

# MACHINE LEARNING

= P.II 1: Machine Learning =

= Warm up (I)

One may measure intelligence through the number of available solutions or new approaches to problems.

= Warm up (II)

Philosophical assumptions:

- that our mind is computationalism-natured
- functionalism: our mind state depends on
  - previous state
  - perception
  - ?

↳ we may replicate and build IA

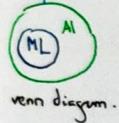
↳ counter-believes

religion, "emotions aren't computational"  $\Rightarrow$  • Weak AI: only some parts are replicable, the rest are imitated  
 • Narrow AI: what we generally refer to when mentioning AI

= Warm up (III)

We don't want new 'colleagues'  $\Rightarrow$  we want tools we may get along with

= Why ML now?



↳ Counterfactual " "  
 ↳ Interventional Causal-Effects: we could model the world into causal-effects, if everything was determined by causality.

Stats + ML

= Key elements for ml.  $\rightarrow$  pattern, data, not down-parallel mathematically

With no data, there's no ML. Without pattern, it'd be random and the result too.

= Types of ML

• Manifold learning: let's say a manifold is a subspace where we may apply the euclidean distance  
 Let's consider a subspace of emotions: interpolation

what's in the middle of (with solution )

The diagram shows five circular nodes representing emotions: a smiley face, a neutral face, a sad face, another neutral face, and a smiley face. A curved orange arrow points from the neutral face to the smiley face, labeled 'interpolation'. A green arrow points from the neutral face to the sad face, labeled 'manifold', with the text 'which we should follow to get the middle point' and a smiley face icon.

> Embeddings. text embeddings      Embeddings represent entities through vectors. (Feature vectors)

• Reinforcement training:  $\sim$  "how parents teach kids" ML models that learn games.

It learns to "win", not to carry out the actions or comprehend the complexity of what it's doing, if that makes sense.

= Let's discuss the following problems

- Supervised.
- Supervised through clustering + unsupervised through TPG.
- Unsupervised, but it could be rewritten as reinforcement. (student idea)
- Reinforcement.

= Let us discuss about the following problems (slide 21)

- Could be both (Regression/Classification): Learn close estimate of the model / Y/N for each possible action or which action is best
- Both: Most Likely zone / Y/N for each bit (subregion of image, sliding window)

= Pragmatic ML pipeline

The model should encode the data.

- Training
- Exploitation (Testing)



= Encoding data

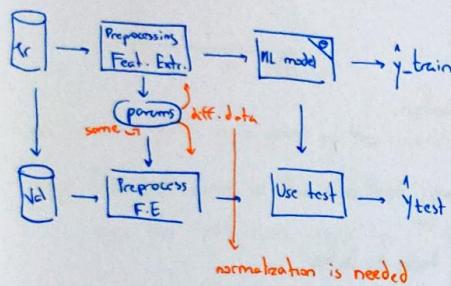
small  $\rightarrow 0$   
medium  $\rightarrow 1$  ✓  $\Rightarrow ?$   
big  $\rightarrow 2$   
↓ categorical  $\downarrow$  order

red  $\rightarrow 0$   
blue  $\rightarrow 1$   
green  $\rightarrow 2$   
↓ cat.  
↑ RGB value could be useful  
no relationship between color and value  
this may not be that good

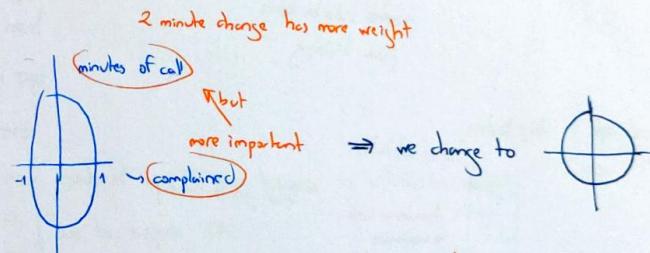
### • One-Hot Encoding:

	color		
	red?	blue?	green?
red	1	0	0
blue	0	1	0
green	0	0	1

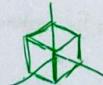
→ set-dummies!  
redundant in this case



• Normalization:  $\bar{x} = \frac{x - \mu_x}{\sigma_x}$  Standardization  
ONLY WITH FEW OUTLIERS  
alt.  
 $\bar{x} = \frac{x - \min(x)}{\max(x) - \min(x)}$



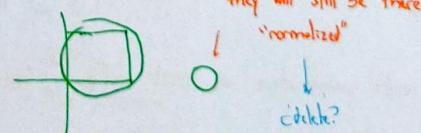
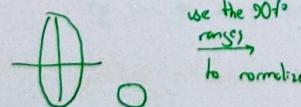
other examples: change to unit cube



$\bar{x}_{mc} = \frac{x_{mc} - \min(x_{mc})}{\max(x_{mc}) - \min(x_{mc})}$   
more to origin  
squish to (0,1)

this range squish might be mostly ruined by outliers.

• Identify outliers: imagine 10<sup>10</sup> outliers



important (21)  
 $F(\Delta) = \underline{\Delta} \underline{\Delta}$   
 $F(\square) = \underline{\square} \underline{\square}$  collisions  
 $F(\circ) = \underline{\circ} \underline{\circ}$  possible

• Feature Hashing: when embedding are not useful.

HASH: deterministic  $F(\text{obj}) \rightarrow$  gives always some number for an object within a customizable range  
↓ use as an index to store the object.

must be pseudo-random.

good practice for it to be in terms of the value to hash (deterministic) and the range (to fit in the boxes  $\rightarrow \circ$ )

We shouldn't use too much space: waste

not good encoding (as stated before)

⇒ representation vector: count the occurrences of embedding  $\leftarrow$  each possible hash value

sparse, efficient, complexity  $\downarrow$  collisions,  $\sim$  uniform, no metric notion, set dimensionality of embedding space

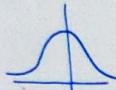
unit ball  $d=2$   $d=3$

...  $d=n$

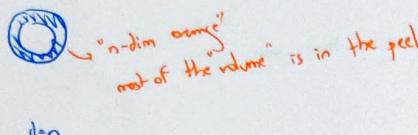


Gauss distribution,



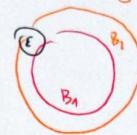
$\sigma$

$d=1$



$d=n$

hyperplane



$$V_{B_2} - V_{B_1} = \dots \gg \gg \gg \gg \gg \gg B_1$$

= Performance =

LOO represents an intersection of odds (in this location, this gender...)

↓ small piece where the value is.

Iterative LOO may take ages  $\Rightarrow$  limit iterations to  $K$  usually 10-fold or 2 times 5-fold  
split dataset into  $K$  pieces  $\rightarrow$  or as big as your data  $\downarrow$   $K$ -fold cross validation.

each piece used for testing in each it

but training:

uses  $K-1$  data

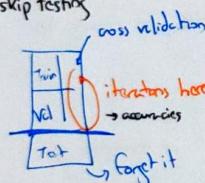
pure training

$\downarrow$   $K$  accuracies  $\rightarrow$  mean ...

$\downarrow$  but calculate: usamos todos los datos.

validation? opt. 1: skip testing

opt. 2:



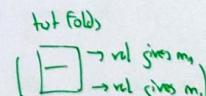
(used in business)

opt. 3 (academic): nested cross-validation: two loops: one for testing and the other for validation.

outer-loop: decides which data is taken for test (2-fold test)

inner-loop: validates with 3-fold val ...

not example } don't work to decide which model works better  
it just tells us the best acc. we may get



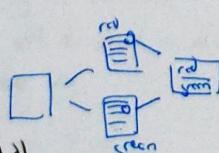
the model parameters are reset in each iteration (no accumulative training)

↳ we don't put the params together in the end, we just work with the accs.

then we train the model with all data and "most likely" its accuracy its going to be at least the accuracy in testing (as it was tested with less data)

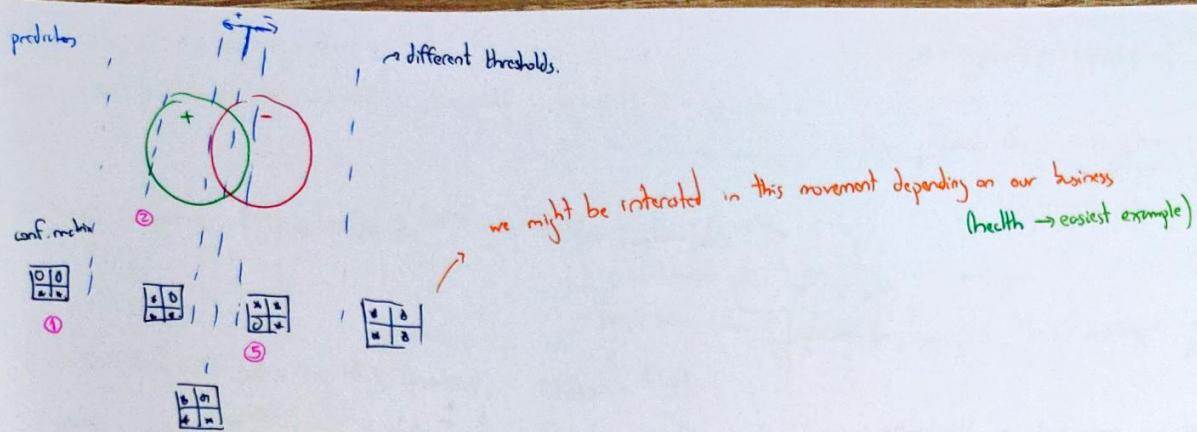
can use  
use it. to test  
(can mix, no he applies)

with most acc

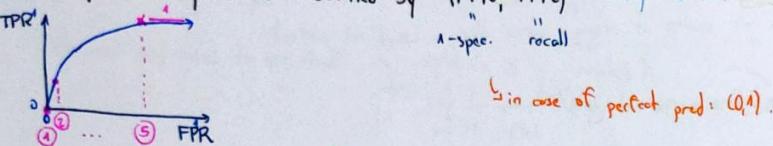


Stacked cross-validation

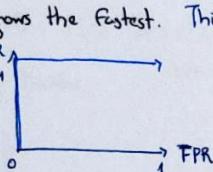
(multiple when some classes are misrepresented)



- **ROC:** can be represented by the curve defined by  $(FPR, TPR)$ . They are not complementary.



- An **operating point** is a point on the **operating curve** (which is chosen as threshold)
  - **prop:** The best curve is the one that grows the fastest. This may be identified by calculating the area. ①

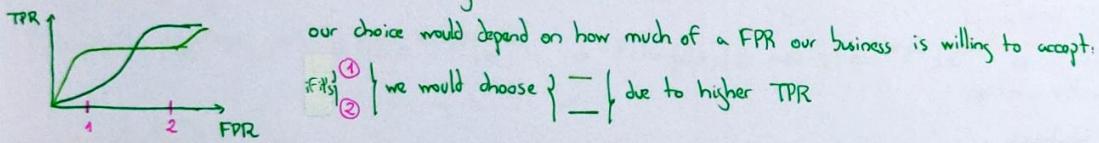


- ④ This might not be our decision criteria: imagine two ROCs



our choice would depend on how much of a FPR our business is willing to accept.

if it's ① we would choose ② due to higher TPR



= Campaign exercise:  $\begin{cases} \text{cost} = \text{all positives} \times 10 \\ \text{ben.} = \rightarrow \text{I guess!} \text{ TP} \times 100 \end{cases}$

$$\text{Cost: } \frac{\alpha}{10} \frac{\text{units}}{\text{person}} \times \frac{(\text{TP} + \text{FP})}{\text{ppl}} \text{ ppl}$$

per model predicted will recover

Beneficio:  $B \times TP \times 100$

models predict given in a document can be compared:  $\frac{TP}{TP + FN}$  (precision) |  $\frac{TP}{TP + FP}$  (recall)

(Loss would be  $Loss = 100 * (TP + FN)$ )

↳ **new que**  **learning and no campaign** 

Está en 0.9° aquí, ya que schumers al 10%

## =Pill 4: Opening The Blackbox (Directioning ML methods)

$$\hat{y} = h(x) = h(x; \theta)$$

↓  
 production  
 (output)

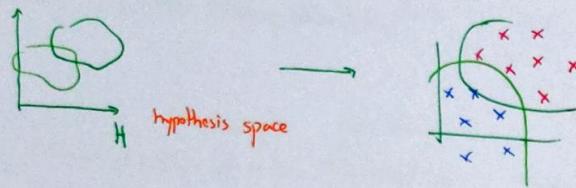
↓  
 model

↓  
 input

↓  
 set of parameters

Suppose that  $\hat{y} = h(x; \theta) = \theta_0 + \theta_1 x + \theta_2 x^2$

 parameter space is not generally used.



