

FIRST STEPS: LINEAR CLASSIFIERS

- discussed here:
 - Perceptron
 - Adaline (Adaptive Linear Neuron)
- good entry point to modern machine learning as it is fairly simple but generalised to deep networks
- idea goes back to McCulloch & Pitts 1943

"A logical calculus of the ideas immanent in nervous activity"

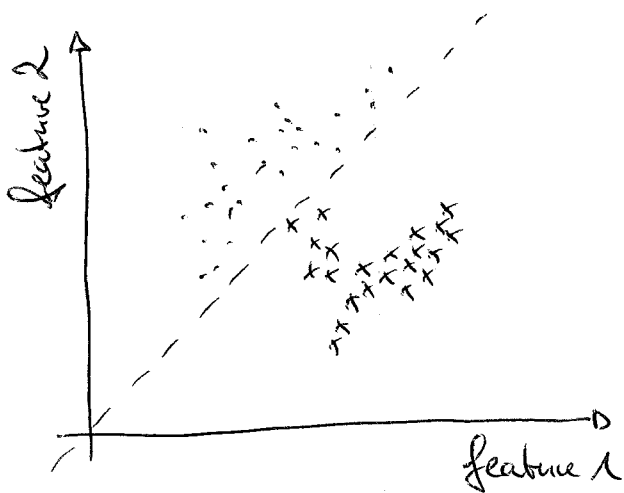
- Rosenblatt, 1957:

"The perceptron, a perceiving and recognizing automaton"

- Widrow, 1960:

"A adaptive 'Adaline' neuron using chemical membranes"

BASIC IDEA



- these data points are to be classified by assigning each data point $x^{(i)}$ a label $y^{(i)} \in \{-1, 1\}$
say $Y = (y^{(i)})_{1 \leq i \leq N}$

- problem of binary classification:

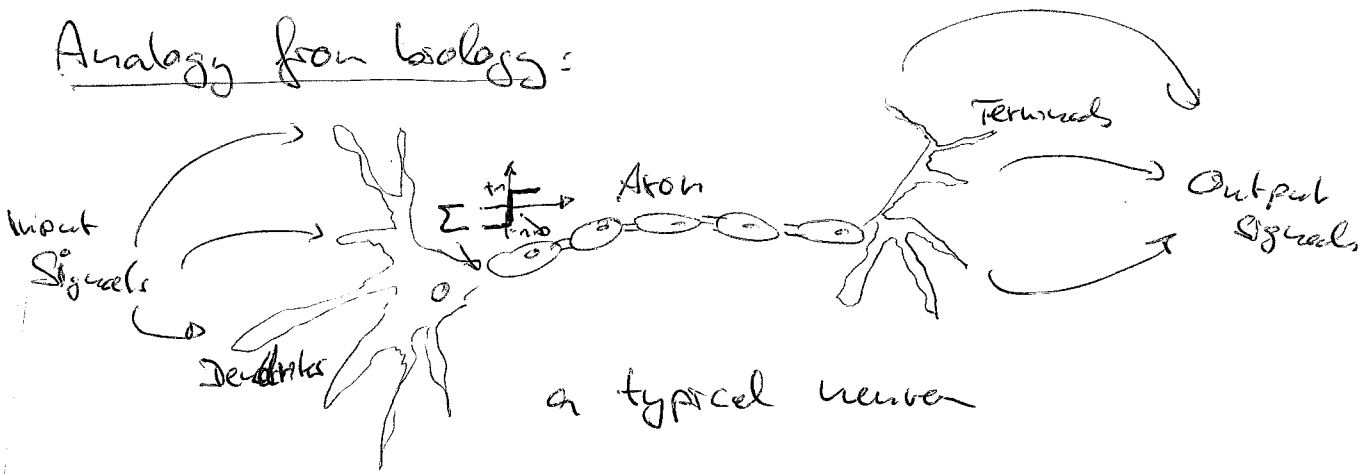
• we look for a function

$$f: \mathbb{R}^n \rightarrow \{-1, 1\}$$

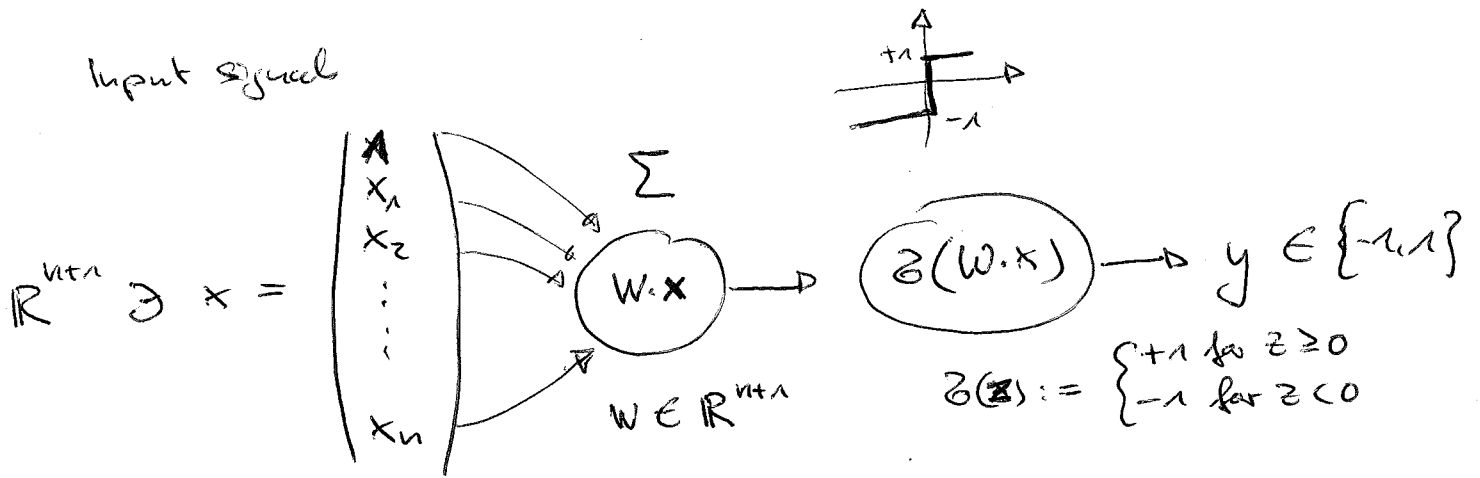
• we want the machine to find f by supervised training

- each data point consists of
for example two features $X = \begin{pmatrix} x^{(i)} \\ \vdots \\ x^{(i)} \end{pmatrix}_{1 \leq i \leq N}$
 \mathbb{R}^2

Analogy from biology:



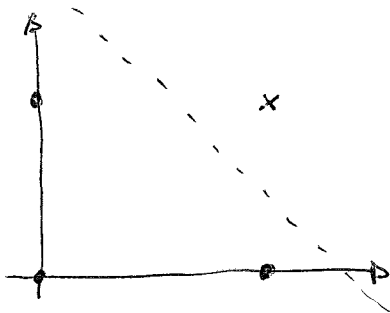
A mathematical model:



Example:

we need a linear map $z = W \cdot x$,
for example

$$W = \begin{pmatrix} -1.5 \\ 1 \\ 1 \end{pmatrix}$$



AND-gate

x_1	x_2	AND
0	0	0
0	1	0
1	0	0
1	1	1

$$W \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = -1.5 < 0$$

$$W \begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix} = W \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -0.5 < 0$$

$$W \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = 0.5 > 0$$

$$\Rightarrow f(x) = z(W \cdot x)$$

is a representation of the AND-gate

HW: Try to implement all 16 logical gates by the perceptron. Which fail?

Next we want an algorithm that allows to find a good choice of weights:

- they should classify the training data most accurately
- however, should also generalize to new data in a good way

$\Delta x \in \mathbb{R}^{n+1}$
 because
 $x = \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$

Learning rule

- given: training data $(x^{(i)}, y^{(i)})_{i \in \mathbb{N}}$

\mathbb{R}^{n+1} $\{-1, 1\}$
 w
 $(x^{(i)}, y^{(i)})$

Alg (Perceptron training):

~~Start: initialize a weight vector $w^{(0)} \in \mathbb{R}^{n+1}$ either zero or at random and $k=0$~~
~~Step:~~

- initialize weight vector $w \in \mathbb{R}^{n+1}$ to zero or at random
- for each training sample $(x^{(i)}, y^{(i)})$
 - a) compute output $y = \sigma(w \cdot x^{(i)})$
 - b) compare y and $y^{(i)}$
 - if equal do nothing
 - else update the weight vector w in a good way

For the perceptron the following update rule was suggested:

$$\Delta = y^{(i)} - y$$

$$\Delta w = \eta \Delta x^{(i)}$$

$\eta \in \mathbb{R}$ is called the learning rate

Why is this a good learning task?

