + \text{overfit-penalty}$$

validation estimates this $\rightarrow E_{out}(h) = E_{in}(h) + \text{overfit-penalty}$
reg estimates this

Problem: K is not given. Each point in validation is a point less for training.

K points validation // \mathcal{D}_{val}
 $N-K$ points training // \mathcal{D}_{train}

• Test set is unbiased; validation set has optimistic bias.

• Example: Two hyp $h_1, h_2 / E_{out}(h_1) = E_{out}(h_2) = 1/2$
Pick $h \in \{h_1, h_2\} / e = \min(e_1, e_2)$
 $\mathbb{E}(e) < 1/2$ (optimistic bias)

M models H_1, \dots, H_M

Use \mathcal{D}_{train} to learn g_m from each model.

$$E_m = E_{val}(g_m)$$

we pick $m = m^*$ with the smallest E_m .

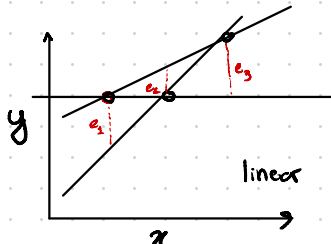
\mathcal{D}_{val} is used for "training" on the final test model

$$\mathcal{H}_{val} = \{\mathcal{g}_1, \dots, \mathcal{g}_M\}$$

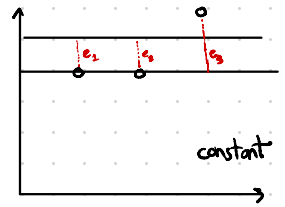
$$E_{out}(\mathcal{g}_{m^*}) \leq E_{val}(\mathcal{g}_{m^*}) + \mathcal{O}\left(\sqrt{\frac{\ln M}{n}}\right)$$

regularization λ , early-stopping T
 \rightarrow each give a degree of freedom: \mathcal{D}_{val}

Data contamination $E_{in}, E_{test}, E_{val}$
 > Optimistic (deceptive) bias in estimating E_{out}
 Training set: totally contaminated
 Test set: totally clean
 Validation set: "slightly contaminated"



$$E_{cv} = \frac{1}{3}(e_1 + e_2 + e_3)$$



$$E_{cv} = \frac{1}{3}(e_1 + e_2 + e_3)$$

Cross-Validation

$$E_{out}(g) \approx E_{out}(g^*) \approx E_{val}(g^*)$$

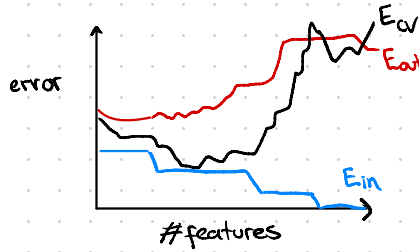
(small K) (large K)

> We want small and big K .
 $N-1$ points for training, 1 for validation
 $D_n = (x_1, y_1) \dots (x_{n-1}, y_{n-1}) \dots (x_n, y_n)$
 Fitted hypothesis learned from $D_n: g_n$
 $e_n = E_{val}(g_n) = e(g(x_n), y_n)$

$$\text{cross val error: } E_{cv} = \frac{1}{N} \sum_{n=1}^N e_n$$

catch: e_n 's aren't indep!
 but not that for avg...

> constant model has less E_{cv} !

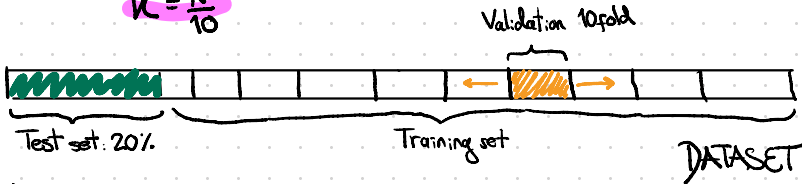


In practice the leave one out is inefficient as N increases

Sol: Take a chunk! $\rightarrow \frac{N}{K}$ training sessions on $N-K$ points each.

Rule of thumb: 10-fold cross validation

$$K = \frac{N}{10}$$

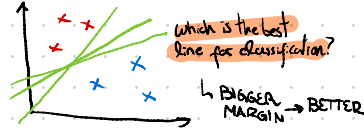


Validation is not data snooping
 as we have accounted for it.