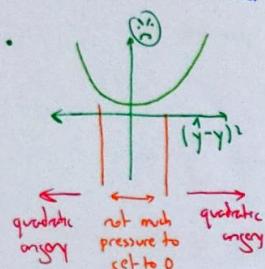
We are minimizing respect to or the following error in predictions notion:

↓  
loss generally

$$\frac{1}{N} \sum_{i=1}^N e(\hat{y}_i, y_i) + \dots$$

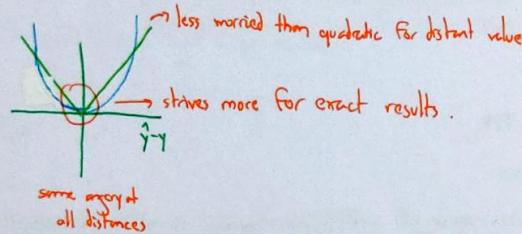
we may add other terms in the future

It is convenient to model an irritable function ~ make it "angry" when the value is not correct



"The more it differs, the angrier it gets"  $\rightarrow \text{LES} \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$

Linear Angry



some angry at all distances



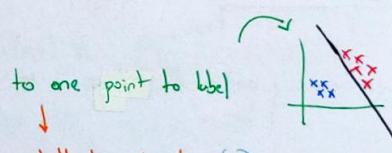
= Human learning algorithm exercise:

model class: linear functions.

bss: error loss  $\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$

→ this might give a fine that is too close to one point to label

✓ indicator: classical bss in classifier.



Remember:  $h(x) = w_1 x + w_0 = [w_1 \ w_0] \begin{bmatrix} x \\ 1 \end{bmatrix} = w^T \tilde{x}$  extended data that integrates the offset.

$$\text{minimize } \frac{1}{N} \sum_{i=1}^N (y_i - w^T \tilde{x})^2 = \frac{1}{N} (y - w^T \tilde{x})^T (y - w^T \tilde{x})$$

// suitable because

it's regression  
in this example

vectorized form  
multi derivative easy.

We need to get the gradient of the bss respect to parameters

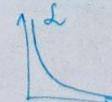
$$\nabla_w \mathcal{L} = -2 \tilde{x}^T (y - w^T \tilde{x})$$

vectorial gradient, independent from dimensionality

curve of descent gradient: convergence curve.

$$\nabla_w \mathcal{L} = \left[ \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_0} \right]$$

Expected loss behaviour



= Feasibility of the learning problem

$P(|x - \mu| \geq \epsilon) \leq \delta \Rightarrow \phi \rightarrow$  Hoeffding's Inequality (no explicit de donde side - yet - odds of mistake)

para  $N$  ha de ser muy grande para compensar una  $\epsilon$  muy pequeña

• Hoeffding's inequality:  $P(|x - \mu| > \epsilon) \leq 2e^{-2\epsilon^2 N} \Rightarrow P(|E_{in} - E_{out}| \geq \epsilon) \leq e^{-\epsilon^2 N}$

→ out-of sample error: expected error

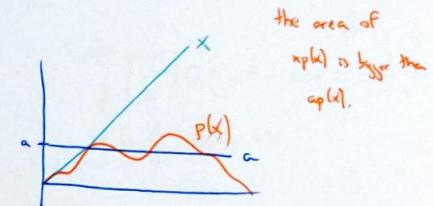
on sample error (freq of hypothesis being wrong)

- Markov's inequality: let  $x$  be (R.V) (positive),  $P(x \geq a) \leq \frac{E[x]}{a}$

proof:  $E[x] = \int_{-\infty}^{\infty} x p(x) dx = \int_0^{\infty} x p(x) dx = \int_0^a x p(x) dx + \int_a^{\infty} x p(x) dx \geq \int_a^{\infty} x p(x) dx$

We now only need to see that  $\int_a^{\infty} x p(x) dx \geq \int_a^{\infty} a p(x) dx$

$$\Rightarrow \int_a^{\infty} x p(x) dx \geq \int_a^{\infty} a p(x) dx = a \int_a^{\infty} p(x) dx = a P(x \geq a) \quad \blacksquare$$



- Chernoff bound: positive RV  $x$  and  $\epsilon > 0$ :  $P(|x - \mu| \geq \epsilon) \leq \min_{\lambda > 0} [E[e^{\lambda(x-\mu)}]] e^{-\lambda\epsilon}$

proof:  $P(|x - \mu| \geq \epsilon) \leq \frac{E[|x - \mu|]}{\epsilon} \text{ by Markov}$

the prob. is the same  $\Downarrow \Leftrightarrow$

$$P(e^{\lambda(x-\mu)} \geq \epsilon)$$

also  $\Downarrow \Leftrightarrow \lambda > 0$

$$P(|x - \mu| \geq \epsilon) = P(e^{\lambda(x-\mu)} \geq \epsilon) \leq \underbrace{E[e^{\lambda(x-\mu)}]}_{\text{For any } \lambda} e^{-\lambda\epsilon} \Rightarrow P(|x - \mu| \geq \epsilon) \leq \min_{\lambda > 0} [E[e^{\lambda(x-\mu)}]] e^{-\lambda\epsilon} \quad \blacksquare$$

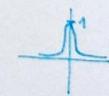
- Simplified Hoeffding's Inequality: let  $z_i$  be independent bounded RV such that  $z_i \in [a, b] \ \forall i$ :  $P\left(\frac{1}{N} \sum_{i=1}^N (x_i - \mu_i) > \epsilon\right) \leq e^{-\frac{\epsilon^2 N}{c}}$  where  $c = 8(b-a)$

to proof:  $P\left(\frac{1}{N} \sum_{i=1}^N z_i > \epsilon\right) \leq e^{-\frac{N\epsilon^2}{2}}$

distribution st.  $F(x) = \begin{cases} 1/2 & \text{if } x=1 \\ 0 & \text{otherwise} \end{cases}$

proof (for a binary case):  $z_i$  are Rademacher RV: in terms of Dirac delta function  $\rightarrow$  value 1 or -1 with prob 0.5

$$P(z_i) = \frac{1}{2} [\delta(z_i=1) + \delta(z_i=-1)] \rightarrow \text{Fair coin} \quad (\text{check what happens with } 0,1)$$



$$\Rightarrow P\left(\frac{1}{N} \sum_{i=1}^N z_i > \epsilon\right) = P\left(\sum_{i=1}^N z_i > N\epsilon\right) \leq E[e^{\lambda \sum z_i}] e^{-\lambda N\epsilon} = E\left[\prod_{i=1}^N e^{\lambda z_i}\right] e^{-\lambda N\epsilon} = \prod_{i=1}^N E[e^{\lambda z_i}] e^{-\lambda N\epsilon}$$

$$= [E[e^{\lambda z_i}]]^N e^{-\lambda N\epsilon} \leq \min_{\lambda > 0} e^{\lambda^2/2 - \lambda N\epsilon} = e^{-\epsilon^2 N/2} \quad \text{④}$$

$$\Rightarrow E[e^{\lambda z_i}] = \frac{1}{2} e^{\lambda} + \frac{1}{2} e^{-\lambda} = \frac{1}{2} (e^{\lambda} + e^{-\lambda}) = \frac{1}{2} \left( \sum_{i=1}^{\infty} \frac{\lambda^i}{i!} + \sum_{i=1}^{\infty} \frac{(-\lambda)^i}{i!} \right) = \frac{1}{2} \sum_{i=1}^{\infty} \frac{\lambda^i}{(2i)!} \leq \frac{\lambda^{\infty}}{2} = \sum_{i=1}^{\infty} \frac{(\lambda/2)^i}{i!} = e^{\lambda^2/2}$$

$$e^{\lambda} = \sum_{i=0}^{\infty} \frac{\lambda^i}{i!}$$

odd summands are hence zero

$$(2i)! \geq 2^i i!$$

From  $0 \dots N$  it's always the same

then  $2$  is always smaller than  $2^i$

(from  $i+1 \dots 2i: i \rightarrow 2 \dots 2^i$ )

$$\text{④ } \frac{\epsilon^2 N}{2} - \frac{\epsilon^2 N}{2} = -\frac{\epsilon^2 N}{2} \neq \left(\frac{N\epsilon}{2} - N\epsilon\right) = \phi$$

$$\lambda = \epsilon$$

- Modification: use max to bound the worst case scenario.

$\uparrow$   $H$ , better.

$\uparrow$   $H$ , worse.  $\rightarrow$  we will give more notions of the complexity of hypothesis space.

$\downarrow$  cardinality of hypothesis space.

- We may define that  $P[|E_{\text{out}} - E_{\text{in}}| \geq \epsilon] \leq 2|H| e^{-2N\epsilon} = \delta \rightarrow \textcircled{1}$

a.k.a with  $P=1-\delta$ ,  $|E_{\text{out}} - E_{\text{in}}| \leq \epsilon$  holds (complementary)

$$\textcircled{1} \Rightarrow \epsilon = \sqrt{\frac{-\log \delta/2|H|}{2N}} = \sqrt{\frac{\log |H| + \log \delta/2}{2N}} \approx O\left(\frac{|H|}{N}\right) \rightarrow \delta \text{ has a minimum effect.}$$

$\downarrow$   $C = |H| \circ C = f(\lambda_N)$   $\rightarrow$  said in posterior class (not sure)

$E_{\text{out}} > E_{\text{in}}$

$\downarrow$   $\rightarrow$  training error: we can control by optimizing (and overfitting and other terrible things)

generalization error: no control  $\leftarrow$  we expect that reducing training error

will reduce generalization error, we can't assure it though.

$$\left\{ E_{\text{in}} \leq E_{\text{out}} \leq E_{\text{in}} + O\left(\frac{|H|}{N}\right) \right\}$$

$\rightarrow$  one of the most important equations in training error. (equivalent to  $E_{\text{out}}(f) \approx E_{\text{in}}(f)$ )

$\downarrow$   $O$  in best case scenario  $\Rightarrow E_{\text{in}} = E_{\text{out}}$

= Using the bound: example on how to use  $|E_{\text{out}} - E_{\text{in}}| \leq \epsilon$  to know the needed amount of samples.

= The last consideration: we must move from function to probabilities some blue  $\rightarrow$  sky sea we now will consider the odds of that specific value to be  $\left\{ \begin{array}{l} \text{sky} \\ \text{sea} \end{array} \right\}$

•  $P(y^*|x^*, \bar{X}, \bar{Y})$  is the target distribution (what we want to learn)

$\downarrow$   $\begin{array}{l} \text{To known data} \\ \text{probability of the prediction} \\ \text{unseen input} \end{array}$

- We may consider that  $P(y^*|x^*) = \underbrace{\mathbb{E}[y|x]}_{\text{expected value}} + \underbrace{g(x)}_{\text{noise}}$  where the noise is  $g(x) = y - \mathbb{E}[y|x]$

$\downarrow$   $\begin{array}{l} \text{noise} \\ \text{noise} \\ y \\ \mathbb{E}[y|x] \end{array}$

-  $P(x)$  is the input distribution that quantifies the relative importance of  $x$ .

- note that  $P(x, y) = P(x)P(y|x) \neq P(y|x)$

$\downarrow$  most math. models aspire to learn this. Harder, though, because you need  $P(x)$  explicitly.

• Learning is to minimize the generalization error ( $E_{\text{out}} \rightarrow 0$ )  $\rightarrow$  done through our control on  $E_{\text{in}}$

- conditions:  $\textcircled{1} E_{\text{out}} \approx E_{\text{in}}$

$\textcircled{2} E_{\text{in}} \rightarrow 0$  : achieved with a training model on the data I have.

bad news: in linear models:  $|H| = \infty$

= Vapnik-Chervonenkis dimension

• Shatter: in a binary classification problem, if a model of a given family is able to classify it correctly

- VC-dim of a perceptron in 2D?