• $\Delta = 0 \Rightarrow$ no error $\Rightarrow \Delta W = 0$

• $\Delta \neq 0 \Rightarrow$ error:

a) $\Delta = 2 \Rightarrow$ wanted $y = 1$ but got $y = -1$

$\Rightarrow \Delta W > 0$ which corrects in the right direction

b) $\Delta = -2 \Rightarrow$ wanted $y = -1$ but got $y = 1$

$\Rightarrow \Delta W < 0$ which again corrects

in the right direction

≡

Convergence in case training data is linearly separable

DEF: Two sets A, B of points in \mathbb{R}^n are

called: i) linearly separable iff $\exists w \in \mathbb{R}^n, w_0 \in \mathbb{R}$:

$$\forall a \in A: w a + w_0 \geq 0$$

$$\forall b \in B: w b + w_0 < 0$$

ii) absolutely linearly separable iff $\exists w \in \mathbb{R}^n, w_0 \in \mathbb{R}$

$$\forall a \in A: w a + w_0 > 0$$

$$\forall b \in B: w b + w_0 < 0$$

with $w_0 = 0$

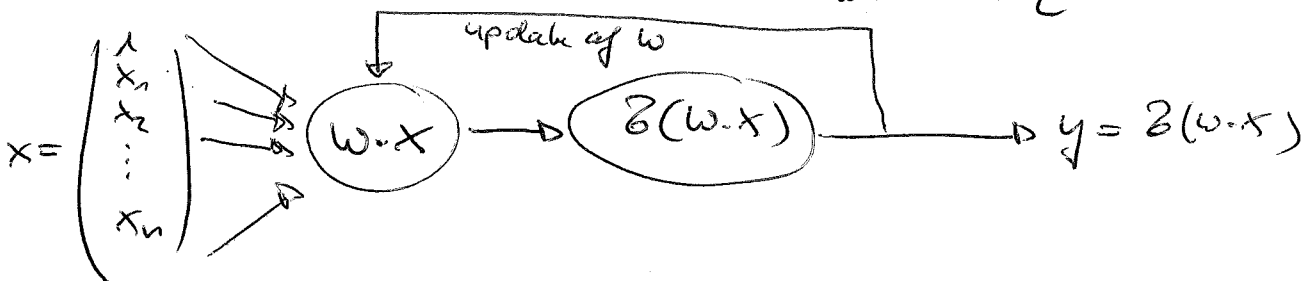
THM: ~~If A, B are linearly separable and finite, then~~

~~the perceptron learning algorithm converges~~

ALG: start: $t := 0, w^{(0)} \in \mathbb{R}^n$ randomized

step: $x \in A, B$

$$w \mapsto w + \eta (y^{(i)} - \sigma(w \cdot x^{(i)})) x^{(i)}$$



ALG (Perceptron training)

Start: initialize weight vector $w^{(0)}$ at random, $t := 0$

Step: choose $x \in A \cup B$ at random

$$x \in A \wedge w^{(t)} \cdot x \geq 0 \quad \text{goto } \underline{\text{Step}}$$

$$x \in A \wedge w^{(t)} \cdot x \leq 0 \quad \text{goto } \underline{\text{Add}}$$

$$x \in B \wedge w^{(t)} \cdot x \leq 0 \quad \text{goto } \underline{\text{Step}}$$

$$x \in B \wedge w^{(t)} \cdot x \geq 0 \quad \text{goto } \underline{\text{Sub}}$$

Add: $w^{(t+1)} := w^{(t)} + x$, $t := t + 1$ goto Step

Sub: $w^{(t+1)} := w^{(t)} - x$, $t := t + 1$ goto Step

Thm: If A, B are finite and linearly separable, then the above algorithm converges in finite many steps.

Prop: A, B finite, then A, B lin. sep $\Leftrightarrow A, B$ abs. sep.

Proof: A, B lin. sep. \Rightarrow $\exists w$:
 $\forall a \in A, w \cdot a \geq 0$
 $\forall b \in B, w \cdot b < 0$

$$\varepsilon := -\max \{ w \cdot b \mid b \in B \} > 0$$

$$w' := w + \begin{pmatrix} \varepsilon/2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\Rightarrow \forall a \in A: w' \cdot a = \overbrace{w \cdot a}^{\geq 0} + \frac{\varepsilon}{2} > 0$$

$$\forall b \in B: w' \cdot b = \overbrace{w \cdot b}^{\leq -\varepsilon} + \frac{\varepsilon}{2} < 0$$

□

Proof of Thm:

- WLOG: we may assume that the point vectors in A and B are normalized as $w \cdot x \geq 0 \iff w \cdot \frac{x}{\|x\|} \geq 0$
- since A, B are lin. sep., they are abs. lin. sep., and therefore there is a normalized w^* for which
 - $\forall a \in A, w^* a \geq 0$
 - $\forall b \in B, w^* b < 0$
- $P := A \cup (-1)B$

abs. lin. sep!

Say after some time $w^{(t+1)}$ is computed

\Rightarrow at time t $p \in P$ was incorrectly classified

$$\text{and } w^{(t+1)} := w^{(t)} + p$$

$$\cos \varphi = \frac{w^* \cdot w^{(t+1)}}{\|w^{(t+1)}\|}$$

$$\text{but 1) } w^* w^{(t+1)} = w^* w^{(t)} + w^* p \\ \geq w^* w^{(t)} + \delta$$

$$\text{for } \delta := \min \{ w^* p \mid \forall p \in P \} \geq 0$$

by induction we get

$$w^* w^{(t+1)} \geq w^* w^0 + (t+1) \delta \quad \begin{array}{l} \text{misclassified} \\ < 0 \end{array} \quad \begin{array}{l} \text{normed} \\ \text{norm} \end{array}$$

$$\text{but 2) } \|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2 \overbrace{w^{(t)} \cdot p}^{< 0} + \|p\|^2 \\ \leq \|w^{(t)}\|^2 + 1$$

again by induction we get

$$\|w^{(t+1)}\|^2 \leq \|w^{(0)}\|^2 + (t+1)$$

Collecting these inequalities we get

$$1 \geq \cos \varphi \Rightarrow \frac{w^* \cdot w^{(0)} + (t+1) \delta}{\sqrt{\|w^{(0)}\|^2 + (t+1)}} \sim \sqrt{t+1} \delta$$

grows like
indefinitely

⇒ only finite many updates are possible

⇒ Alg converges in finite many steps \square

↳ Python code

HW: 1) Implement by yourself. ~~Useful data sets.~~
 2) Visualize the weights per epoch
~~3) learn logical gates~~ 3) learn logical gates

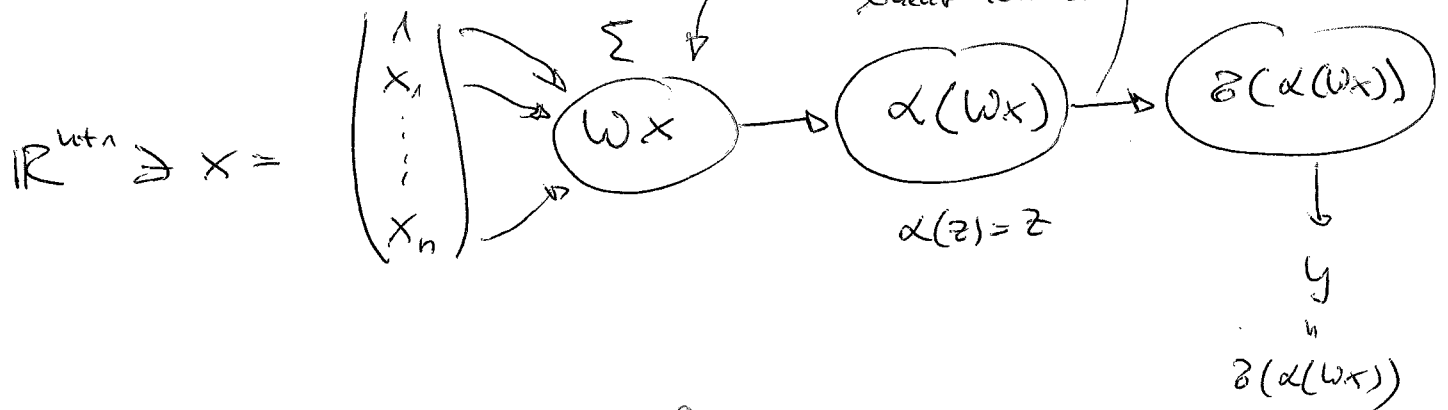
PROBLEMS OF THE PERCEPTRON

- 1) Convergence is not guaranteed in general only when A, B are lin. sep. ⇒ threshold for max. updates
- 2) If A, B are not lin. sep. at least one time every epoch w will be updated and w might start oscillate and fail to find a good decision boundary
- 3) If A, B are lin. sep. w stops being updated the moment everything in training data is classified correctly