SVM

Support Vector machine.

- Maximize margin.
- The solution.
- Nonlinear transforms. (When data is not lin sep.)



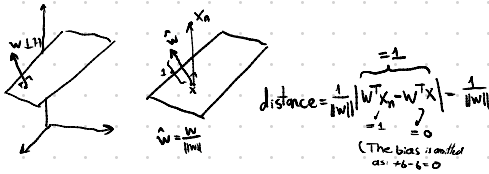
> Best model before deep learning.

Dichotomies with fat margin
 • If we restrict the margin, we get a smaller VC dimension. (Restrict num. of dichotomies)

Let x_n be the nearest point to the plane $w^T x + b = 0$

Technicalities \rightarrow Normalize w :
 $|w^T x_n + b| = 1$
 Take out w_0 param (as we did here)
 $w^T x + b = 0$ (w/ x_0)

To compute the distance:



The optimization problem

$$(1) \begin{cases} \max \frac{1}{\|w\|} \\ \text{st} \min_{n=1, \dots, N} |w^T x_n + b| = 1 \end{cases}$$

problem analysis.

$$(1. Alt) \begin{cases} \min \frac{1}{2} w^T w \\ \text{st} y_n (w^T x_n + b) \geq 1 \quad \forall n \end{cases}$$

this is prdble.

Obs: $|w^T x_n + b| = y_n (w^T x_n + b)$ // we only take w that classifies correctly.

This is a Constrained Optimization problem!

But we need equalities... $\ddot{\sim}$ \rightarrow Add slack variables!

This is similar to the regularization problem

	TO OPTIMIZE	CONSTRAINT
Regularization	E_{in}	$w^T w$
SVM	$w^T w$	E_{in}

We thus formulate the Lagrangian of an problem

$$(1) \begin{cases} \min \frac{1}{2} w^T w \\ \text{st} y_n (w^T x_n + b) \geq 1 \end{cases} \xrightarrow{L} \min [L(w, b, \alpha)] = \frac{1}{2} w^T w - \sum \alpha_n (y_n (w^T x_n + b) - 1)$$

w.r.t to w and b and maximise w.r.t each $\alpha_n \geq 0$

$$\nabla_w L = w - \sum \alpha_n y_n x_n = 0 \Leftrightarrow w = \sum \alpha_n y_n x_n \quad (1)$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_n y_n = 0 \Leftrightarrow \sum \alpha_n y_n = 0 \quad (2)$$

We write the dual formulation of the problem

Substitute (1) & (2) in $L(w, b, \alpha)$

$$L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m y_n y_m x_n^T x_m$$

where maximize w.r.t α subject to $\alpha_n \geq 0 \quad \forall n = 1, \dots, N$
 low w.r.t respect to α and $\sum \alpha_n y_n = 0$

The solution to the problem (quadratic programming)

We have $\max_{\alpha} \sum \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m$

prog. level pers. for fun. are quadratics.

$$\min_{\alpha} \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m - \sum \alpha_n$$

It's just $\begin{cases} \min_{\alpha} \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \\ \text{st } y^T \alpha = 0, \alpha \geq 0 \end{cases}$

$$\begin{cases} \min_{\alpha} \frac{1}{2} \alpha^T \begin{bmatrix} y_1 y_1 x_1^T x_1 & \dots & y_1 x_N x_1^T x_N \\ \vdots & \ddots & \vdots \\ y_N y_1 x_1^T x_1 & \dots & y_N y_N x_N^T x_N \end{bmatrix} \alpha + (-\mathbf{1})^T \alpha \\ \text{st } y^T \alpha = 0 \end{cases}$$

Obs: It's a convex function. No problem in optimising.

$$0 \leq \alpha_n \leq \infty$$

Obs: When N is too large (N > 10000) we need to use other heuristics to solve this.