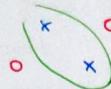
$$\left| \begin{array}{l} + 2 \text{ entries} \\ \hline \end{array} \right. \Rightarrow \left| \begin{array}{l} 0 \\ \hline x \end{array} \right.$$

$$\left| \begin{array}{l} + 3 \text{ entries} \\ \hline \end{array} \right. \Rightarrow \left| \begin{array}{l} 0 \\ \hline x \\ 0 \end{array} \right. \text{unpassable test} \Rightarrow \text{VC-dim of } | = 2$$

Now let's add tilt to the classifier:

4 entries:

Intuitively, an ellipsis will have higher VC:



- As an intuition,  $|H| \approx d_{vc}$   $\Rightarrow$  Final result formula (substitute  $H$ )
- $d_{vc} \approx 10 \text{ DF}$ , where DF := degrees of freedom.

3D hyperplane: "if you can do it in 2D you can do it in 3"  $\rightarrow$  3 entries +  $\Rightarrow$  VC-dim  $\geq 3$

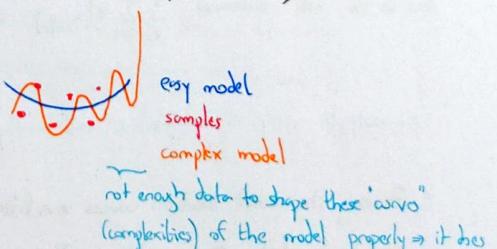
4 entries: In fact  $VC\text{-dim} = 4$ .

- Note that  $\dim \text{hyperplane } d_{vc}$
- |        |               |          |
|--------|---------------|----------|
| 2      | $\rightarrow$ | 3        |
| 3      | $\rightarrow$ | 4        |
| :      | :             | :        |
| $\{ N$ | $\rightarrow$ | $d+1 \}$ |

= Pill 6: Overfitting =

- Remember: you can regress from  $x_1, \dots, x_n$  or from  $x_1, \dots, x_1^m, x_2, \dots, x_2^m, \dots, x_n, \dots, x_n^m$  where m is lower than the dim of the model. ↴ training set extension ↴  $x_1^0, \dots, x_n^0$  should be included  $\{1, \dots, 1\}$  ↴ bias representation.
- It is not a good idea always to increase the order (degree of our model)  $\rightarrow$  we need enough data for that.  $\Rightarrow$  overfitting.\*

- \*  $E_{in}^{10} \leq E_{in}^3 \leq E_{out}^{10} \leq E_{out}^3$  instead of  $E_{in}^{10} \leq E_{in}^3$  because at some point  $E_{out}^{10}$  stopped following  $E_{in}^{10}$



In terms of classification, 0 errors  $\Rightarrow$  forcing to do so, will cause to bad generalization.

"purple islands"  $\Rightarrow$  reduce complexity, to avoid forcing

add data to be "closer to infinity" (sense of generalization)

↳ but with 0 error is terrible,  $\Rightarrow$  the islands are very small  $\Rightarrow$  no influence area to take correct decisions.

- bias: difference between the limits of training and testing error rates.

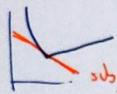
more complexity  $\Rightarrow$  the curves get closer together with less bias.

Attention:  $E_{in}(c) + O\left(\frac{c}{N}\right)$  we make it converge to a solution with worse errors.

- We want:  $E_{\text{out}} \approx E_{\text{in}} \Rightarrow$  balance between C and N at  $O\left(\frac{C}{N}\right)$
- $E_{\text{in}} \rightarrow 0 \Rightarrow 11C$   $\leftarrow$   $\nabla C$  not always a good solution.

= Pill 7: Stochastic Subgradient Methods = (next page)

Subgradient for



subgradient that supports (closest hyperplane to the point)

When we choose one of the violated constraints, we do it in order to move towards the feasible set.

But these steps are not guaranteed to be always forward.

If we find  $f_0(x^{t+1}) > f_0(x^t) \rightarrow$  Keep going but keep track best solution POCKET STRATEGY

Stochastic methods let you get past saddle points and local minima (there's some randomness)  
↳ most robust strategy so far.

momentum  $\rightarrow$  imitate the movement of letting a ball go.

Adagrad: avoid not arriving ever to the minimum (as it happens to desc. grad) by increasing the step in plain areas flat

If gradient was small in past it.  $\rightarrow D^t$  will increase the speed.

= Pill 8: Regularization - Combating Overfitting =

We know that  $E_{\text{out}} \leq E_{\text{in}}(C) + O\left(\frac{\sqrt{C}}{N}\right)$  let's add a term related to complexity  
 And so far, we minimize  $\frac{1}{N} \sum_{i=1}^N e(x_i, y_i) \Rightarrow$  minimize  $\frac{1}{N} \sum_{i=1}^N e(x_i, y_i) + \lambda C$

We should be worried if adaptative and analytical solutions are not close.

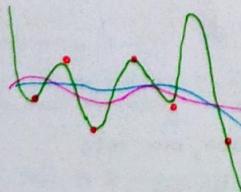
- Getting the exact solution causes overfitting  $\Rightarrow$  it forces to pass through all the entries.  $\Rightarrow$  EARLY STOPPING

④ Iterate during few epochs so the method does not memorize.

$x = \text{optimum} \rightarrow \text{"solves the problem"} \rightarrow \text{error} = 0 \rightarrow \text{overfitting.}$



learns / trains

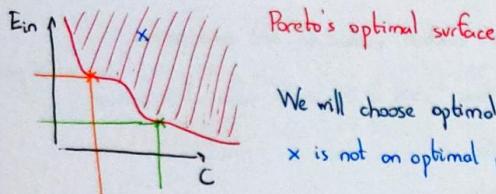


Early stopping model.

Stoch. Grad. Methods

Other alternative: Multi Objective Problem

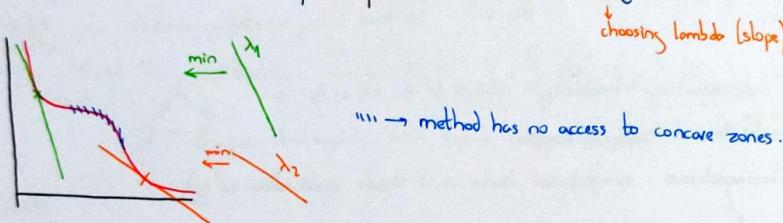
$(\frac{1}{N} \sum_{i=1}^N e(x_i, y_i), C) \rightarrow \text{minimize both at the same time}$



We will choose optimal points (they don't have adequate points in their 3rd quad:  $\times \times$ )  
 $\times$  is not an optimal point.

We should return the whole surface (you may choose from there a balance between training error and complexity) but that's not feasible:

- **Scalarization:** plot a line that describes the optimal points. minimize  $f_1 + \lambda f_2$

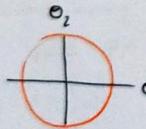


We know now about the  $\lambda$  we added to the minimization. What about  $C$ ?

Actually, let's consider minimize  $\frac{1}{N} \sum_{i=1}^N e(y_i, h(x_i, \theta)) + \lambda C(\theta)$

$$h(x, \theta) = \sum_{j=1}^J \theta_j h_j(x) \quad \begin{array}{l} \text{sparse approx. } \|\theta\|_1 \\ \text{"smooth" approx. } \|\theta\|_2 \end{array} \quad \left\{ \begin{array}{l} \|\theta\|_1 \text{ low} \\ \|\theta\|_1 + \|\theta\|_2 + \epsilon + \delta = 1 \\ \|\theta\|_2 \text{ low} \\ \epsilon, \delta \neq 0 \end{array} \right.$$

• **Regularization:**

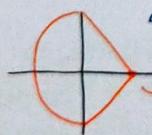


set a problem where  $\min \sum$

$$\text{st } \|\theta\|_1 = 1$$

con esta restricción, lo que le sumamos a una coordenada se la quitamos a la otra.

Si una sube mucho  $\rightarrow$  muchas ceros.



changes (constraints) reduce the angles.

many problems have this as solution.

problem is minimize  $\sum e(y_i - \theta^T x_i) + \lambda \|\theta\|_1$

$$h(x, \theta) = \theta^T x = \sum_{i=1}^J \theta_i x_i \quad \begin{array}{l} \text{d dimensions} \\ \text{Features} \end{array}$$

$$+ \lambda \|\theta\|_2$$

↳ add this for a lasso.

many zeros  $\rightarrow$  we are selecting features with impact on the output.  
TEST-WORTHY

• **Data augmentation:** (remember it also works to avoid overfitting) - add noise:  $\tilde{x} = x + \eta$ ;  $\tilde{x} \sim \mathcal{N}$  with  $\eta \sim \mathcal{N}$  (normal dist.)

• **Dropout:** set to 0 Features at random (changes pixels to 0 in an image)  $\tilde{x} = x \cdot 0$  (Feature wise)

= ^ midterm exam contents =

= Decision Theory =

- An **agent** is someone or something choosing options (**actions**) based on context evaluation and consequences.

- state-consequence matrix:

Actions / State of the world	OK	Traffic-Jam
A1. Take the bus	arrive early	very late
A2. Walk	arrive late	arrive late

state of the world might not be certain.  
it general it won't be

we expect **preferences** to give us an ordering based on a numerical value (**utility**):  $A \geq B, A \rightarrow B$

A) Garlic soup  $ABC \leftarrow A \Rightarrow \{AB \leftarrow A\}$

B) Carrot soup  $ABC \leftarrow A \Rightarrow \{AC \leftarrow A\}$

C) Pumpkin cream  $BC \leftarrow ?$  we have to say something about it, e.g.  $BC \leftarrow BC$  "tie/ide"

④ Irrationalism: not respecting transitivity:  $ABC \leftarrow A \Rightarrow AB \leftarrow B$  or  $\begin{array}{c} A \\ \swarrow \quad \searrow \\ B \rightarrow C \end{array}$   
we won't consider this: it should not happen.

> completeness vs incompleteness: completeness shows us if there's irrationalism or not

$\begin{array}{c} A \\ \swarrow \quad \searrow \\ B \rightarrow C \end{array} \rightarrow$  could be rational or not.

- **utility**:  $u(A) \geq u(B) \Leftrightarrow A \geq B$  (it is a function that respects the order)

Let's take into account the state of the world's:  $A / S_{0-W} \mid P(S_i) \quad P(S_j)$

- **lottery**:  $L_1 = \langle (p_1, x_1) \rangle$  where  $p_j$  are probabilities and  $x_j$  are prizes.

$L_2 = \langle (p_2, x_2) \rangle$

compounded lotteries:  $\langle (q_1, x_1), \dots, (q_n, x_n) \rangle$  where  $\sum_i q_i = 1$

$$\langle (50\%, \$10), (50\%, \$100) \rangle \xrightarrow{\text{assessment}} \text{expected value (prize)}: \sum_i x_i q_i = E(L)$$

$L_1 \geq L_2 \Leftrightarrow E(L_1) \geq E(L_2)$

We want to assess  $A_1 > A_2$ :  $A / S_{0-W} \mid P(S_1) \quad P(S_2)$  { probabilities

$A_1$	+1	-2
$A_2$	-1	-1

{ prizes

$\rightarrow A_1 \Rightarrow L_1 = \langle (P(S_1), 1), (P(S_2), -2) \rangle$

$\rightarrow A_2 \Rightarrow L_2 = \langle (P(S_1), -1), (P(S_2), -1) \rangle$

Lottery

$\Rightarrow E(U(A_j)) = \sum_{i=1}^n u(c_{ij}) P(S_i)$  : expected value of the utility of action  $A_j$

apply to conf. matrix

predicts state of the world (let's suppose if it's 1 or N)

		state of the world	
		P	N
not confusion matrix	Classifier chooses P	P	FP
		A	B
" "	" "	FN	TN
		C	D

we could add utilities here

preferences (default):  $A \sim D, B \sim C, T \geq F: A=D=1, B=C=0$

adding our understanding of the SoW:  $A_1: P(y=P|x) \quad P(y=N|x)$   
 $A_2: P(N) \quad \downarrow \quad \downarrow$   
 $A_2: \quad N \quad \text{prob of our classifier to say} \quad \text{prob of it saying Neg.}$   
 $\quad \quad \quad \text{Pos given data}$

using lotteries:  $L_1 = \langle P(y=+|x), 1 \rangle = EU(+)=P(y=+|x)$   
 $L_2 = \langle P(y=-|x), 1 \rangle = EU(-)=P(y=-|x)$   $\Rightarrow$  IF  $EU(+)\geq EU(-) \Rightarrow$  classifier uses +  $\Rightarrow$

• Assumption 1:  $\exists$  set  $\Omega$  with  $X \times Y \subseteq \Omega$  a  $\sigma$ -algebra  $\mathcal{F} \subseteq 2^\Omega$   $(x, y) \sim P_{\text{BB}}$  and a probabilistic measure  $P: \mathcal{F} \rightarrow [0, 1]$   
 $(M)$  That is  $\exists$  probabilistic space  $(\Omega, \mathcal{F}, P)$  such that  $D \subseteq \Omega$ .