HW: 1) Find such case (maybe by putting additional data points by hand)

ADAPTIVE LINEAR NEURON

Input Signals



Advantage:

- The activation function α is differentiable.
- Using $\alpha(w \cdot x)$ instead of $\alpha^2(w \cdot x)$ as indicator for an update allows to employ analytical optimization theory

Learning rule

def. a loss function

$$L(w) := \frac{1}{2} \sum_i [y^{(i)} - \alpha(w \cdot x^{(i)})]^2$$

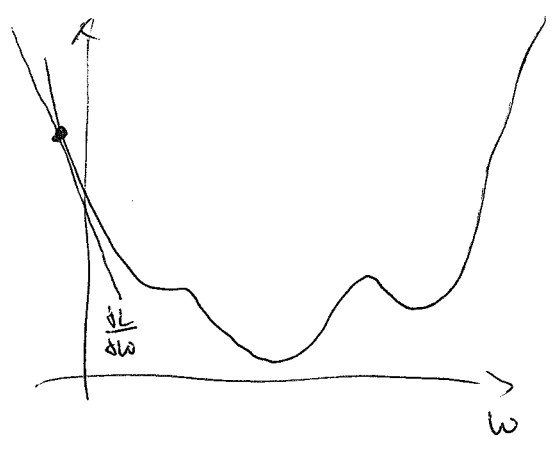
instead of minimizing directly the number of misclassifications we may minimize $L(w)$

compute: $\frac{\partial L}{\partial w} = \sum_i (y^{(i)} - \alpha(w \cdot x^{(i)})) (-1) x^{(i)}$

~~then~~ $f(x) = \sum_i (y^i - \alpha(w \cdot x^i))^2$

$$f(x) = (y - x)^2$$

$$\begin{aligned} \frac{\partial L(w)}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i (y^i - \alpha(w \cdot x^i))^2 \\ &= \sum_i (y^i - \alpha(w \cdot x^i)) \left(\frac{\partial [\alpha(w \cdot x^i)]}{\partial w} \right) \\ &= \sum_i (y^i - \alpha(w \cdot x^i)) \alpha'(w \cdot x^i) x_j^i \end{aligned}$$



to get closer to a local minimum we may use the gradient $\frac{dL}{dw}$ to give a good direction

$$\Delta w = -\eta \frac{\partial L(w)}{\partial w}$$

for the situation

$$L(w + \Delta w) = L(w) + \frac{\partial L(w)}{\partial w} \cdot \Delta w + O(\Delta w^2)$$

choosing $\Delta w = -\eta \frac{\partial L(w)}{\partial w}$ gives

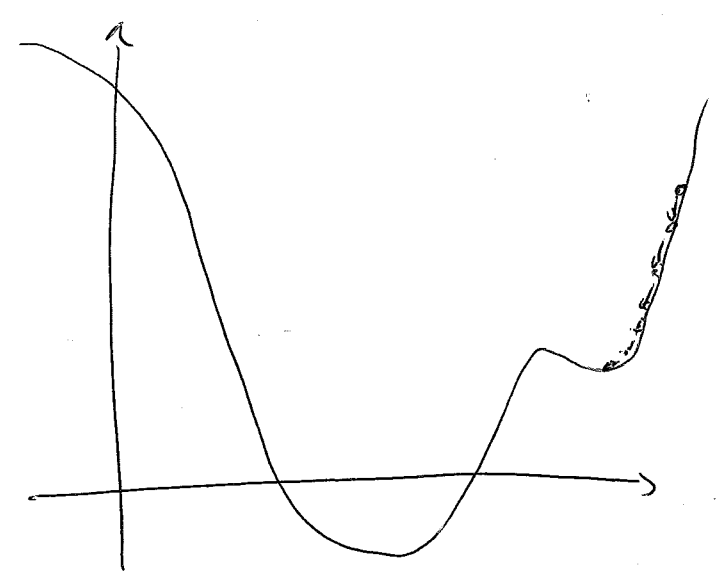
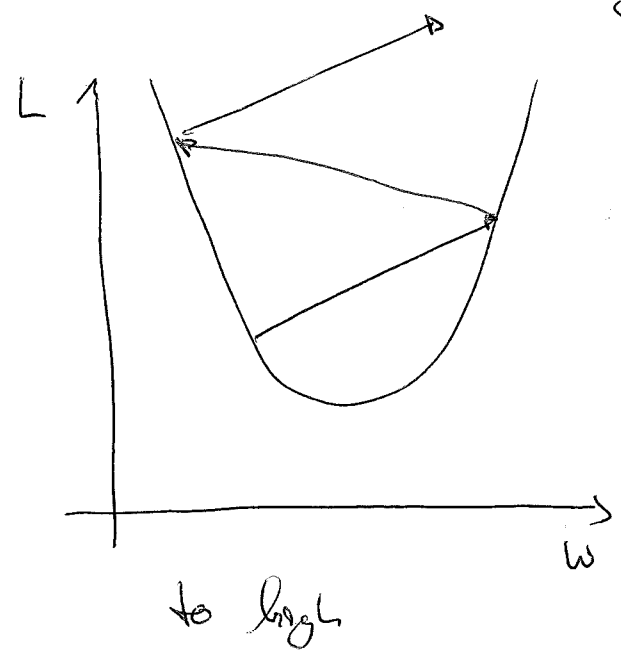
$$L(w + \Delta w) = L(w) - \eta \left(\frac{\partial L(w)}{\partial w} \right)^2 + O(\Delta w^2)$$

linear perturbation has the right direction

Differences to Perceptron : not online but batch updates
 • even in sep. case, no convergence guaranteed?
 HW: ~~try other L and alpha~~
 keep it yourself.
 • How could online learning be performed

↳ Python code

How to choose the learning rate?



HW: 1) Find such inappropriate learning rate. 2) Try other L and alpha? 3) Consider Ad. Bd examples

A word about the preparation of the data

• the choice of learning rate depends on the fluctuation of the data

⇒ usually good practice to "normalize" the data

$$\tilde{x}^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$$

$$\text{for } \mu = \frac{1}{N} \sum_i x^{(i)}, \quad \sigma = \sqrt{\frac{1}{N} \sum_i (x^{(i)} - \mu)^2}$$

Compare the 2 sets data normalized or not.
NE (we're in the library)

HW: ~~code for other data!~~ 1) Prepare it: train, test, all, normalize and train a model. 2) choose $\eta = \frac{1}{N}$ and discuss

Online Learning & Stochastic Gradient Descent

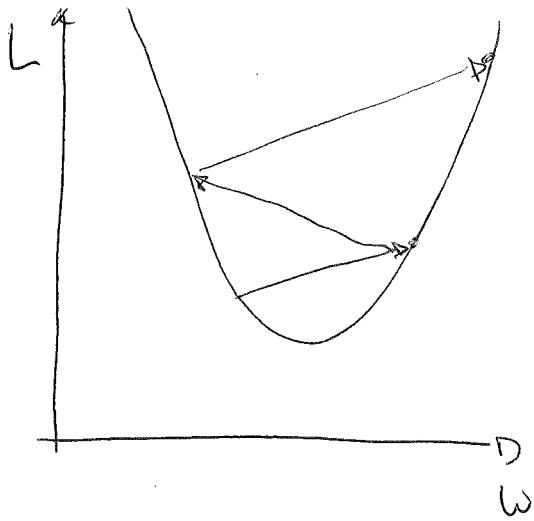
- Adaline so far compute ΔW depending on all training data - can be expensive!
- Perceptron performed "online" learning after each newly encountered training data point
- ~~Perceptron~~ This could also be done for Adaline in which case it is called "stochastic gradient descent"
- again ~~we don't~~ got to avoid to strong dependence on the sequence of learning or can do "mini-batch^{st.} gradient descent".

HW: 1) Make Adaline an online learner
2) Make Perceptron a batch learner using average
3) Make both mini-batch learner.

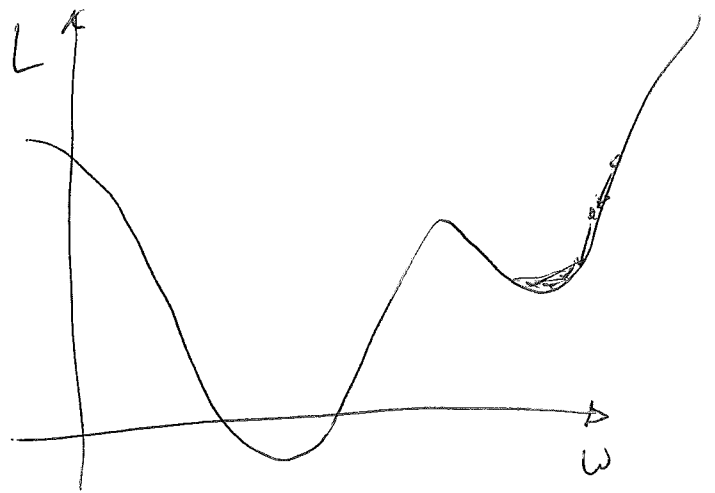
Choosing a good learning rate: gradient descent

IV

$$w \mapsto \tilde{w} := w - \eta \frac{\partial L(w)}{\partial w}, \quad L(w) := \sum_{i=1}^n |y^{(i)} - \alpha(w \cdot x^{(i)})|^2$$



too high



too low

HW: implement on test case for both situations in 1d, 2d and plot.