Using Quadratic programming, we get the solution $\alpha = \alpha_1, \dots, \alpha_N$

$$W = \sum \alpha_n y_n x_n$$

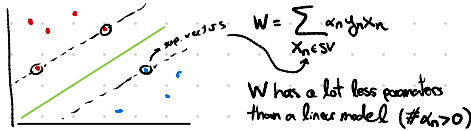
KKT conditions:

For $n=1, \dots, N$

$$\alpha_n (y_n (w^T x_n + b) - 1) = 0$$

There are going to be some $\alpha_n > 0 \Rightarrow$ Their respective x_n are the support vectors

Support Vectors



Now we solve for b using any SV

$$y_n (w^T x_n + b) = 1$$

non-Linear Transformations $X \rightarrow Z$

We work with z instead of x : $z_i =$

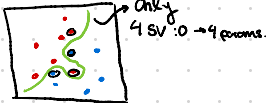
$$\rightarrow \mathcal{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m z_n^T z_m$$

Obs: $\# \alpha_n$ has nothing to do with the dimension of x or z , only N : the data points.

The support vectors in Z space

• Support vectors live in Z space

• In X space, are images of SV



$$\mathbb{E}[E_{out}] \leq \frac{\mathbb{E}[\# \text{ of SV's}]}{N-1}$$

SVM allows you to go very sophisticated without fully paying the price of it.

Kernel Methods

1. Kernel trick
2. Support vector

> Introducing: Infinite dim. spaces!

→ The kernel trick.

$$L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m Z_n^T Z_m$$

st $\alpha_n \geq 0, \sum_{n=1}^N \alpha_n y_n = 0$

$$g(x) = \text{sign}(w^T Z + b), \text{ where } w = \sum_{Z_n \text{ is SV}} \alpha_n y_n Z_n \text{ and } b: y_m (w^T Z_m + b) = 1$$

Given two points x and $x' \in \mathcal{X}$, get $Z^T Z'$?

$$Z^T Z' = K(x, x') \quad (\text{the kernel})$$

↳ like what's that.

Example: $x = (x_1, x_2) \rightarrow 2d$ and $x' = (x'_1, x'_2)$
 $Z = \mathbb{R}(x) = (1, x_1, x_2, \sqrt{2}x_1, \sqrt{2}x_2)$
 $K(x, x') = Z^T Z' = 1 + x_1 x'_1 + x_2 x'_2 + \sqrt{2}x_1 x'_1 + \sqrt{2}x_2 x'_2$

Let's try: can we compute $K(x, x')$ without transforming x and x' ?

Ex: Consider: $K(x, x') = (1 + x^T x')^2 = (1 + x_1 x'_1 + x_2 x'_2)^2$
 $= 1 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 + 2x_2 x'_2 + 2x_1 x'_1 x_2 x'_2$
 This is an inner p.
 $(1, x_1, x_2, \sqrt{2}x_1, \sqrt{2}x_2, \sqrt{2}x_1 x_2)$
 $(1, x'_1, x'_2, \sqrt{2}x'_1, \sqrt{2}x'_2, \sqrt{2}x'_1 x'_2)$

The polynomial kernel

$\mathcal{X} = \mathbb{R}^d, \Phi: \mathcal{X} \rightarrow \mathcal{Z}$ is a polynomial of order Q

The equiv. kernel $K(x, x') = (1 + x^T x')^Q = (1 + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^Q$
 It can be very ugly for $d=10, Q=100$

We can adjust coeff. const: $K(x, x') = (c + x^T x' + b)^Q$

Now, let's try without knowing what Z is.

Let $K(x, x')$ be the inner product in some space \mathcal{Z}

Ex: $K(x, x') = \exp(-\gamma \|x - x'\|^2)$

The corresponding Z of this kernel is infinite dimensional.

We compute a finite inner product from an infinite dimensional space!

Take an example: $K(x, x') = e^{-(x-x')^2} = e^{-x^2} e^{x^2} = \sum_{k=0}^{\infty} \frac{(-x^2)^k}{k!} \sum_{l=0}^{\infty} \frac{(x'^2)^l}{l!}$
 $= \sum_{k=0}^{\infty} \sum_{l=0}^{\infty} \frac{(-1)^k}{k! l!} x^{2k} x'^{2l}$
 (Taylor series, for exp)

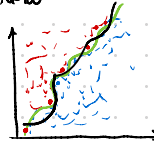
How do we use it?