Let  $D = \{(x_i, y_i)\} \quad i=1, \dots, N$

↳ can compress data into a model (basis of physics)

$P(y^*|x^*, D) \xrightarrow[\text{ML}]{\text{discriminative}}$  Model  $P(y|D)$  directly  $\xrightarrow[\text{A2}]{\text{A2}}$   $P(y^*|x^*, D) = \int_M p(y^*|x^*, M) p(M|x^*, D) dM = \int_M p(y^*|x^*, D, M) p(M|x^*, D) dM \Rightarrow \text{②}$   
 $\downarrow$   
 $\xrightarrow[\text{ML}]{\text{generative}} M$  non-parametric models: KNN, Kernels, Gaussian Processes.  
 $\hookrightarrow$  the more samples, the more parameters they need.  
 $P(x, y)$ : Probabilistic  $P(y^*|x^*, M) \quad p(M|D)$

• Assumption 2:  $P(y^*|x^*, M, D) = P(y^*|x^*, n)$  compression (we can get rid of data when we have our model)

• Assumption 3:  $P(M|x^*, D) = P(M|D)$

④  $\Rightarrow \int_M p(y^*|x^*, M) p(M|D) dM \xrightarrow[\text{ML}]{\text{classical}}$   $p(y^*|x^*, D) = p(y^*|x^*, n)$   
 $\downarrow$  Ensemble learning  $\downarrow$  AH, A5  
 $\downarrow$  SS substituting  $p(M|D)$  by  $\delta(M - M^*)$   $\xrightarrow[\text{Test}]{\text{Test}}$   
 $\sum_{i=1}^k p(y^*|x^*, M_i) w_i$   
 $\xrightarrow{P(M_i|D)}$

• Assumption 4: winner takes it all  $p(M|D) = \delta(M - M^*)$

• Assumption 5: "Maximum A Posteriori" model.  $M^* = \underset{\text{training}}{\text{argmax}} P(M|D) = \underset{n}{\text{argmax}} \frac{P(D|M) P(M)}{P(D)} = \delta(M - M^*) = \begin{cases} 1 & \text{if } M = M^* \\ 0 & \text{otherwise} \end{cases}$

$$= \prod_{i=1}^N p((x_i, y_i)|M) \cdot P(M) = \underset{M}{\text{argmax}} \prod_{i=1}^N p(y_i|x_i, M) p(x_i|M) p(M) = \underset{M}{\text{argmax}} P(M) \prod_{i=1}^N p(y_i|x_i, M) = \text{②}$$

• Assumption 7: Exponential Family:  $p(y_i|x_i, M) = e^{-\lambda_2 F_2(y_i, x_i, M)} \rightarrow P(M) = e^{-\lambda_2 F_2(M)}$

$$\text{④} = \underset{M}{\text{argmax}} e^{-\lambda_2 F_2(M)} \prod_{i=1}^N e^{-\lambda_1 F_1(y_i, x_i, M)} = \underset{M}{\text{argmax}} \underbrace{-\lambda_2 F_2(M)}_{\|w_M\|} - \underbrace{\sum_{i=1}^N \lambda_1 F_1(y_i, x_i, M)}_{\text{loss func.}} = \underset{M}{\text{argmin}} \lambda_1 \sum_{i=1}^N \mathcal{L}(y_i, x_i, M) + (\lambda_2 F_2(M)) \xrightarrow{\text{regularization}}$$

Empirical risk minimization extended

Let's say now:  $M^* = \underset{M}{\operatorname{argmax}} P(M|D) = \underset{M}{\operatorname{argmax}} \prod_i P(M|_{(x_i, y_i)}) = \text{(missing)}$

= PdL 10: Linear Models =

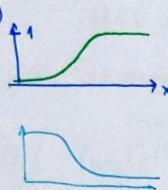
• Logistic regression:

We start from  $M^* = \underset{M}{\operatorname{argmax}} P(M|D)$  as we want to maximize  $\Rightarrow M^* = \underset{M}{\operatorname{argmax}} \prod_i P(M|_{D(i)})$

$$= \underset{M}{\operatorname{argmax}} \prod_i P(D(i)|M) * P(M) = \underset{M}{\operatorname{argmax}} \prod_i P(y_i|_{x_i, M}) P(x_i|M) P(M) \quad \text{④}$$

"we can drop this model's prob." → why?

I want to model  $P(y_i=1|x_i, M)$  as

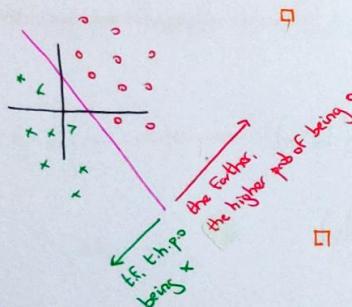


$$1 - P(y_i=1|x_i, M) = P(y_i=0|x_i, M)$$

our data is {0, 1} → remember "linear classifier"

↳ we will use a regression in order to be able to get the middle part.

The classifier we are building will be



what about these samples?  
very high prob in zones we don't know anything about: ASSUMPTION

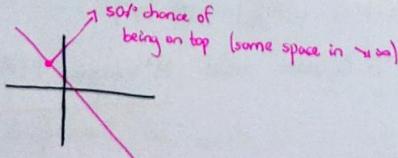
It makes sense: our model is too simple

$$\textcircled{4} = \underset{M}{\operatorname{argmax}} \prod_i P(y_i=1|x_i, M)^{y_i} P(y_i=0|x_i, M)^{1-y_i} \quad \text{up to this point: it is true for any binary classifier.}$$

$$= \underset{M}{\operatorname{argmax}} \prod_i P(y_i=1|x_i, M)^{y_i} (1 - P(y_i=1|x_i, M))^{1-y_i} = \underset{M}{\operatorname{argmax}} \sum_i y_i \log P(y_i=1|x_i, M) + (1-y_i) \log (1 - P(y_i=1|x_i, M))$$

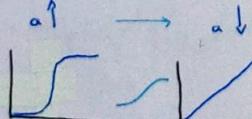
↓  
skipping a few steps

$$\Rightarrow P(y_i=1|x_i, M) = \frac{1}{1 + e^{(a^T x_i + b)}}$$



-  $b$  moves the boundaries "left to right"

-  $a$  defines how sharp the curve is

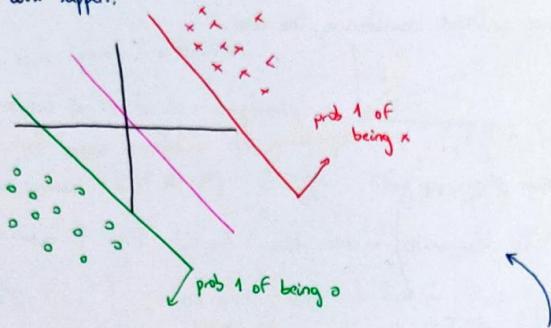


it defines how confident we are (less intermediate zone)

- we will be using the logit Function for prediction:  $P(x_i|y) = P(x_i) = \frac{1}{1 + e^{(a^T x_i + b)}}$  → decided thru  $\hat{y} = \underset{y}{\operatorname{argmax}} P(y|x_i, D)$

optimization: gradient descent with its loss (in notebook)

IF the sets are separable, this will happen:

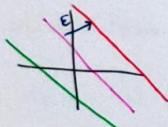


as we are maximizing (the number of predictions with high confidence). Of course, this might be zoomed in/out.

- First approach to modeling the separating hyperplane.

$$\begin{aligned} & \text{minimize } 1 \\ & \text{subject to } a^T x + b > 0 \text{ if } y=1 \\ & \quad a^T x + b < 0 \text{ if } y=0 \end{aligned}$$

Adding a threshold  $\varepsilon$  to



$$\text{we may get } \varepsilon(a^T x + b) > 0 \Rightarrow \varepsilon(a^T x + b) \geq \varepsilon \Rightarrow a^T x + b \geq 1$$

Therefore:

$$\text{minimize } 1$$

$$\begin{aligned} & \text{subject to } y_i(a^T x_i + b) \geq 1 \text{ with } y_i \in \{-1, 1\} \\ & \quad \underbrace{f(x_i)}_{\text{subject to.}} \\ & \text{we want to maximize } d(p \rightarrow \pi) = \frac{a^T x + b}{\|a\|} \end{aligned}$$

[...]

$$\begin{aligned} & \text{Soft-margin SVM: minimize } \|a\|_2 + C \sum \xi_i \quad \text{error} \\ & \text{subject to } a^T x_i + b \geq 1 - \xi_i \text{ if } y_i = 'x' \quad \text{labels.} \\ & \quad a^T x_i + b \leq -1 + \xi_i \text{ if } y_i = 'o' \end{aligned}$$

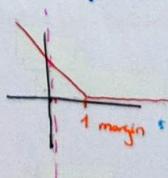
So now our problem is

$$\begin{aligned} & \text{minimize } \|a\|_2 + C \sum \xi_i \\ & \text{subject to } y_i(a^T x_i + b) \geq 1 - \xi_i \Rightarrow \xi_i \geq 1 - y_i(a^T x_i + b) \quad \xi_i \geq 0 \end{aligned}$$

$$\xi_i \geq \max(0, 1 - y_i(a^T x_i + b))$$

Turns out then that the solution should satisfy the equality

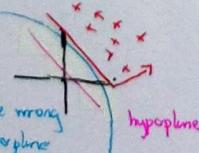
$$\text{minimize } \|a\|_2 + C \sum \max(0, 1 - y_i(a^T x_i + b)) \rightarrow \text{Hinge Loss}$$



unconstrained problem solved by constrained methods.

\* Only points that irritate me:

aka those points that are wrong or those close to hyperplane



separable problem  $\Rightarrow$  can  $C$  be 0? NO.

Because we are maximizing the distance without considering the error.

rewriting:

$$\text{minimize } \|\omega\|_2 + C \left( \sum_i \xi_i + \sum_j \xi_j \right) \rightarrow \text{separating the validation errors for neg and pos. value.}$$

subject to

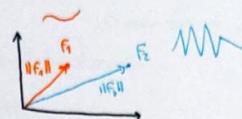
$$\begin{cases} (\omega^T x_i + b) \geq 1 - \xi_i & \text{if } y_i = +1 \\ (\omega^T x_i + b) \leq -1 + \xi_j & \text{if } y_j = -1 \end{cases}$$

$\|\omega\|_2 + C \left( \sum_i \xi_i + \sum_j \xi_j \right)$  to give more 'anger' to one of the errors (FP/FN)

## = Kernel Methods = Pill 11 (16)

We will consider a space where each point is a function.

We will assume that the norm is related to the function complexity.



Functional Analysis  $\rightarrow$  function spaces  $(\mathbb{R}, (f: \mathbb{R} \rightarrow \mathbb{R}), \mathcal{C}^1(\mathbb{D}))$ . They are vector spaces.

We are looking to bound the norm of  $\|f\|$ , for all  $f$  with similar behaviour.