The choice of learning rate depends on the fluctuation of training data:

\Rightarrow good practice to normalize data

$$\tilde{x}^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$$

$$\text{for } \mu = \frac{1}{N} \sum_{i=1}^N x^{(i)}, \quad \sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu)^2}$$

HW: • Try with Iris data set, random data, or some standard training data from UCI ML Lib.

Outline of Batch Learning

- Adaline so far computes ΔW given the batch of all training data — batch learning
- Perceptron performed an update after each misclassification — online learning
- In many situations it is less expensive and one gains convergence speed by not using the whole batch for an update

→ Stochastic gradient descent
"online"

$$L(w) = \sum_{i=1}^n L_i(w) \quad L_i(w) = \frac{1}{2} |w \cdot x^{(i)} - y^{(i)}|^2$$

$$w \mapsto \tilde{w} := w - \eta \frac{\partial L_i(w)}{\partial w} \text{ for sample } i$$

- A compromise between these two algorithms is "mini-batch" gradient descent.

$n=1, \dots, M$ training data points divided into distinct subsets $I_1 \dots I_m$ and m -times one performs the update

$$w \mapsto \tilde{w} := w - \eta \frac{\partial \sum_{i \in I_k} L_i(w)}{\partial w}$$

- to avoid a strong dependence on the sequence of these data-points one can average by taking

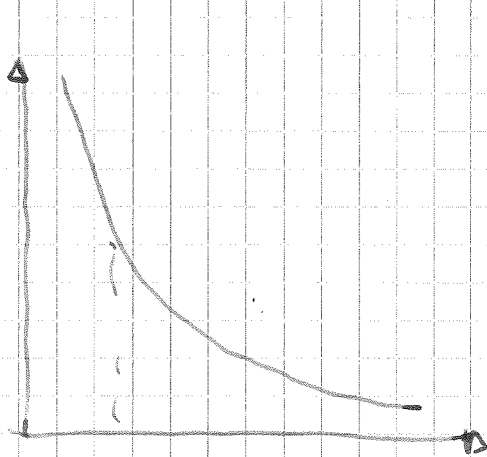
$$\eta = (\# \text{ sample points})^{-1}$$

- HW:
- 1) Make Adaline an online learner
 - 2) Make Perceptron a batch learner with averaged update
 - 3) Make both mini-batch learners.

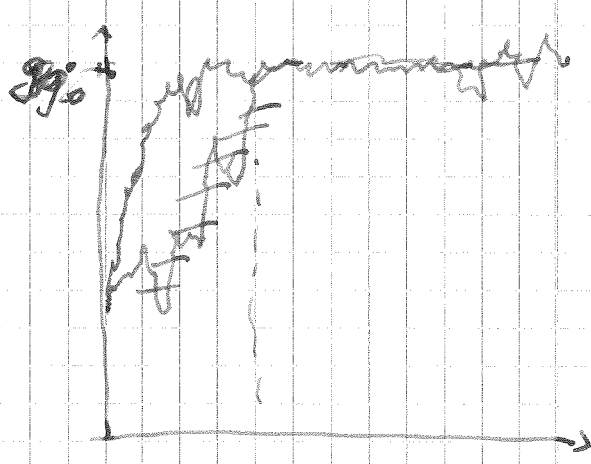
Overfitting

- as the convergence is not guaranteed
1. general one must find a good criterion when to stop the learning process
- for Adaline there is the danger of "overfitting":
 - practicing the same training data over and over again W gets better and better adapted to the training data
 - however, that usually means that W generates worse and worse to unseen data
 - overfitting training data could turn out to be a perverse way to fit slow the training data

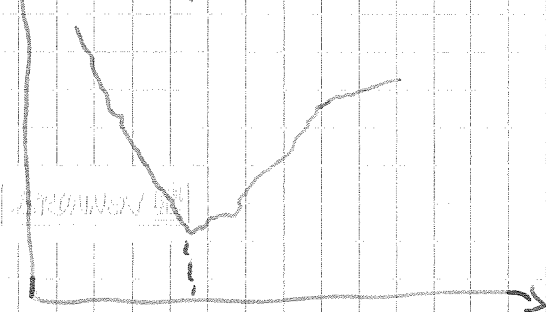
cost of training data



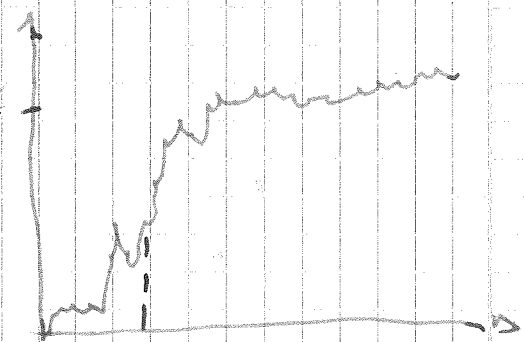
accuracy on training



cost of test data



80%



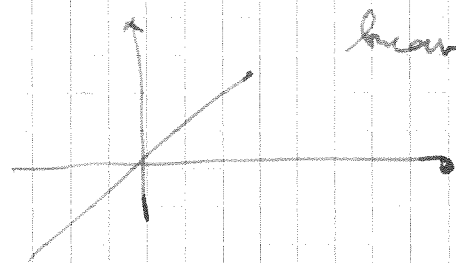
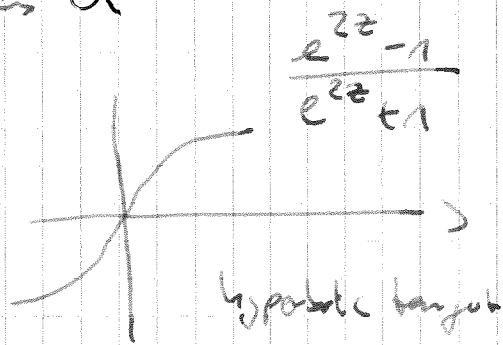
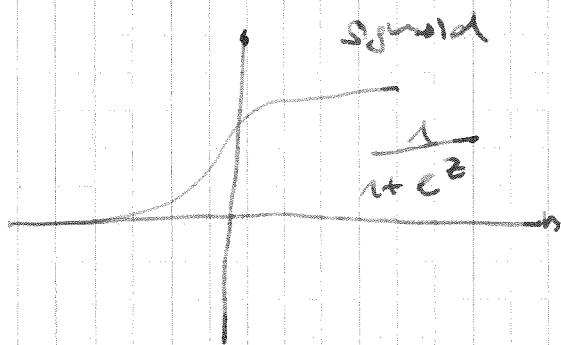
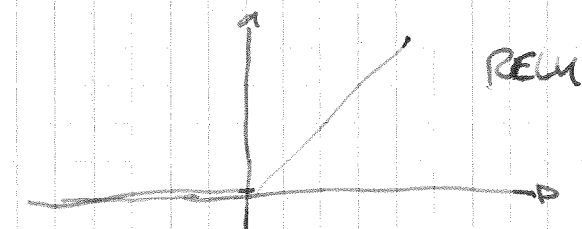
What's next?