Instead of passing y_i, y_j, x_i, x_j to the quadratic programming, we pass $y_i y_j K(x_i, x_j)$

Express $g(x) = \text{sign}(w^T Z + b)$ in terms of K ↳ Depends on Z_n .

$$w = \sum_{Z_n \text{ is SV}} \alpha_n y_n Z_n \rightarrow g(x) = \text{sign} \left(\sum_{\alpha_n > 0} \alpha_n y_n K(x_n, x) + b \right)$$

where $b = y_m - \sum_{\alpha_n > 0} \alpha_n y_n K(x_n, x_m)$ for any $\alpha_n > 0$.



Kernel in action:

How can we know, given a kernel candidate, if its respective Z exists.

- 3 approaches
- ↳ Construction. ↳ when we want to test a new kernel.
 - ↳ Math prop of the kernel (Mercer's conditions)
 - ↳ Who cares? I never visit it. ↳ sometimes you succeed and

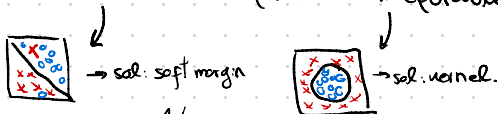
Design your own kernel

$K(x, x')$ is a valid kernel \iff It is symmetric \wedge $\begin{bmatrix} K(x_1, x_1) & \dots & K(x_1, x_N) \\ \vdots & \ddots & \vdots \\ K(x_N, x_1) & \dots & K(x_N, x_N) \end{bmatrix}$ is positive semi-definite $\forall x_1, \dots, x_N$

C'est les conditions de Mercer.

Soft Margin SVM

Data can be SLIGHTLY non separable to inseparable.



Margin Violation: A point may violate the margin \rightarrow the cushion...

Quantify: $y_n(w^T x_n + b) \geq 1 - \epsilon_n$, $\epsilon_n \geq 0$

Total violation: $\sum_{n=1}^N \epsilon_n$

Constant that defines the importance between w and $\sum \epsilon_n$

The new optimization: $\begin{cases} \min w^T w + C \sum \epsilon_n \\ \text{st } y_n(w^T x_n + b) \geq 1 - \epsilon_n \\ \epsilon_n \geq 0 \end{cases}$

We go to the Lagrangian again...

$$L(w, b, \epsilon, \alpha, \beta) = \frac{1}{2} w^T w + C \sum_{n=1}^N \epsilon_n - \sum_{n=1}^N \alpha_n (y_n(w^T x_n + b) - 1 + \epsilon_n) - \sum_{n=1}^N \beta_n \epsilon_n$$

$$\nabla_w L = w - \sum \alpha_n y_n x_n = 0$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_n y_n = 0$$

$$\frac{\partial L}{\partial \epsilon_n} = C - \alpha_n - \beta_n = 0$$

$$\alpha_n \leq C$$

So the solution is: $\begin{cases} \max L(\alpha) = \sum \alpha_n - \frac{1}{2} \sum_{n=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m \\ \text{st } 0 \leq \alpha_n \leq C, \sum \alpha_n y_n = 0 \end{cases}$

$$\rightarrow w = \sum \alpha_n y_n x_n$$

which minimises: $\frac{1}{2} w^T w + C \sum \epsilon_n$

Types of support vectors

- margin support vectors ($0 < \alpha_n < C$)
 $\hookrightarrow y_n(w^T x_n + b) = 1$ ($\epsilon_n = 0$)
- non-margin support vectors ($\alpha_n = C$)
 $y_n(w^T x_n + b) < 1$ ($\epsilon_n > 0$)

C is defined using cross-validation

Two technical observations

- $\hookrightarrow 1$ - What if data non-linearly sep?
- $\hookrightarrow 2$ - \bar{x} : what about the bias (w_0)?
All goes to 1 and $w_0 \rightarrow 0$

Radial Basis Functions (RBF)

Each $(x_n, y_n) \in D$ influences $h(x)$ based on $\|x - x_n\|$ ^{radial}
 ↳ closer \rightarrow + influence.

Standard form:
$$h(x) = \sum_{n=1}^N w_n e^{-\gamma \|x - x_n\|^2}$$
 ↳ basis function.