• Hilbert Space:  $\rightarrow H = (V, \langle \cdot, \cdot \rangle)$  also: norm concept is linked to Hilbert Spaces.  $\|v\| = \sqrt{\langle v, v \rangle}$

- A Hilbert Space  $H$  is a reproducing Kernel Hilbert space (RKHS)

Also, Cauchy-Schwarz:  $|\langle v, v \rangle| = \|v\| \|v\|$

IF  $H$  is a RKHS:  $\mathcal{F}_x[f] = f(x) = \langle K_x, f \rangle$

notion of complexity.

$$\alpha = \{f(x_i) - y_i\}^2$$

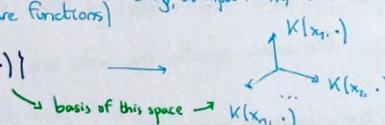
Now we may introduce this notion into regularization:  $f^* = \underset{f \in H}{\operatorname{argmin}} \frac{1}{n} \sum_i \mathcal{L}(f(x_i), y_i) + \lambda \|f\|_H^2$

complexity penalty

• Representer's Theorem: the solution of former problem has the following form:  $f^*(\cdot) = \sum_i \alpha_i K(x_i, \cdot)$

(So far, all we know about Kernels is that they are functions)  $\xrightarrow{\text{2 arg, as input. Represents similarity.}}$

proof:  $S = \operatorname{span} \{K(x_1, \cdot), K(x_2, \cdot), \dots, K(x_n, \cdot)\}$



$$\textcircled{1} \quad F \perp S \Leftrightarrow f(x_i) = 0 \quad i=1, \dots, N$$

proof:  $f(x) = \langle K(x, \cdot), f \rangle = \phi \quad i=1, \dots, N \rightarrow$  if it is orthogonal to the whole space, all its evaluations must be 0.

$$\textcircled{2} \quad f(x_i) = f_S(x_i) \quad \uparrow \quad \text{a function can be decomposed into a hyperplane and its perpendicular.}$$

proof:  $f(x_i) = f_{\perp}(x_i) + f_S(x_i) \rightarrow$  minimizing  $f(x_i) \equiv \min_{\perp} f_S(x_i)$   $\textcircled{3}$

$$\textcircled{3} \quad \|f\|_H \geq \|f_S\|_H$$

proof:  $\|f\|_H \geq \|f_{\perp}\| + \|f_S\| + \langle f_{\perp}, f_S \rangle$

$$\textcircled{4} \rightarrow f^* = \underset{f \in H}{\operatorname{argmin}} \frac{1}{n} \sum_i \mathcal{L}(f_S(x_i), y_i) + \lambda \|f_S\|_H^2$$

$$f(z) = \alpha_1 K(x_1, z) + \alpha_2 K(x_2, z) + \dots \quad \rightarrow \quad f_S^*(z) = \sum_{i=1}^N \alpha_i K(x_i, z)$$

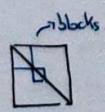
$\rightarrow$  we passed from smth very abstract to smth in function of  $x_i$ .

{...}

- RKHS needs positive semidefinite Kernel functions.

What are Kernels?

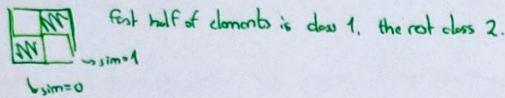
• Gram Matrix:



→ evaluate Kernel on whole dataset  $G = (K(x_i, x_j))_{i,j}$

$G \in \mathbb{R}^{N \times N}$

- Kernel (similarity matrix):  $K(x_i, x_j) = 1$  if  $y_i = y_j$



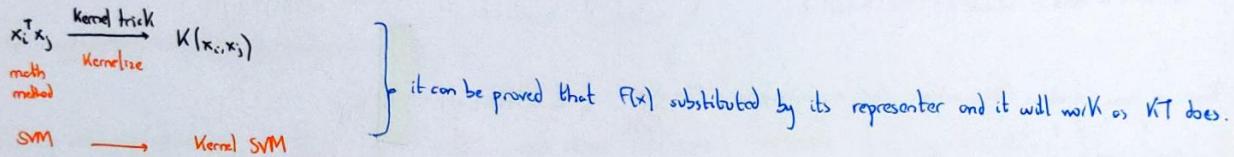
It allows us to classify any kind of data in any form (not only vectors)

Question:  $K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{\sigma^2}}$  → you can input anything as long as it is inside a distance (graph, e.g.)

- The idea comes from  $y^T y$  and from the goal to form something close to the similarity matrix.

- Kernel are linear ( $K = \alpha K$ ;  $K = K_1 + K_2$ ) and also,  $K = K_1 \cdot K_2$

(Kernel trick)



= Decision Trees = (Pill 12: Ensembles)

In decision trees, the top nodes encode the most information.

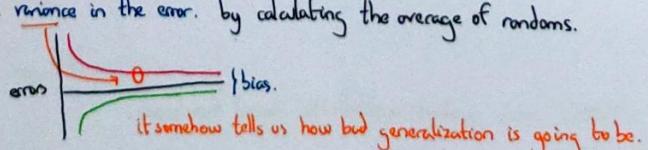
↳ the deeper you go, the more 'specific' the nodes are.

- What's the best decision tree? Many ways to answer this:
  - $\nwarrow$  a value of a feature is tested
  - In. Nodes = [ ] → contains feature + threshold
  - $N \leftarrow$  retrieve first (pop)
  - $\text{data} \leftarrow N$
  - $c \leftarrow$  find splitting point.
  - After taking a step, the next one to be chosen will be the one that maximizes  $\text{Gain} = \text{Parent} - \underbrace{\text{Avg}}_{\text{error}}(\text{children}) - \underbrace{\text{Avg}}_{\text{errors}}(\text{children})$
  - Add 2 more nodes to the list.
  - Repeat.

Out. Decision Tree (Classification)

- Avoiding overfitting: ensemble learning → reduces the variance in the error, by calculating the average of randoms.

↓  
train a set of models  
aggregate



notion: if one classifier is unlucky at classifying a sample → maybe the rest is luckier.

it makes little sense to have 1 model for huge datasets → submodels.

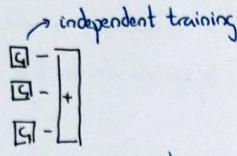
divide + conquer (have many classifiers)

it makes no sense either to have one classifier in cases of data fusion (heterogeneous sources)

- diversity: "as how you choose 2nd opinions from people with different backgrounds"

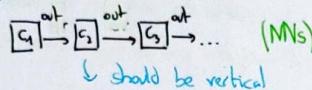
Ensemble through: 1) Bagging

(Bootstrapping  
aggregation)



2) Boosting relies on waste  $G_1 \xrightarrow{\text{err}} G_2 \xrightarrow{\text{err}} G_3 \xrightarrow{\text{err}} \dots$  they all depend on each other. link through residual statistical resampling

3) Stacking output  $\rightarrow$  input



gradient boosting  
adaptive boosting

① ② aggregation done through majority voting

③ Stage wise approximation:

$$H(x_i) = \sum_{t=1}^T \alpha_t h_t(x_i) \quad \text{with classic LSE} \rightarrow \text{minimize } \frac{1}{N} \sum_{i=1}^N (y_i - H(x_i))^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \sum_{t=1}^T \alpha_t h_t(x_i))^2$$

linear comb. of models

we will minimize model by model, aka, at time  $T$ ; minimize  $\frac{1}{N} \sum_{i=1}^N (y_i - \sum_{t=1}^{T-1} \alpha_t h_t(x_i) - \alpha_T h_T(x_i))^2$

already fixed

The new added model  $T$  is trying to decrease the residue of the rest.

↳ Specialized in correcting the rest of models.

- reductionist frameworks for multi-class problems

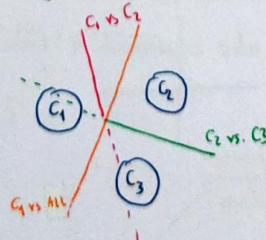
	H MCClass	scalation ability
SVM	N	Naturally?
NN	Y	
Probabilistic	X	'NO' Kinda last layer changes dim.
Dec. Tree	Y	Bad. Linear (you need to model one prob per class)
K-MN	Y	Yes
Log. Regression	N	Yes

↳ we only care about the labels of close samples to the query samples.

'1 against all': divide-n-conquer in multiclass

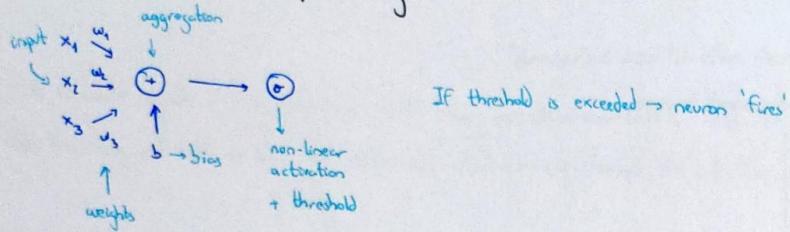


1 vs 1 but aggregated.



'1 against 1': consider all pairs of classes and then aggregate.

= PdL 13b: A story of deep learning =



- sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$



(Fully connected layers of neurons)

- To define a NN:
  - 1) Architecture (Hypothesis space)
  - 2) Loss function.
  - 3) Optimization algorithm

> Number problem: 10 class problem

Remember from DL: FCN are a subtype of FCNN.  
c = conv.  
= dense in Keras.

- architecture definition:
  - 1) Sequential: limited to stacking layers.
  - 2) Graph RESNet (the input of a layer is not always the output of the latter)

We want a probability of the input being a specific number  $\Rightarrow$  we need normalization: SOFTMAX  $\rightarrow$  at the end.

$$p_i \in [0, 1] \quad \sum p_i = 1$$

- it is called a softmax layer because it emphasizes big values and deemphasizes small values.

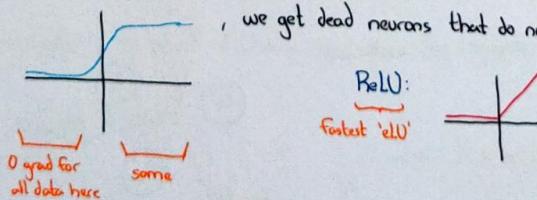
It is not needed for the NN to work, but it is handy for the exponentials in the loss function.  $\rightarrow$  easier computation.

- cross-entropy loss function.

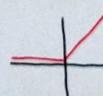
We expect more layers to lead to better results. Spoiler: it doesn't. At some cases it gets even worse.

overfitting  
vanishing gradients.

- Vanishing gradients: due to sigmoid



ReLU:  
fastest 'elu'



many elus: ReLU, Leaky ReLU, ...  
recomendacion del chef

Can we add now  $10^{300}$  layers? No.

- Internal covariate shift: The distribution moves through each layer:

We need Batch Normalization (normalize data according to stats in the batch)  $\rightarrow 0$  means centering the data

'approx'  $\rightarrow$  we're only using a batch.

Will it work to add 30 sigmoid layers with batch normalization? YES. Also, 100% accuracy.

No training and test so far  $\xrightarrow{\text{add split}}$  shit results. And we only changed the data distribution.  $\rightarrow$  overfitting / few samples  $\rightarrow$  data augmentation \*  
high complexity

- change the split ratio
- transformation (img, flipping)
- random noise  $\rightarrow$  dropout layers. some values are randomly set to 0.

$\hookrightarrow$  only in images is good to drop some pixels due to high correlation between neighbours.

- Layers in deep learning is all about Knowledge representation. The last layer is a linear classifier, the rest of layers transform for the last layer.

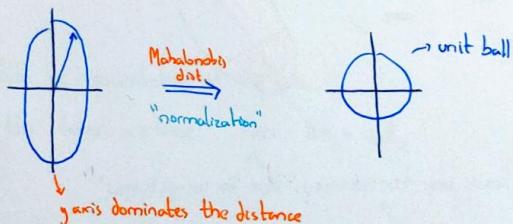
$\downarrow$  can be changed for a real classifier. (random forest)

increase in accuracy  $\rightarrow$  the transformation was good.

(...)  $\rightarrow$  Autoencoders, whole pill.

= Pill 12 Gaussian discrimination and processes =

We can compensate axes with large variance.