- freedom to incorporate additional terms in the loss function

- e.g. $\|w\|_p$ norms to control what kind of weights are preferred

- different distance relations such as L^p norms

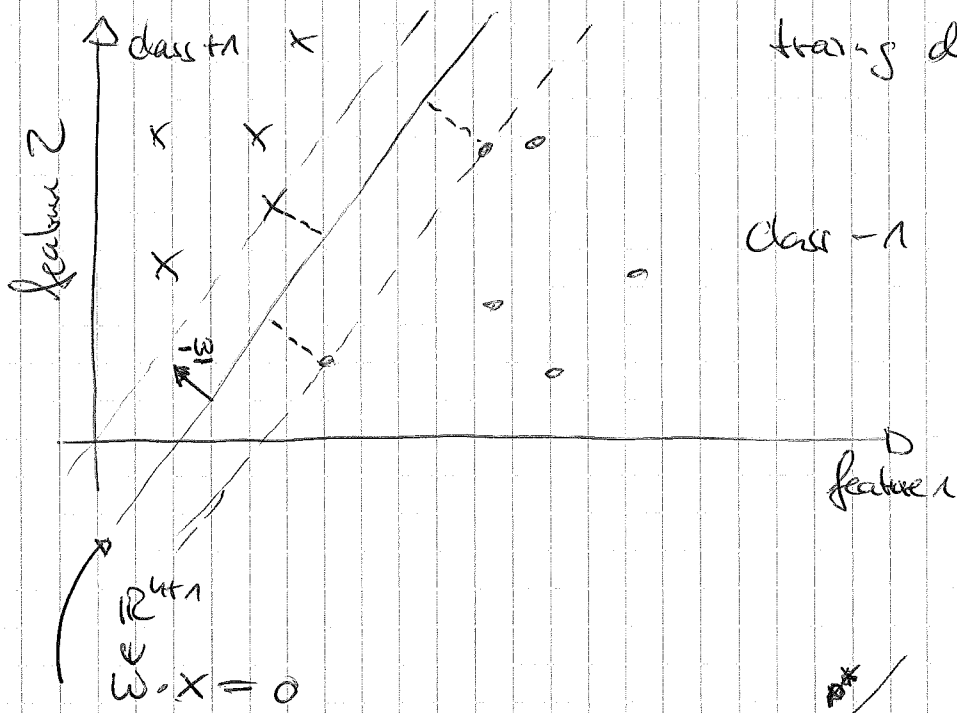
- different activation functions α



etc...

HW: implement some possibilities and discuss the behavior on diff. log. training data.

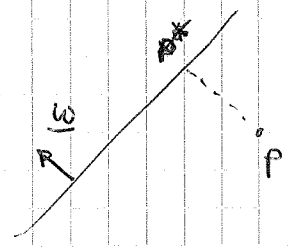
SUPPORT VECTOR MACHINES (lin. sep. case)



training data $(x^{(i)}, y^{(i)})_{i=1, \dots, n}$

class -1

$$\begin{aligned} \text{dist}(p, w) &= ? \\ &= \|p - p^*\| \\ &= \|\alpha \underline{w}\| \\ &= \left\| \frac{w \cdot p}{\|w\|^2} \underline{w} \right\| = \frac{|w \cdot p|}{\|w\|} \end{aligned}$$



$$\begin{aligned} p^* &= p + \alpha \underline{w} \\ \text{for } w \cdot p^* &= 0 \Leftrightarrow w \cdot p + \alpha \|w\|^2 = 0 \\ &\Leftrightarrow \alpha = -\frac{w \cdot p}{\|w\|^2} \end{aligned}$$

Want to find hyperplane w that maximizes the margin:

$$f_w \Leftrightarrow \min_{i=1, \dots, n} \frac{|w \cdot x^{(i)}|}{\|w\|}$$

- for this purpose note that $w \cdot x = 0$ is scale independent \Rightarrow we may scale w s.t.

$$\min_{i=1, \dots, n} |w \cdot x^{(i)}| = 1$$

- This means there could be points x_+ and x_- in the training data s.t.

$$w \cdot x_+ = +1$$

$$w \cdot x_- = -1$$

- With this we find

$$\xi_{\text{opt}} = \frac{1}{\|w\|}$$

Optimization
Problem

We may thus def the optimal choice w in terms of an optimization problem as follows:

I) maximum $\frac{1}{\|w\|}$ subject to $\begin{cases} w \cdot x^{(i)} \geq +1 & \text{if } y^{(i)} = +1 \\ w \cdot x^{(i)} \leq -1 & \text{if } y^{(i)} = -1 \end{cases}$
 $w \in \mathbb{R}^{d+1}$ for $i=1, \dots, M$

or equivalently

II) minimum $\frac{1}{2} \|w\|^2$ subject to $y^{(i)} w \cdot x^{(i)} \geq 1$
 $w \in \mathbb{R}^{d+1}$ for $i=1, \dots, M$

Prop. If the training data is finite and lin. sep. there is a unique solution.

Proof: The objective function $F(w) = \frac{1}{2} \|w\|^2$

is smooth, $\nabla_w F(w) = w$, Hessian $\nabla_w^2 F(w) = \mathbb{1}$

$\Rightarrow F$ strictly convex

$\Leftrightarrow \forall \alpha \in (0, 1), \forall w \neq \tilde{w}$:

$$F(\alpha w + (1-\alpha)\tilde{w}) < \alpha F(w) + (1-\alpha)F(\tilde{w})$$

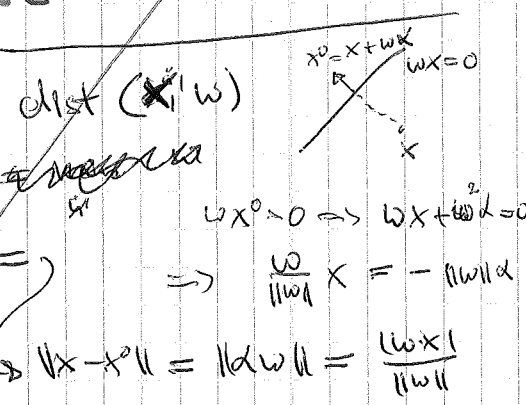
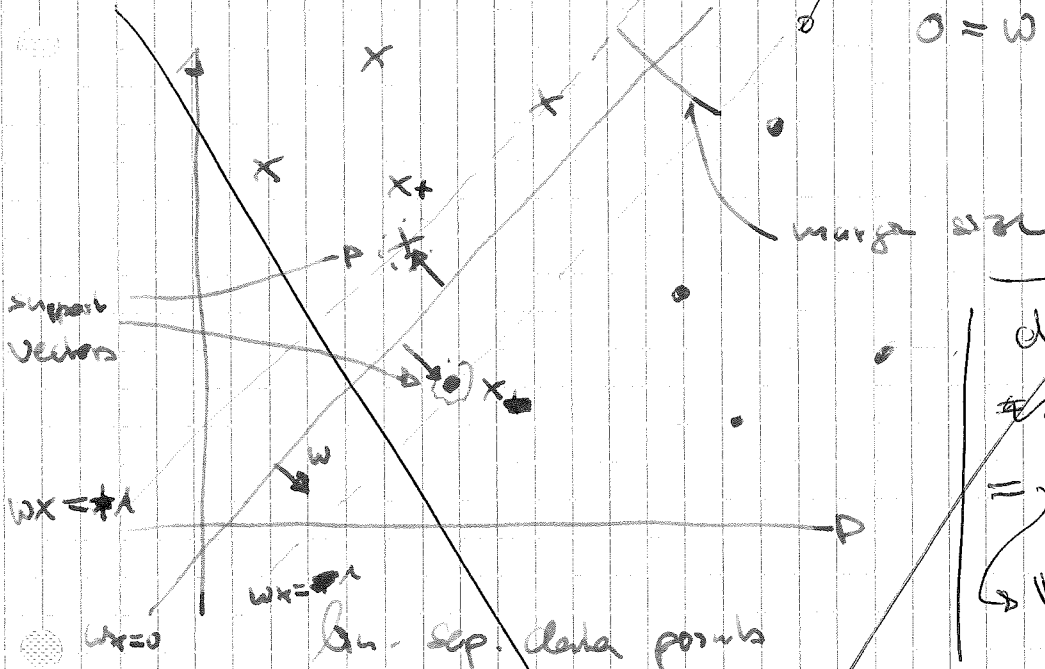
Support Vector machines

Scale s.t. $w \cdot x_+ = +1$
 $w \cdot x_- = -1$

optimal hyperplane

$$0 = w \cdot x = w_0 + \underline{w} \cdot \underline{x}$$

VII



- since $w \cdot x = 0$ also $\frac{w}{\|w\|} \cdot x = 0$
and we may normalize w at any other scalar
- choose a scalar s.t. $w \cdot x_+ = +1$
 $w \cdot x_- = -1$

$$\frac{w \cdot (x_+ - x_-)}{\|w\|^2} = \frac{2}{\|w\|} = \text{size of margin}$$