The Learning Algorithm

(Finding w_i).

* I have N datapoints, I'm trying to learn N params

We want $E_{in} = 0$, $h(x_n) = y_n$
 $\rightarrow \sum_{n=1}^N w_n e^{-\gamma \|x_n - x_n\|^2} = y_n$
 ↳ N eq., N unknown!

$$\begin{bmatrix} \exp(-\gamma \|x_1 - x_1\|^2) & \dots & \exp(-\gamma \|x_1 - x_N\|^2) \\ \vdots & & \vdots \\ \exp(-\gamma \|x_N - x_1\|^2) & \dots & \exp(-\gamma \|x_N - x_N\|^2) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Φ w y

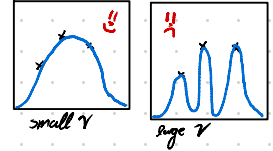
• If Φ is inv $\rightarrow w = \Phi^{-1} y$ (perfect interpolation)

• Often RBF is used for classification.

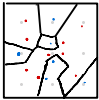
$$h(x) = \text{sgn} \left(\sum_{n=1}^N w_n e^{-\gamma \|x_n - x\|^2} \right)$$

↳ min $(s - y)^2$ on $(y = \pm 1)$
 $h(x) = \text{sign}(s)$

Effect of γ



RBF ~ KNN



• Instead of using a gaussian for the basis function. Lets use a function that is like a cylinder. ($u = 1$)
 The bigger u , the smoother should be the basis function

RBF with K centers

N parameters, w_1, \dots, w_N based on N datapoints.

Sol: Use $K \ll N$ centers: μ_1, \dots, μ_K instead of x_1, \dots, x_N

$$h(x) = \sum_{k=1}^K w_k \exp(-\gamma \|x - \mu_k\|^2)$$

- How can I choose μ_k ? \rightarrow Min distance between each x_n and its closest centroid μ_k .
- How can I choose the weights?

↳ **Clust K-means clustering.**

Split x_1, \dots, x_N into S_1, \dots, S_K

minimise $\sum_{k=1}^K \sum_{x_n \in S_k} \|x_n - \mu_k\|^2$

• Example of unsupervised learning
 • NP-hard problem

Find the minimum is NP-hard, we need an algorithm to get a decent sol

↳ **Lloyd's Algorithm:** Iteratively minimise:

$$\sum_{k=1}^K \sum_{x_n \in S_k} \|x_n - \mu_k\|^2 \text{ w.r.t } \mu_k$$

$$\mu_k \leftarrow \frac{1}{|S_k|} \sum_{x_n \in S_k} x_n$$

$$S_k \leftarrow \{x_n / \|x_n - \mu_k\| \leq \text{all } \|x_n - \mu_{l \neq k}\|\}$$

We converge to a local minimum, which depends on the initial clustering.

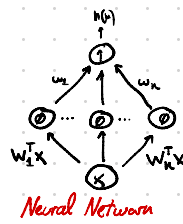
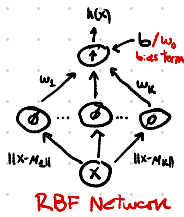
$\sum w_k e^{-\gamma \|x_n - \mu_k\|^2}$
 N eq., $K \ll N$ unknown

$$\begin{bmatrix} \exp(-\gamma \|x_1 - \mu_1\|^2) & \dots & \exp(-\gamma \|x_1 - \mu_K\|^2) \\ \vdots & & \vdots \\ \exp(-\gamma \|x_N - \mu_1\|^2) & \dots & \exp(-\gamma \|x_N - \mu_K\|^2) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix} \approx \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

If $\Phi^T \Phi$ inv $\rightarrow w = (\Phi^T \Phi)^{-1} \Phi^T y$

RBF Network

Features: $e^{-\gamma \|x - M_k\|^2}$
 Nonlinear transform depends on γ



Choosing γ

$$h(x) = \sum_{k=1}^K w_k e^{-\gamma \|x - M_k\|^2}$$

Expectation maximisation

Iterative approach (~EM algorithm in mixture of Gaussians)

- 1 - Fix γ , solve for W
- 2 - Fix W , minimise error w.r.t γ

Obs: We can have a different γ_k for each M_k .