We will suppose we have quadratic forms  $\Rightarrow$  all eigenvalues  $> 0$

PCA: The eigenvalues will appear in order of maximum variance.

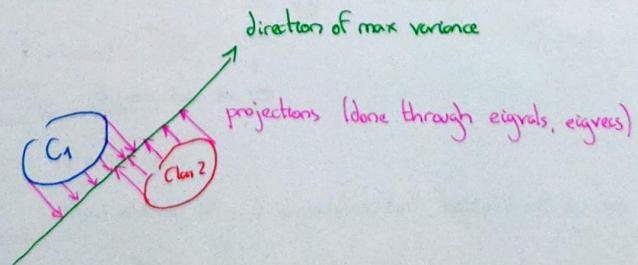
assumption

$\Rightarrow$  let's represent all the dataset as a normal distribution and get the direction with the most variation

we want to use a linear transformation to select  $d'$  dimensions onto which project our  $d$ -dimensional dataset

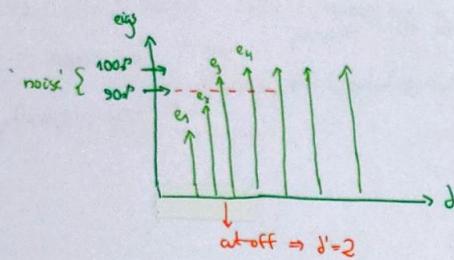
$\downarrow$  with  $d' \ll d$

$d'$  is project or depending on 'energy'  
(Keep a % of variance of data)



- how to fix  $\delta'$ ? As a hyperparameter, in relation to the ratio dim/noise.

Imagine we want to recover 90% of the energy:



note: remember to center by subtracting the mean

$$\sum \underbrace{\mathbf{B}}_{\substack{\text{signals} \\ \text{eigenvectors}}} \mathbf{B}^T \xrightarrow{\text{centered}}$$

$\downarrow$  selection axes of max. eigenvectors  $\Rightarrow$

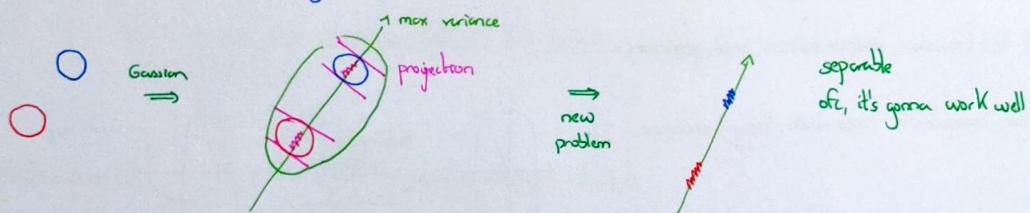
$$\mathbf{B}' \in \mathbb{R}^{d' \times d'}$$

$$\left. \begin{array}{l} \mathbf{x} \in \mathbb{R}^{35} \\ \mathbf{B}' \in \mathbb{R}^{30 \times 5} \\ \mathbf{x}' \in \mathbb{R}^{d'} \end{array} \right\} \boxed{\mathbf{x}' \mathbf{B}' = \mathbf{x}'}$$

it results that working with the selection is the same as working with the whole dataset.

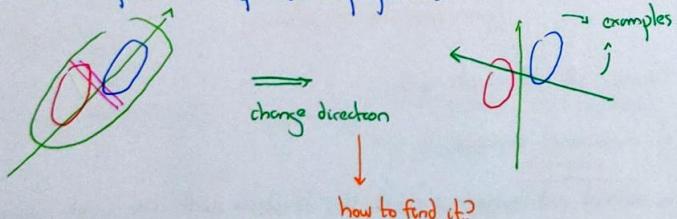
note: the cut-off for  $d' = 5$  can be seen with the graph in cell 37

Problem 1:



A problem can turn from sep. to non-sep. due to projections.

Problem 2:



new problem: Linear discriminating analysis.

LDA: our new  $\mathbf{B}$  is gonna be:

$$\mathbf{B} = \frac{\mathbf{S}_b}{\mathbf{S}_w} \rightarrow \text{between classes scatter matrix: } \underbrace{\mathbf{S}_b}_{\text{difference between the mean of both classes}}$$

$$\mathbf{S}_w \rightarrow \text{within " " " " : } \underbrace{\mathbf{S}_w}_{\text{distance from class instances to class center}}$$

$$\Downarrow$$

$$\max_{\min} \rightarrow \underset{\min}{\text{maximize}} \mathbf{B}$$

we want to max. dist. line from center of both classes

minimum  $\rightarrow$  compression

LDA  $\rightarrow$  Line in problem 2 of this cbr.

note: we get the direction, that maximizes  $\mathbf{B}$ . All parallels work.

I would like to thank you for your help

The introduction of the use of honey had an extreme influence

One thing you can do is to make sure that you are not the only one who is writing, so if there are two or three people, then you can all write, because you can all have different ideas and different ways of writing.

Wij vinden de hand <sup>van</sup> de vader in de <sup>hand</sup> van de zoon de <sup>hand</sup> van de moeder.

1000

My last few words are these to those who have been  
"Baptized into Christ Jesus our Saviour, you will be born again."

May 1990 1990

One more to consider, as in the top section, *and the*



• We can use **adjectives**, **adverbs**, or **adverbs of emphasis** (when there is **emphasis**), as there are **adverbs** of **emphasis**? Can you give an example with an **adverb of emphasis**?

www.pearson.com

Each stage will have a ~~mark~~, all marks in sequence to the aim.

• We will use this fact to be enough to prove that  $\mathcal{C}_1$  and  $\mathcal{C}_2$  are isomorphic.

卷之三

www.english-test.net

Impressionism, in



"Zembla"

22 *Empidonax hammondi* (Vig.)

1770

It might be especially striking how a few changes to the school (with its goals) → simple school again → simple by the ... in consequence.

How many clusters do we need to compute?

We may get less or more clusters than data sources (e.g.)

There are automatic ways to determine the number of clusters (we will see that)

We can also find correct solutions that do not fit: • we see '2' clusters and we know there are 3

the solution may split the 'wrong' cluster.

• divide horizontally when it's vertical, or viceversa.

- This happens because K-means assumes all directions weigh the same (Gaussian)

$r_{ik} = \begin{cases} 1 \\ 0 \end{cases}$  is too black or white  $\Rightarrow$  make all elements be involved in responsibility: all elements will eventually pull if there are several together.

↓ change it to something that relates distance and similarity.

**SOFT K-MEANS**  $e^{-\gamma d(x_i, c_k)} \rightarrow$  actually a Gaussian  
weight  $\rightarrow$  it represents the 'zoom effect' previously mentioned.

- Imagine we want way too many clusters. Depending on the zoom, centroids will eventually collapse.

↑ use this as a toy and discover number of clusters (high computational needs)

• Mixture of Gaussians: learn  $\beta = \frac{1}{\sigma^2}$

$$\text{minimize} \sum_i \sum_k e^{-\beta \|x_i - \mu_k\|^2} \|x_i - \mu_k\|^2$$

$$\text{minimize} \sum_i e^{-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k)}$$

covariance matrix

we can use Mahalanobis distance (not all directions weigh the same)

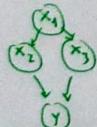
• Maximum Likelihood

Letting each Gaussian have different covariance matrices allows us to have concentric clusters.



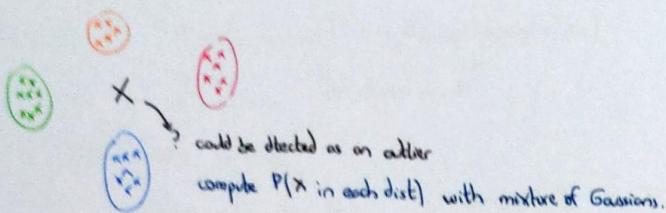
= PdII 14b: Unsupervised learning + Clustering =

• PGM:  $\underbrace{P(y, X)}_{\text{joint}} \rightarrow P(y|X)$



$$P(y, x_1, x_2, x_3, x_4) = P(y|x_1, x_2, x_3, x_4) = \underbrace{P(y|x_1, x_3)}_{\substack{\text{most complex term} \\ \text{no edge} \\ (\text{independency})}} \underbrace{P(x_1|x_2, x_3)}_{P(x_1)} \underbrace{P(x_2|x_3)}_{P(x_2)} \underbrace{P(x_3)}_{P(x_3)}$$

Assume all data comes from distributions:



• **Clustering (K-means):** novel  $\Leftrightarrow$  outliers

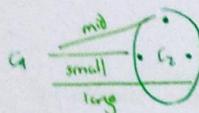
in terms of 'unusual'

Control countries may make relevant data less relevant if they are too many: 1) one may be an outlier

2) they could be in a cluster of only control countries.

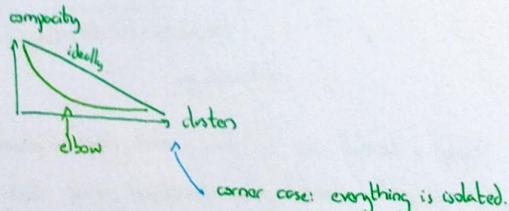
Solutions: - delete from dataset and reinsert it on top. (after clustering)

It might be interesting to locate where in the cluster an element is (checking the distance to other clusters' centroids)



- K-means depends on the K number of clusters. How to get K? **Elbow/Vknee technique** (notion of compactness is related to the average distance to the centroids)

corner case: all entries in a point  $\rightarrow$  1 cluster



It only gives an idea of a number of clusters that may be good.

remember: don't use K-means (soft K-means, etc.)

- **Hierarchical clustering:** dividir maximizando distancia entre clusters más cercano      # inicial de clusters: 1      # entries      } costful

- **DBSCAN** (density-based S.C.A.N.)  $\rightarrow$  gets funny shaped clusters.

It uses local measures, rather than global ones.

It can detect points in high density areas, low density and outliers.

CORE      boundaries

Select a point, consider a  $\epsilon$ -neighbour, share label with all  $\epsilon$ -neighbour. Select a point from the neighbourhood and repeat.

-  $\epsilon$  and minpoints. (...)

- Spectral clustering:  
eigenvectors of a similarity matrix

## = Pill 2: Dirty Data =

- Split columns  $df[\text{cols}] = df[\text{col}].str.split(sep, expand=True)$ 
  - $\text{list of rows}$
  - $\text{sep: column names}$
  - $\text{str: column to split}$
- Drop useless columns  $df.drop$
- Replace all missing values for NaN (\*will be extended)  $df = df.replace(..., np.Nan)$ 
  - $\downarrow$  opt: restrict by columns
- Delete waste from numerical columns:  $df.replace(regex, "", regex=True)$
- Handling missing data: *inputting*
  - ↳ add pair wise deletion to notes: if we are analysing a subset, only delete samples from that subset.
  - ↳ average + dummy variable indicating if value was originally missing.
- More info  $\Rightarrow$  more dimensionality  $\not\Rightarrow$  Better results
  - Especially in ML

- Approaches to outliers:
  - Statistical description: assume data follows a distribution. Those samples that do not  $\rightarrow$  outliers.
  - Geometric: drawing in notes
  - Distance to neighbours: KNN may be used.
  - Percentile funct:  $\hat{x} = \frac{x - \text{prc}(x, \sigma)}{\text{prc}(x, 100 - \sigma) - \text{prc}(x, \sigma)}$ 
    - percentile we consider as outliers ( $10\sigma^2$ )
    - $\downarrow$  np.percentile

- Questions:
- 1) In general it is good idea to drop samples/missing data (only not when dataset is small)
  - 2) You can capture non-linear behaviour with linear regressors using dataset extension.
  - 3) Only use standardization with few outliers.
  - 4) Normalization is useful to compare different ranges.
  - 5) Categorical variables  $\rightarrow$  One HtOT (Don't do green=1, blue=2...)
  - 6) Hashing trick does not maintain distance 
  - 7) Hashing is  $\sim$  uniform.
  - 8) In distance based classification algorithms, it is a bad idea to add random features.

## = Pill 3: Performance =

- Cross-validation: all points will be used for testing. (LOO with  $n$  it).
- K-fold cross validation as  $L^{\text{Many}} O \approx LKO$  with  $K$  it.
- Unbalanced datasets with imbalance between positives and negatives.