Thus, an optimal choice for w can be formulated by

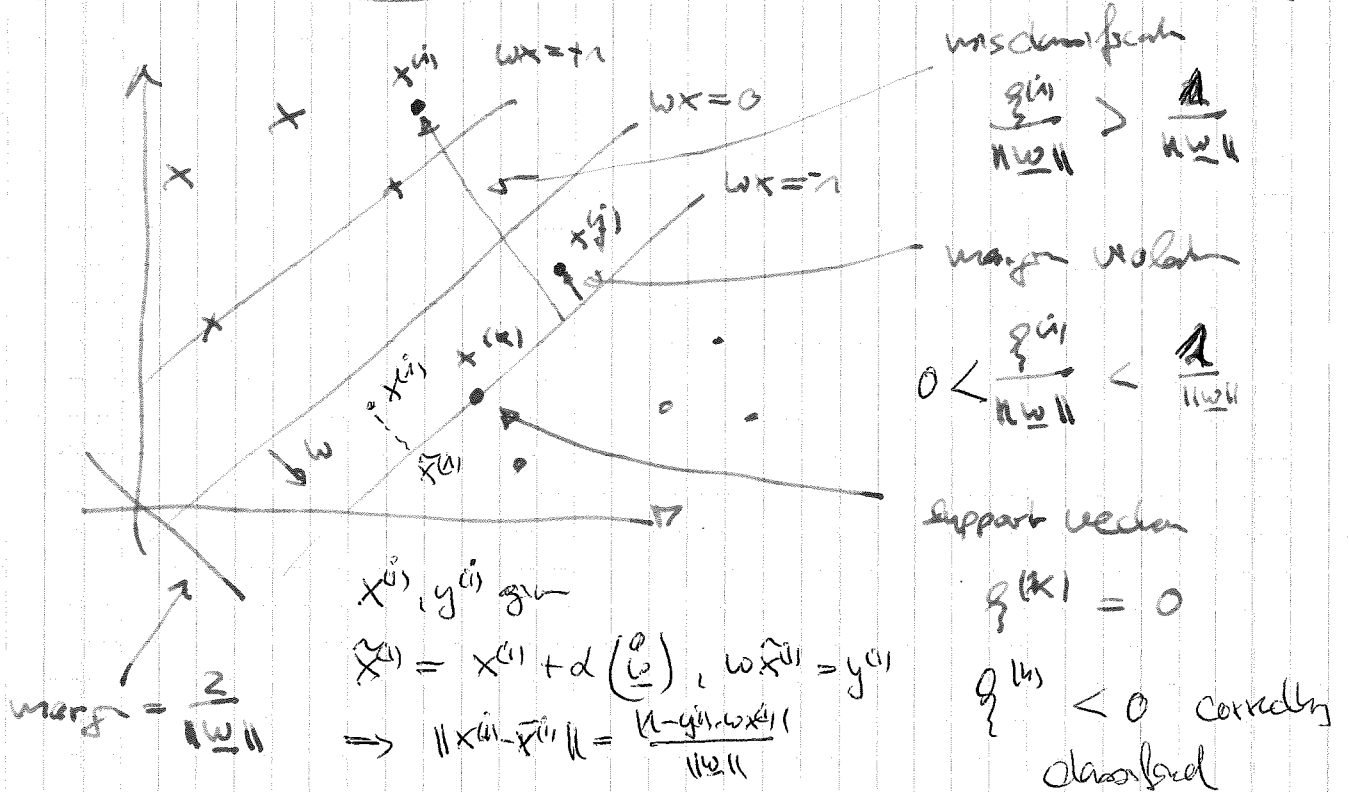
- maximum $\frac{2}{\|w\|}$ subject to $w x^{(i)} \geq 1$ for $y^{(i)} = 1$
 $w x^{(i)} \leq -1$ for $y^{(i)} = -1$
for $i = 1 \dots n$
- or equivalently
minimum $\|w\|^2$ subject to $y^{(i)} w x^{(i)} \geq 1$
for $i = 1 \dots n$

This is a quadratic optimization problem
subject to linear constraints

→ there is a unique minimum

$F(w) = \frac{1}{2} \|w\|^2$ obj. fun. $\nabla F = w$, $\nabla^2 F = I \Rightarrow$ strictly convex, constraint affine
 \Rightarrow local min = global min and since strictly convex, min is unique

The general case for only approximately lin. sep. data



Soft margin optimization problem

minimum
 $w \in \mathbb{R}^{n+1}$
 $\xi^{(i)} \in \mathbb{R}^+$

$$\|w\|^2 + \mu \sum_{i=1}^n \xi^{(i)}$$

subject to $y^{(i)} (w \cdot x) \geq 1 - \xi^{(i)}$

- every constraint can be satisfied for $\xi^{(i)}$ s.t. large
- $\mu \in \mathbb{R}^+$ is a regularization parameter that controls the weight of the margin
- still a quadratic optimization problem \Rightarrow unique solution with lin. constraints

OPTIMIZATION THEORY



History: 1) 1629: Fermat's theorem

$$f(x) \rightarrow \text{extrema}$$

THM: Let $f: \mathbb{R}^n \rightarrow \mathbb{R}$ diff at $x^* \in \mathbb{R}^n$. If x^* is a local extremum, then $f'(x^*) = 0$. (stab. point).

Then discussion of boundary, non-diff., stationary points.

2) 1788: Lagrange's theorem

$$f_0(x) \rightarrow \text{minimum subject to equality constraints} \\ f_1(x), \dots, f_m(x) = 0$$

THM: Let $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ be cont. diff in the vicinity of a local extremum, then one can find $z^* \in \mathbb{R}^{m+1}$ s.t.

$$L(x^*, z^*) = \sum_{i=0}^m z_i^* f_i(x^*), \quad z^* \neq 0$$

fulfills

$$\nabla_{(x,z)} L(x, z) \Big|_{(x,z)=(x^*, z^*)} = 0 \quad (1)$$

$\exists \lambda_i$ $\exists f_i(x^*)$ $i=1, \dots, m$ are linear λ -dependent $\lambda_0 \neq 0$.

After solving (1) we find stab. point to discuss.

3) 1851: Karush-Kuhn-Tucker theorem

$$f_0(x) \rightarrow \text{minimum subject to constraints} \quad \left. \begin{array}{l} f_1(x), \dots, f_m(x) \leq 0 \end{array} \right\} \text{minimization program}$$

Thm: Let $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$, $i=0, \dots, m$, be convex. If x^* is a minimum of above program, then $\exists z^* \in \mathbb{R}^{m+1}$, $z^* \neq 0$, s.t.

i) $\inf_x L(x, z^*) = L(x^*, z^*)$

ii) $z_i^* \geq 0$ for $i=0, \dots, m$

iii) $z_i^* f_i(x^*) = 0$ $i=1 \dots m$

a) If $z_0^* \neq 0$, (i)-(iii) $\Rightarrow x^*$ a solution of unconstrained problem.

(Slater condition) b) If $\exists \bar{x}$ s.t. $f_i(\bar{x}) < 0$ $i=1 \dots m \Rightarrow z_0^* \neq 0$.

c) If the Slater cond. is satisfied ~~##~~ ~~##~~ the Lagrange function

$$L(x, z) = f_0(x) + \sum_{i=1}^m z_i f_i(x)$$

has a saddle point at $(x^*, z^*/z_0^*)$, i.e.,

$$L(x, \frac{z^*}{z_0^*}) \geq L(x^*, \frac{z^*}{z_0^*}) \geq L(x^*, z/z_0^*) \quad \forall x, z, z \neq 0$$

Existence of optimal solutions

II

Optimization program:

$$(P) \quad f_0(x) \rightarrow \text{maximum subject to} \quad \begin{aligned} f_i(x) &\leq 0 & i=1 \dots p \\ f_j(x) &= 0 & j=p+1 \dots m \\ x &\in C \subseteq \mathbb{R}^n \end{aligned}$$

Thm: Let C be closed, $f_i: \mathbb{R}^n \rightarrow \mathbb{R}$ cont. and f_0 coercive over $S := \{x \in C \mid f_i(x) \leq 0, f_j(x) = 0 \text{ for } i=1 \dots p, j=p+1 \dots m\}$

$\Leftrightarrow \forall (x_n)_{n \in \mathbb{N}}$ in S with $\|x_n\| \xrightarrow{n \rightarrow \infty} \infty$ also $|f_0(x_n)| \rightarrow \infty$,

then (P) has at least one optimal solution.

Proof: • C closed, f_i cont $\Rightarrow S$ closed

• let $(x_n)_{n \in \mathbb{N}}$ be a sequence in S s.t.

$$f_0(x_n) \xrightarrow{n \rightarrow \infty} \inf_{x \in S} f_0(x)$$

• since f_0 is coercive and $(f_0(x_n))_{n \in \mathbb{N}}$ obviously bounded also $(x_n)_{n \in \mathbb{N}}$ is bounded

• $\Rightarrow \exists (x_{n_k})_{k \in \mathbb{N}}$ s.t. $x_{n_k} \xrightarrow{k \rightarrow \infty} x^* \in S$
since S was closed

• but f_0 cont. $\Rightarrow f_0(x^*) = \lim_{k \rightarrow \infty} f_0(x_{n_k}) = \inf_{x \in S} f_0(x)$

Therefore, at least there is one x^* satisfying (P). \square

Thm: Let $f_i, i=1, \dots, m$, convex, then:

i) f_0 convex \Rightarrow local min = global min

ii) f_0 strictly convex \Rightarrow min. unique.