RBF & Regularization

RBF can be derived purely for regularization

$$\sum_{k=1}^N (h(x_k) - y_k)^2 + \lambda \int_{-\infty}^{\infty} \left(\frac{dh}{dx} \right)^2 dx$$

This is the smoothest interpolation. Which solution is RBF approx.

SVM kernel implements:

$$\text{sign} \left(\sum_{k: y_k > 0} \alpha_k y_k e^{-\gamma \|x - x_k\|^2} + b \right)$$

RBF implements

$$\text{sign} \left(\sum_{k=1}^K w_k e^{-\gamma \|x - M_k\|^2} + b \right)$$

Three Learning Principles

- Occam's razor
- Sampling bias
- Data snooping

> An explanation of the data should be as simple as possible, but no simpler (A. Einstein)

> The razor of Occam: The simplest model that fits all the data is also the most plausible. describe with as few bits as possible.

↳ What is simple?
 ↳ Why simple → better?

Simpler → better E-out performance

Measures of complexity

- Complexity of h → MDL (minimum desc length)
 Order of a polynomial.
- Complexity of H → Entropy
 VC-dimension

Fewer simple hyp than complex ones are less likely to fit the data: $\frac{M(h)}{2^N}$
 So when it happens it's more significant.

• Sampling bias: If the data is sampled in a biased way, learning will produce a similarly biased outcome.

When we think of simple, it's in terms of h . Proof use simple in terms of H ...

If you have a population that was not used in training, we have a sampling bias if they appear on testing.

• Data snooping → If the data has affected any step in the learning process, its ability to assess the outcome has been COMPROMISED

* Most common trap - many ways to slip.

They are related to the principle of counting:

λ bits specify $h \rightarrow h$ is one of 2^λ elements of H

Exception: SVM; looses complex but is one of few.

Examples of Snooping

1. → Looking at the Dataset: $Z = (z_1, x_1, x_2, x_3, x_4, x_5, x_6)$. Why not $z = (z_1, x_1^2, x_2^2)$? The data fits good (just a circle)
 Note: you can take into account info you have of f (sym in x_2, x_3), while you not look at D .

2. → Trying one model after the other **ON THE SAME DATASET**, you will eventually succeed.

> If you torture the data long enough, it will confess.

- VC dimension of the total learning model.
- May include what others tried! // This paper says that SVM with $\gamma = 1/2$ is perfect!

• The problem is matching for a PARTICULAR dataset

1.5: Example: $\Delta_{1-20} \Delta_{1-1} \rightarrow \Delta_{1-0} \Rightarrow$ Normalize data, split D_{train}, D_{test}
from the past values of now!

Wrong order! ✓
 You give the mean and var of the test set to the algo.



Bayesian Learning

We extend probability to calculate $P(D|h=f)$

$$\rightarrow P(h=f|D) = \frac{P(D|h=f)P(h=f)}{P(D)} \propto P(D|h=f)P(h=f)$$

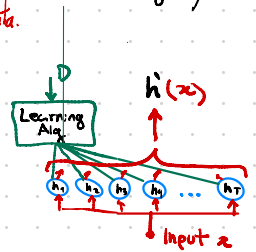
posterior likelihood prior

• Key: The prior is an assumption that can be wrong

Bayesian learning is justified when

- The prior is valid **trumps all other methods**
- The prior is irrelevant $\rightarrow \Pi_0 \dots \Pi_{100000}$ (doesn't matter Π_0 that much)

Knowing the prior we can find the most probable h given the data.
 Can we really calculate $P(h=f)$?



Aggregation

Combining many models h_1, \dots, h_T that were trained on D .

- Tactics:
 - Take an average. (Reg.)
 - Take a vote. (Classif.)
- > ensemble learning

Aggregation can be done

1. After the fact → combines existing solutions
2. Before the fact → creates solutions to be combined

↳ Boosting: create h_1, \dots, h_T which are decorrelated with the previous h that were created at each step
 (AdaBoost)
 We care about misclassifications in D
 Each model is weighted:

$$g(x) = \sum_{t=1}^T \alpha_t h_t(x)$$

We choose α_t s to minimize the error on the aggregated dataset