• Confusion Matrix:

		Truth		
		Positive	Negative	
Prediction	Positive	TP	FP	→ precision: $\frac{TP}{TP+FP}$
	Negative	FN	TN	→ negative predicted value: $\frac{TN}{TN+FN}$
		↓ sensitivity/recall	↓ specificity	
accuracy:		$\frac{TP+TN}{TP+TN+FP+FN}$	$\frac{TP}{TP+FN}$	$\frac{TN}{TN+FP}$
				estimativa: $\frac{T \text{ de la fila / columna}}{\text{suma de la fila / columna}}$

Si churn = +.

Correctly predicted churned ppl? Recall

From the ppl that predictedly churn, how many actually do? Precision.

$$-F1 \text{ Score} = 2 \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2TP}{TP+FN+TP+FP} = \frac{TP}{TP + \frac{1}{2}(FP+FN)}$$

• Other alternatives: undersample majority class

oversample minority → interpolator.

class weights.

change performance metric.

split majority in subclasses → train a classifier with each. → aggregation.

= Pill 4: Supervised ML =

• Hypothesis space/model class: Family of mathematical models to consider. → Linear:  $f((w_0, w_1) = w_1 x + w_0)$

• Problem model: optimization, minimization of error function

• Learning algorithm: optimization method or algorithm that fits the model to the data.

↳ target: simplicity/speed/interpretability vs accuracy ↓ (trade-off between por la ultima) → Gradient descent.

Linear regression model:

$$\tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix} \Rightarrow \hat{y} = \tilde{x}^T w \rightarrow \text{minimize} \underbrace{\frac{1}{N} (y - \tilde{x}^T w)^T (y - \tilde{x}^T w)}_{\text{LSE}} \text{ with } \nabla_w \text{d} = -\frac{2}{N} \tilde{x} (y - \tilde{x}^T w)$$

= Pill 5: Feasibility =

Is learning possible if  $f$  is unknown? No, we can't know if  $f$  gives a certain value in unseen data

But yes, through probability → red/green balls table-out

might give more green

while there are more red in bin

⇒ All combinations are possible, but not equally probable. We hope the sample to be similar to the total.

•  $\epsilon = E_{in}$  'in sample error' → # times hypothesis getting it wrong.

•  $\mu = E_{out}$  'out of sample error' → expected

$\mu$  in Hoeffding

Learning is assuming  $E_{in}$  indicates  $E_{out}$  and finding the hypothesis with least  $E_{in}$ . ( $E_{out} \rightarrow 0$  through  $E_{in} \approx E_{out}$ )

- Hoeffding taking into account most case scenario among hypothesis is in terms of complexity of hyp. space.

We move to Probabilities as the target may not always be a function (sea, sky detection  $\leftarrow$  2 vals for blue)

- VP-dom: measure of complexity

Entries on questions: 1) Overfitting when  $E_{in} \downarrow E_{out} \uparrow$  ( $E_{in} \neq E_{out}$ )  
 2) Overfitting depends on complexity and samples (complexity  $\uparrow$  and samples  $\uparrow$  to avoid)

= Pill 6: Overfitting =

- We have to match the data complexity, not the underlying true model complexity (always start from data)
- Increasing complexity to reduce bias and make test and train error rates converge closer generally causes worse error rates.
 

$\uparrow$   
err. rates

fixed. complexity

bias seen in

$\downarrow$   
err. rates

num samples

given for a fixed complexity.

= Pill 7: Stochastic Subgradient Methods =

They satisfy that at  $x^*$  for all  $z$

$$f(z) \geq f(x) + \vec{S}(z-x) \rightarrow \text{hyperplane with normal } (S_i, -1) \text{ supports } f \text{ at } (x, f(x))$$

- $f$  convex and differentiable  $\Rightarrow$  its gradient is subgradient.
- The subgradient may exist even if  $f$  is not differentiable.
- $x^*$  is a minimizer of  $f$ , convex  $\Leftrightarrow f$  subdifferentiable at  $x^*$  and  $0 \in \partial f(x^*) \Rightarrow f(z) \geq f(x^*)$ .  
there is at least 1 subgradient      set of subgradients
- If  $f$  is differentiable at  $x^*$ , the equation is reduced to  $\nabla f(x^*) = 0$ .

- They are useful to solve problems like minimize  $f_0(x)$   
subject to  $f_i(x) \leq 0 \quad i=1, \dots, m$  with  $x^{(k+1)} = x^{(k)} - \alpha_k S^{(k)}$

$$\text{where } S^{(k)} = \begin{cases} \partial f_0(x^{(k)}) & \text{if } f_i(x^{(k)}) \leq 0 \quad i=1, \dots, m \\ \partial f_i(x^{(k)}) & \text{if } f_i(x^{(k)}) > 0 \end{cases}$$

→ There is a version with noise.

$$\text{- momentum: } x^{(k+1)} = x^{(k)} + m^{(k+1)} \quad \text{with } m^{(k+1)} = \gamma m^{(k)} - \alpha_k S^{(k)}$$

$$\text{- autograd: } x^{(k+1)} = x^{(k)} - \alpha_k D^{-1} S^{(k)} \quad \text{with } D = \sqrt{\sum_{i=1}^m \text{diag}(S^{(i)})}$$

$$\text{also could be } D^{(k+1)} = D^{(k)} + \text{diag}(S^{(k)})^2 \rightarrow x^{(k+1)} = x^{(k)} - \alpha_k (D + \epsilon)^{-1} S^{(k)}$$

- adam: adaptive momentum (combination of both)  $\rightarrow$  I wouldn't even bother...

= Pill 1: First Steps =

- Machine Learning:
  - Improve SW performance based on previous experience.
  - Methods that detect patterns in data to predict future data or take decisions under uncertainty.

- Types of ML
  - Predictive / supervised: training <sup>with target</sup> labelled dataset  $\Rightarrow$  regression, classification
  - Descriptive / unsupervised: given data set find information about structure  $\Rightarrow$  density, clustering, dim reduction, manifolds.
  - Reinforcement
- Curse of dimensionality: too complex

- ④ Quiz questions :
- 1) Subgradient generaliza gradiente para no diferenciales.
  - 2) Subgradients don't always decrease.
  - 3) Allow approximations with ~~loss~~ noise.
  - 4) Subgradient = Gradient si convexa (+ differentiable?)
  - 5) Adagrad and momentum don't give better solutions, just faster convergency.
  - 6) Advantage of adagrad over momentum: solution adaptation coordinate-wise.

= Pill 8: Regularization - Combating overfitting =

- 1) Keeping overfitting at check: regularization, build on ensemble, validation set
- 2) Regularization  $\rightarrow$  controls complexity of model
- 3) Correct ways of adding regularization term to objective function:
  - control # coeff.  $\neq 0$
  - control size of coeff. values.
  - control max value coeff.
- 4) Regularization techniques: adding norm to loss function, adding noise to input samples, early stopping.

= Pill 10: Linear Models = • Least Squares regression:  $h(x) = a^T x + b$ ;  $\min_{a,b} \sum_i (y_i - h(x_i; a, b))^2$

• Logistic regression: linear classifier based on maximum likelihood to model classification boundary and certainty of the fit. Fits logistic function to data.

$$p(x_i | y_i) = \frac{1}{1 + e^{-(a^T x_i + b)}}$$

• Support Vector Machines (SVM): discriminative learning

↓  
Focuses on the decision boundary. works with labeled data (the more the better)  
discriminative vs. generative → prob. model of each class; unlabeled; boundary = where a model is more likely.

We assume a function model for the boundary: linear (classical)

Other discriminative linear models: SVM, perceptron, logistic classifier...

It finds the boundary with maximum distance/margin to both classes.

- Support vectors: points right in the margins (remember OPT lab) If they disappear, boundaries change.  
critical subset of points.

It implicitly models noise through boundaries. We expect it to be robust against tiny perturbations.

- Maximum margin classifier → unique solution in separable case.

- Distance to a hyperplane:  $d(x_i, \pi) = \frac{a^T x_i + b}{\|a\|_2}$

goal:	$\max \frac{2}{\ a\ _2}$	$\min \frac{2}{\ a\ _2}$
$\Leftrightarrow$	$a^T x_i + b \geq 1$	al revs
or	$a^T x_i + b \leq -1$	
s.t.	$y_i(a^T x_i + b) \geq 1$	s.t. ...

↓  
closest different classes points are  $2 \cdot \frac{1}{\|a\|_2}$  apart

• Soft-Margin SVM (non-separable): introduce a new set of variables that measure amount of violation in i-th constraint. ( $\xi_i$ )

We want  $\xi_i \geq 0$ .  $\rightarrow \min \left( \frac{2}{\|a\|_2} \right) + C \sum_i \xi_i$  → trade-off between margin and misclassification rate.

s.t.  $y_i(a^T x_i + b) \geq 1 - \xi_i$   
 $\xi_i \geq 0$

- Quiz:
- Logistic regression commonly used for regression  $\times$   $\rightarrow$  classification.
  - " " " probability assignment to each data point  $\checkmark$
  - Concept of margin related to distance (samples, boundary)  $\checkmark$
  - SVM maximizes margins  $\checkmark$
  - Non-separable SVM's C balances robustness and training error  $\checkmark$
  - " " " "  $= 0 \equiv$  Hard margin SVM  $\times$

## = Pill 11: Kernels =

- Linear evaluation functional over the Hilbert space of functions  $\mathcal{H}$   $\mathcal{F}_x: \mathcal{H} \rightarrow \mathbb{R}$   

$$f \mapsto \mathcal{F}_x[f] = f(x)$$
- RKHS: if evaluation functionals are bounded,  $\exists M / |\mathcal{F}_x[f]| = |f(x)| \leq M \|f\|_{\mathcal{H}} \quad \forall f \in \mathcal{H}$ .
  - By the Riesz representation theorem, for each  $x \in X$ ,  $\exists$  function  $K_x \in \mathcal{H}$  (called representer) that satisfies:
$$\mathcal{F}_x[f] = f(x) = \langle K_x, f \rangle = \langle K(x, \cdot), f(\cdot) \rangle$$

Examples of positive semidefinite Kernel function  $K$ :

- linear:  $K(x_i, x_j) = x_i^T x_j$
- polynomial:  $K(x_i, x_j) = (x_i^T x_j + 1)^d$
- gaussian:  $K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{\sigma^2}\right)$

- Gram matrix: shows the degree of shared information among training samples.
  - interpretation:
    - 1) bad Kernel: mostly diagonal  $\rightarrow$  most points are orthogonal to each other, no clusters/structure.
    - 2) good Kernel: structure + clusters.
    - 3) ideal:  $K = yy^T$

- Kernels: tips
  - use your fav. distance  $d(x_i, x_j) \rightarrow K(x_i, x_j) = e^{-d(x_i, x_j)}$
  - " " " non-linear transform:  $K(x_i, x_j) = \varphi(x_i)^T \varphi(x_j)$

- Kernel trick: regularized Least Squares but using a mapping  $\phi(x): \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$  instead of the original input space.
 
$$\Rightarrow \alpha^* = \underset{\alpha}{\operatorname{argmin}} \frac{1}{2} \sum_i (\gamma_i - \alpha_i \phi(x_i))^2 + \frac{\lambda}{2} \|\alpha\|^2$$

$$\Rightarrow \text{model is } \hat{y} = \hat{\alpha}^T \phi(x) = \frac{1}{2} \phi^T \phi \alpha$$

$\hookrightarrow$  same addition as Kernel method with  $K = \phi \phi^T$

- props:
- using Kernels can be seen as mapping data in an arbitrarily high dimensional space spanned by  $\phi(x)$
  - using Kernels allows a compact implicit description of the space generated by  $\phi(x)$  by means of the Gram matrix  $K(x, \cdot) = \phi(x) \phi(\cdot)^T$

- $\dim \{\bar{g}(x)\} = \infty$  for a gaussian Kernel.