Proof: i) \circ let $x^* \in S$ be a local min.

$$\Rightarrow \exists \varepsilon > 0 : \forall x \in B_\varepsilon(x^*) \cap S : f(x^*) \leq f(x)$$

\circ As $f_i, i=1, \dots, m$ are convex, also S is convex

\circ take $y \neq x, y \in S, \lambda \in (0, 1)$ s.t.

$$x^* + \lambda(y - x^*) \in B_\varepsilon(x^*)$$

$$\Rightarrow f(x^*) \leq f(x^* + \lambda(y - x^*)) \leq (1-\lambda)f(x^*) + \lambda f(y)$$

$$\Rightarrow f(x^*) \leq f(y)$$

ii) Same argument with strict inequalities. \square

Proof of our Props about SVTs:

$$f(w) = \frac{1}{2} \|w\|^2 \Rightarrow \nabla_w f(w) = w \quad (\text{gradient})$$

$$\Rightarrow \nabla_w^2 f(w) = \mathbb{1} \quad (\text{Hessian})$$

Hence, f_0 is strictly convex, and also coercive, \mathbb{R}^d closed.

Furthermore, f_i are affine $i=1, \dots, m$, \Rightarrow convex

\Rightarrow above Thms apply and there is a unique solution.

\square

Convex Optimization Problems

Proof of L^* 1-2

Optimization program:

(without $f(x) = g(x)$)

$$(P) \quad \min_{x \in \mathbb{R}^n} f(x) \quad \text{subject to} \quad \begin{aligned} f_i(x) &\leq 0 & i=1 \dots p \\ f_j(x) &= 0 & j=p+1 \dots m \\ x &\in C \end{aligned}$$

DEF: $S := \{x \in C \mid f_i(x) \leq 0, f_j(x) = 0, i=1 \dots p, j=p+1 \dots m\}$

ASSUMPTIONS:

(A1) $C \subseteq \mathbb{R}^n$ convex s.t. $\text{supp } f_i \supseteq C, i=1 \dots m$

(A2) $f_i: \mathbb{R}^n \rightarrow \overline{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$ convex for $i=1 \dots p$

(A3) $f_j: \mathbb{R}^n \rightarrow \mathbb{R}$ affine for $j=p+1 \dots m$

REF: (A1)-(A3) $\Rightarrow S$ convex

and furthermore if C closed & f_i f_i f_i continuous (implied by convexity) $\Rightarrow S$ closed

ASSUMPTIONS:

(A4) for convex S holds:

a) $\exists \bar{x} \in S \cap C^\circ$

b) $\forall f_i, i=1 \dots p$, not affine $\exists x_i \in S: f_i(x_i) < 0$

This is called Slater condition.

THM: (P) fulfilling (A1)-(A3)

• (A4) Slater cond

• $\alpha := \inf \{f(x) \mid x \in S\} \in \mathbb{R}$

$\Rightarrow \exists y \in \mathbb{R}^m, y_i \geq 0, i=1 \dots p$ s.t.

$$f(x) + \sum_{i=1}^p y_i f_i(x) \geq \alpha \quad \text{for all } x \in C$$

REM: Only feasibility of α is assumed not existence of $x^* \in C$ with $f(x^*) = \alpha$, however, if $\alpha^* \in C$ is the optimal value $f(x^*) = \alpha$ holds.

Proof: • Preliminary: (A4) b) $\Rightarrow \exists \hat{x} \in S$ st. $f_i(\hat{x}) < 0$
for all non-active functions, say,
 $i = 1 \dots l$

Proof: $\forall i = 1 \dots l \exists x_i \in S$ with $f_i(x_i) < 0$, $f_k(x_i) \leq 0$
and $f_j(x_i) = 0$ for $k \neq i, k = 1 \dots p$

$\hat{x} := \frac{1}{l} (x_1 + x_2 + \dots + x_l) \in S$ due to convexity

$$\begin{aligned} \text{and } f_i(\hat{x}) &= f_i\left(\frac{1}{l} \sum_{n=1}^l x_n\right) \\ &\leq \frac{1}{l} \sum_{i=1}^l f_i(x_n) = \frac{1}{l} \left(\underbrace{f_i(x_1)}_{\leq 0} + \dots + \underbrace{f_i(x_l)}_{< 0} + \dots + \underbrace{f_i(x_l)}_{\leq 0} \right) \\ &< 0 \quad \square \end{aligned}$$

• To simplify the cases, let us assume $\exists \hat{x} \in S$
with $f_i(\hat{x}) < 0$ for $i = 1 \dots p$ (also for the active function)

(*) • By (A4) a) $\exists \bar{x} \in S \cap C^0$ so that by convexity
 \bar{x} can be selected to fulfill $f_i(\bar{x}) < 0 \quad i = 1 \dots p$

• WLOG $\alpha = 0$

PART 1: Without (A4) we prove that $\exists z \in \mathbb{R}^{m+1}$, $z \neq 0$:

$$z_i \geq 0, i = 1 \dots p \quad \wedge \quad \sum_{i=0}^p z_i f_i(x) \geq 0 \quad \forall x \in C$$

PART 2: With (A4) we show that $z_0 > 0$

$$\Rightarrow y := (z_i / z_0)_{i=1 \dots m}$$

Proof of Part 1s

Proof of
LEMMA 1-2

$$A := \left\{ v \in \mathbb{R}^{m+1} \mid \exists x \in C : \begin{array}{l} v_0 > f_0(x) \\ v_i \geq f_i(x) \quad i=1, \dots, p \\ v_j = f_j(x) \quad j=p+1, \dots, m \end{array} \right\}$$

C convex $\wedge f_i$ convex $i=1, \dots, p$ $\wedge f_j$ affine $j=p+1, \dots, m$
 $\Rightarrow f_j$ convex

$\Rightarrow A$ convex

- $0 \notin A$ because of $\alpha = 0$
- $A \neq \emptyset$ as v_i can be chosen arbitrarily negative

LEMMA: $A \neq \emptyset, 0 \notin A, A$ convex $\subseteq \mathbb{R}^{m+1}$

$\Rightarrow \exists z \in \mathbb{R}^{m+1}, z \neq 0$ s.t. i) $z \cdot v \geq 0 \quad \forall v \in A$
 ii) $\exists \bar{v} \in A, z \cdot \bar{v} > 0$

- But $w \in A \Rightarrow v + \alpha w \in A \quad \forall \alpha \geq 0, w \geq 0$
 and $\lim_{\alpha \rightarrow \infty} z \cdot (v + \alpha w) \geq 0 \Rightarrow z \cdot w \geq 0 \Rightarrow z \geq 0$

for all $x \in C \exists v_\epsilon := \begin{pmatrix} f_0(x) + \epsilon \\ f_1(x) \\ \vdots \\ f_m(x) \end{pmatrix} \in A$

$$0 \leq z \cdot v_\epsilon = z_0 (f_0(x) + \epsilon) + \sum_{i=1}^p z_i f_i(x)$$

$$\epsilon \rightarrow 0 \Rightarrow 0 \leq z_0 f_0(x) + \sum_{i=1}^p z_i f_i(x), \quad z_i \geq 0, z \neq 0$$

Proof of Part 2:

• Let us assume $z_0 = 0$

$$\Rightarrow v = (f(x)+1, f_{j=1-p}^{(x)}, \overset{0, \dots, 0}{\cancel{\dots}}) \in A$$

because of (A*) a) (see *)

$$\Rightarrow z \cdot v \geq 0$$

but since $z_0 = 0$, $f_j(x) < 0$

$$\Rightarrow z_0, z_1, \dots, z_p = 0$$

$$\Rightarrow \text{due to } A: \sum_{j=p+1}^m z_j f_j(x) \geq 0 \quad \forall x \in C$$

and since $\{0\}$ and A well separated $\exists \tilde{x} \in C$:

$$\sum_{j=p+1}^m z_j f_j(\tilde{x}) > 0$$

• because $\tilde{x} \in C^\circ$ we have $\hat{x} - \varepsilon(\tilde{x} - \hat{x}) \in C$

• since f_j affine for $j = p+1, \dots, m$

$$\begin{aligned} f_j(\hat{x} - \varepsilon(\tilde{x} - \hat{x})) &= \underbrace{f_j(\hat{x})}_{=0} - \varepsilon [f_j(\tilde{x}) - f_j(\hat{x})] \\ &= -\varepsilon f_j(\tilde{x}) \end{aligned}$$

$$\Rightarrow \sum_{j=p+1}^m z_j f_j(\hat{x} - \varepsilon(\tilde{x} - \hat{x})) = -\varepsilon \sum_{j=p+1}^m z_j f_j(\tilde{x}) < 0$$

$$\Rightarrow z_0 \neq 0 \text{ but } z_0 > 0. \quad \square$$

Optimality cond. for constraint convex optimization

III
I

AIM: Derive cond. which allow to decide whether a point $x \in S$ is optimal or not.
→ basis for many numerical receipts