The Kernel trick replaces any  $\langle x_i, x_j \rangle$  by a Kernel  $\langle x_i, x_j \rangle = K(x_i, x_j)$

Also, in RKHS,  $f(x) = \langle K(x, \cdot), f(\cdot) \rangle$  and we can replace  $f(x) = \sum_i \alpha_i K(x, x_i) \quad \forall x_i \in D_{\text{train}}$

Online Kernel perceptron:

perceptron: linear model that solves

$$\min \frac{1}{2} \|\omega\|^2$$

$$\text{s.t. } y_i(\omega^T x_i) \geq 0$$

algorithm. 1) select random sample pair  $(x_i, y_i)$

2) If  $y_i \sum_{k=1}^l \alpha_k \langle x_k, x_i \rangle < 0$ :  $\alpha_i \leftarrow \alpha_i + y_i$

$\Rightarrow$  model:  $f(x) = \sum_{k=1}^l \alpha_k \langle x_k, x \rangle$

$\Rightarrow$  generalized alg: 2) turned into  $y_i f(x_i) < 0$ :  $\alpha_i \leftarrow \alpha_i + y_i$  with model  $f(x)$

in a RKHS  $f(x) = K(x, \cdot)^T \alpha$ , thus, if  $y_i \sum_{j=1}^N K(x_j, x_i) \alpha_j < 0$ :  $\alpha_i \leftarrow \alpha_i + y_i$

Online Kernel SVM:

$$\min \frac{1}{2} \|\omega\|_2^2 + C \sum_i (1 - y_i \underbrace{\omega^T x_i}_{\alpha^T K(x_i, X)})$$

Quizz: • Kernels allow to use linear models in non-linear problems with success. ✓

• Kernel trick = replace inner product by a Kernel evaluation. ✓

• Represent Theorem: minimal solution of a regularized loss objective function in RKHS is a linear combination of Kernel evaluations in the dataset. ✓

•  $K(x, y) = x^T y =$  linear Kernel ✓

• Increasing the sigma parameter in a RBF Kernel reduces the non-linear capacity of modelling boundaries. ✓

• A good Kernel has a diagonal matrix ✗

= Pill 12: Ensembles =

• Decision Tree: easy to interpret.

- All data might be boxed.

- There can be more than one solution that feeds the data.

- Splitting criteria: 1) Misclassification error

2) Gini index

3) Cross-entropy / Information gain / Mutual information

} maximize node purity

- stop overfitting: early stopping or post-prune

- **Diversity:** it is needed because we need the errors to happen in different samples.

- how to get it:
  - 1) Different training sets.
  - 2) Different training parameters.
  - 3) Combining different architectures (SVM, decision trees)
  - 4) Training on different features random subspaces or projections.

- Examples of ensemble models:

1) **Bootstrapping aggregation:** resampling the training set with replacement.  $\rightarrow$  bagging  
one classifier per dataset copy + majority voting.

2) **Random forest:** randomization on the feature selected for building each tree in the ensemble.

3) **Incremental learning**  $\uparrow$  **boosting**

gradient boosting

our model is  $H_t(x) = \sum_{t=1}^T a_t h_t(x)$

time  $\downarrow$  models: once they are added they remain constant.

$\mathcal{L}_t(x) = \sum_{i=1}^N (y_i - H_t(x_i))^2$  LSE

Increment  $\Rightarrow$  select  $h_{t+1}$  and optimize  $a_{t+1}$

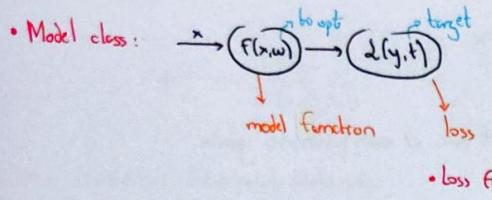
$H_{t+1}$

In  $\mathcal{L}_{t+1} = \sum_{i=1}^N \underbrace{(y_i - H_t(x_i) - h_{t+1}(x_i))}_{\epsilon = r_t(x_i)}^2$ ,  $r_t(x_i)$  represents error or residue of the model at time  $t$  on  $x_i$ .

adaboosting  $\left\{ \text{stage-wise approach also but with } \mathcal{L}(x, y) = \sum_i \epsilon_i y_i H(x_i) \right\}$

- Quiz:
- Decision trees use axis orthogonal hyperplanes ✓
  - Naturally handles multi-class: Decision Tree, KNN, NN, Prob. Classifier, SVM  $\downarrow \times \times \times$
  - Decision trees can fit any dataset ✓
  - Ensemble learning = field of ML studying stacking and aggregation of models. ✓
  - Diversity is key to ensembles. It means that the output of models should be as uncorrelated as possible disregarding accuracy ✗
  - Combining the same model trained in different unrelated datasets is an ensemble producing technique. ✗
  - Boosting is based on the concept of residual: reducing the difference between true label and our current approximation ✓
  - Error correcting output is an ensemble technique that extends methods to multiclass ✓
  - 1-vs-1 and winner-takes-all convert multiclass problems to binary problems. ✗
  - DL is an ensemble technique ✓

= Pill 13a (18): NNs and DL =



• Learning algorithm:  $\min_w \mathcal{L}(f(x, w), t)$

• The flow of the chain rule: basis of backpropagation algorithm and automatic differentiation.  
used in DL training

- computational graph: graphical representation of operations.

$$f(x_1, y_1, z) = (x+y) \cdot z$$

Partial derivative of each node respect to its inputs

$$F = q \cdot z$$

$$\frac{\partial F}{\partial F} = 1 \quad \frac{\partial F}{\partial x} = q \quad \frac{\partial F}{\partial q} = z \quad \frac{\partial q}{\partial x} = 1 \quad \frac{\partial q}{\partial y} = 1$$

Gradient respect to any input: (chain rule)

$$\frac{\partial F}{\partial x} = \frac{\partial F}{\partial q} \frac{\partial q}{\partial x} = z \cdot 1$$

$$\frac{\partial F}{\partial y} = \frac{\partial F}{\partial q} \cdot \frac{\partial q}{\partial y} = z \cdot 1$$

$$\frac{\partial F}{\partial z} = q$$

main idea behind automatic differentiation.

differentials can be computed as we pose the graph.

Given the following chain

$$\frac{\partial F}{\partial x} = \frac{\partial F}{\partial q} \frac{\partial q}{\partial u} \frac{\partial u}{\partial x}$$

Forward

$$\frac{\partial F}{\partial x} = \frac{\partial F}{\partial q} \left( \frac{\partial q}{\partial u} \frac{\partial u}{\partial x} \right)$$

useful for vectorial functions

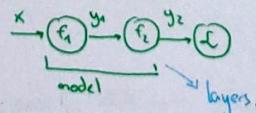
backward

$$\frac{\partial F}{\partial x} = \left( \frac{\partial F}{\partial q} \frac{\partial q}{\partial u} \right) \frac{\partial u}{\partial x}$$

good for gradients

• Deep models: composition or stacking of functions

$$y_1 = f_1(F_1(x))$$



Use of standard chain rule:

$$\frac{\partial l}{\partial w_2} = \frac{\partial l}{\partial y_2} \frac{\partial y_2}{\partial w_2}, \quad \frac{\partial l}{\partial w_1} = \frac{\partial l}{\partial y_2} \frac{\partial y_2}{\partial y_1} \frac{\partial y_1}{\partial w_1}$$

$$\frac{\partial l}{\partial y_1} = \frac{\partial l}{\partial y_2} \frac{\partial y_2}{\partial y_1}$$

$$\text{For } N \text{ layers: } \frac{\partial l}{\partial w_m} = \frac{\partial l}{\partial y_m} \prod_{i=m}^{N-1} \frac{\partial y_{i+1}}{\partial y_i} \cdot \frac{\partial y_m}{\partial w_m}$$

both gradients respect to parameters and layer input are needed.

= Pill 13b: a story of DL =

- A layer defines its output through neurons. Layer with 100 units  $\rightarrow$  100-dimensional vector.
- Inner layers use non-linearity. Outer does not.
- Batch: Approximate true loss by a noisy approximation by considering a subset of data at each parameter update.

• Softmax:  $\sigma(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$

Quiz:

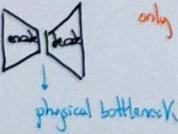
- Automatic differentiation can be regarded as a programmatical way of differentiating  $\checkmark$
- 'Flavors of automatic differentiation': back- and forward  $\checkmark$
- Backpropagation is an algorithm that sequentially uses chain rule to update different layers in a compositional model  $\checkmark$
- In standard NNs or FCNNs, each layer is composed of two ops:
  - non-linear
  - weights  $\times$  output previous layer.  $\checkmark$

= Manifold Learning =

- Manifold Learning: non-linear dimensionality reduction (S-curve). basic hypothesis: data lies in a manifold of intrinsic dimensionality smaller than the original space.
- Multi Dimensional Scaling: constructs constellation of points using info on the distances between samples.
  - Isomap (extension): performs MDS in the geodesic space of the non-linear manifold (not in original space)  
shortest distances along the curved surface represented by geodesic distances.
- Locally Linear Embedding: approximates manifold around each sample and finds the best low-dim representation to reconstruct that sample.  
weights of the embedding vector of a sample are computed from its neighbors.
- Stochastic Neighbor Embedding: (SNE) changes the notion of distance to that of conditional probabilities. (of  $x_i$  selecting  $x_j$  as neighbor under a Gaussian centered in that sample)  $P_{ij}$   
T-SNE uses a T-distribution.

= Unsupervised learning in Keras =

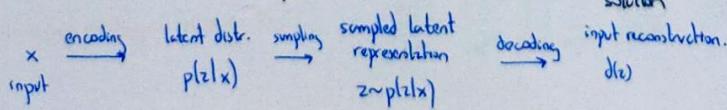
- Autoencoders:  only the relevant info is kept, very compressed



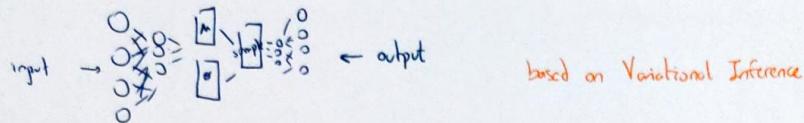
- the middle representation can be used for feature extraction

- logical bottleneck:

- Variational autoencoders: although autoencoders generate samples, the space of these samples seems to be composed of bubbles of information with bad interpolation properties among them.  $\xrightarrow{\text{VAE solution}}$



It needs Kullback-Leibler divergence loss function. Latent-space  $\sim$  Gaussian. 2 vectors should be returned



### • GANs: Generative Adversarial Networks.

Clever trick for guiding a network into mimicking sampled data. Two step process:

- 1) The generator: a network plays this role. Input comes from Normal distribution and the output is expected to be the type of data we desire. Same as decoder
- 2) The discriminator: (trick) it is trained with data from the true source and from the generator (mixed). Needs to learn to tell them apart.

Quizz:

- Isomop preserves Euclidean distance ✓
- Preserving Euclidean distance in manifold learning is generally a bad idea ✓
- LLE approximates the neighbourhood of a sample with a hyperplane ✓
- LLE is best: data  $\times$  projection matrix (in order to apply) ✗
- Manifold  $\rightarrow$  space where Euclidean intuition holds ✓
- Autoencoders need to memorize training data points ✗
- DL learns a good feature set that fits the problem at hand ✓
- Key for a good DL classifier: powerful classifier on last layer ✓
- Autoencoders need a physical or logical bottleneck in order to be effective ✓
- Autoencoders try to learn the identity.  $\rightarrow$  unsupervised ✓
- To get a representation in DL: cut at one layer and get the result ✓
- Representations can be transferred to other domains ✓

= Pill 12: Gaussian discrimination and processes =

- Quadratic forms: all eigenvals  $\neq 0 \rightarrow$  ellipsoid, hyperboloid, imaginary ellipsoid.  
 $\quad \text{all } > 0 \quad \text{all } \neq 0 \quad \text{all } < 0$
- Gaussian distribution analysis (model): assume each class comes from a Gaussian dist.

Quizz:

- The classification boundary in Gaussian distribution is the same as KNN with Mahalanobis distance ✓
- PCA uses info of the classes in a supervised problem ✗
- Eigvect with longest eigenval of the covariance matrix points to the direction with max variance ✓

= P.11 14 : Mixture models =

- Clusters put together similar points and separate different ones.