Question: "the solvability of a system of equalities and inequalities".

$$(P) \quad f(x) \rightarrow \text{minimum} \quad \text{subject to} \quad \begin{array}{ll} f_i(x) \leq 0 & i=1..p \\ f_j(x) = 0 & j=p+1..m \\ x \in C \end{array}$$

$$\text{DEF: } S := \{x \in C \mid f_i(x) \leq 0, f_j(x) = 0 \quad i=1..p, j=p+1..m\}$$

ASSUMPTIONS:

$$(A_{\text{convex}}) = \begin{cases} (A1) & C \subseteq \mathbb{R}^n \text{ convex, } C \subseteq \text{supp } f_i \quad i=1..m \\ (A2) & f_i: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\} \text{ convex for } i=1..p \\ (A3) & f_j: \mathbb{R}^n \rightarrow \mathbb{R} \text{ affine for } j=p+1..m \end{cases}$$

REF: $(A_{\text{convex}}) \Rightarrow S$ convex and closed (by constraints)

ASSUMPTION:

$$(A_{\text{ Slater}}) = \begin{cases} (A4) & \text{a) } \exists \bar{x} \in S \cap \overset{\circ}{C} \\ & \text{b) } \forall i=1..p, f_i \text{ non-affine } \exists x_i \in S: f_i(x) < 0 \end{cases}$$

The strategy will use the hyperplane separation theorems for convex sets to arrive at the lemmas:

LEM 1: Let (P) fulfill (A_{convex}) and let

$$\alpha := \inf_{x \in S} f_0(x) \in \mathbb{R} \quad (\text{finiteness only w/ existence of an optimal } x^* \in S)$$

then $\exists z \in \mathbb{R}^{m+1}$, $z \neq 0$, $z_i \geq 0$, $i=1 \dots p$ s.t.

$$z_0 [f_0(x) - \alpha] + \sum_{i=1}^m z_i f_i(x) \geq 0 \quad \forall x \in C.$$

LEM 2: Assumptions as in LEM 1 but in addition $(A_{\text{ Slater}})$,

then $\exists y \in \mathbb{R}^m$, $y_i \geq 0$, $i=1 \dots p$ s.t.

$$[f_0(x) - \alpha] + \sum_{i=1}^m y_i f_i(x) \geq 0 \quad \forall x \in C.$$

Implications: Let x^* be a optimal solution of (P)

$$\rightarrow \alpha = f_0(x^*)$$

$(A_{\text{convex}}) \xrightarrow{\text{LEM 1}} \exists z \in \mathbb{R}^{m+1}$, $z \neq 0$, $z_i \geq 0$ for $i=1 \dots m$ s.t.

$$z_0 [f_0(x) - f_0(x^*)] + \sum_{i=1}^m z_i f_i(x) \geq 0 \quad (*)$$

Suppose f_i , $i=0 \dots m$, are diff, then

$$\begin{cases} \text{(I)} & \sum_{i=0}^m z_i \nabla f_i(x^*) = 0 \\ \text{(II)} & f_i(x^*) z_i = 0, \quad f_i(x^*) \leq 0, \quad z_i \geq 0 \quad i=1 \dots p \\ \text{(III)} & f_j(x^*) = 0 \quad j=p+1 \dots m \end{cases}$$

Because: (I)

IV
II

$$\phi(x) = z_0 (f_0(x) - f_0(x^*)) + \sum_{i=1}^m z_i f_i(x)$$

is convex & diff & $\phi(x^*) \leq 0$

• but also (*) $\Rightarrow \phi(x) \geq 0$

$\Rightarrow x^*$ gives rise to min. of $\phi(x)$

$$\Rightarrow \nabla_x \phi(x) \Big|_{x=x^*} = 0 \quad (\text{Fenchel's theorem})$$

which implies I)

$$\text{II) (*)} \Rightarrow z_i \geq 0 \quad i=0, \dots, p, \quad f_i(x^*) \leq 0 \quad i=p+1, \dots, m$$

$$\Rightarrow \nexists z_k f_k(x^*) \neq 0 \Rightarrow z_k f_k(x^*) < 0$$

$$\Rightarrow \phi(x^*) < 0 \quad \nexists \Rightarrow z_k f_k(x^*) = 0$$

III) direct for (P).

& (Assato)

By LEM 2^y, the same holds for $y = \frac{z^*}{z_0}$ because $z_0 \neq 0$.

THM: (P) fulfills (Assato) & (Assato). If $x^* \in S$ is optimal and f_i are diff at $x^* \quad i=0, \dots, p$, then $\exists y \in \mathbb{R}^m, y_i \geq 0$:

(KKT conditions)

$$\text{(I)} \quad \sum_{i=1}^m y_i \nabla f_i(x^*) + \nabla f_0(x^*) = 0$$

$$\text{(II)} \quad f_i(x^*) y_i = 0, \quad f_i(x^*) \leq 0, \quad y_i \geq 0 \quad i=p+1, \dots, m$$

$$\text{(III)} \quad f_j(x^*) = 0 \quad j=p+1, \dots, m$$

Example 2 lin. sep. case SVM

$$f_0(w) = \frac{1}{2} \|w\|^2$$

$$f_i(w) := 1 - y^{(i)} (w \cdot x^{(i)})$$

say $w^* \in \mathbb{R}^{n+1}$

We know there is a unique optimal solution for (P)
 In this case $\Rightarrow \exists \alpha \in \mathbb{R}^n, \alpha_i \geq 0$

$$(I) \quad \nabla_w f_0(w^*) + \sum_{i=1}^n \alpha_i \nabla_w f_i(w^*) = 0$$

$$\Rightarrow w^* = \sum_{i=1}^n \alpha_i y^{(i)} x^{(i)}$$

$$(II) \quad \alpha_i f_i(w^*) = 0 \quad i=1 \dots n, \alpha_i \geq 0$$

~~the~~ ~~optimal~~

In conclusion, w^* can be given by training data by (I)
 but only some points actually contribute, those for
 which $f_i(w^*) = 0$ otherwise (II) implies $\alpha_i = 0$!

That is the reason for the name SVM.

Example 1: A) $f_0(x) = (x-5)^2$

$$f_1(x) = 6 - x$$

$$\text{KKT: } f_0'(x) + z_1 f_1'(x) = 0$$

$$\Leftrightarrow 2(x-5) - z_1 = 0$$

$$z_1 \geq 0, f_1(x) z_1 \geq 0, f_1(x) \leq 0$$

$$z_1 > 0 \Rightarrow f_1(x) = 0 \Rightarrow x = 6 \quad \checkmark$$

$$z_1 = 0 \Rightarrow 2(x-5) = 0 \Rightarrow x = 5, f_1(x) = 1 \quad \checkmark$$

B) $f_0(x) = (x-5)^2$

$$f_1(x) = \text{~~6-x~~} 3-x$$

$$\text{KKT: } 2(x-5) - z_1 = 0$$

$$z_1 \geq 0, f_1(x) \leq 0$$

$$z_1 f_1(x) \leq 0$$

$$z_1 > 0: x = 3$$

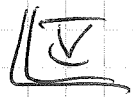
$$\Rightarrow 2(3-5) - z_1 = 0$$

$$\Rightarrow z_1 = -4 \quad \checkmark$$

$$z_1 = 0: x = 5$$



Lagrange function



Given a nice interpretation of the results so far

(3)

Ex: To solve $\inf_{x \in C} \{ f_0(x) \mid f_i(x) \leq 0 \}$ (P)

we may look at an family of auxiliary problems: for $y \geq 0$ and solve

$$\inf_{x \in C} \{ f_0(x) + y f_1(x) \} \quad (P_y)$$

Then y is the "weight" of the constraint:

Let $x^*(y)$ be optimal for (P_y) and y given, then

- $x^*(0)$ may violate $f_1(x^*(0)) \leq 0$
- for y very large one may have $f_1(x^*(0)) < 0$ but $x^*(y)$ may be far from optimal
- there should be an intermediate value \bar{y} s.t. $f_1(x^*(\bar{y})) = 0$
 $\Rightarrow x^*(\bar{y})$ is also a solution of (P).

DEF: i) Let $D = \{ y \in \mathbb{R}^m \mid y_i \geq 0 \ i=1, \dots, m \}$.

We call $L: C \times D \rightarrow \mathbb{R}$ Lagrange function of (P):

$$L(x, y) := f_0(x) + \sum_{i=1}^m y_i f_i(x)$$

ii) We call $(\bar{x}, \bar{y}) \in C \times D$ saddle point of L
iff

$$L(x, \bar{y}) \geq L(\bar{x}, \bar{y}) \geq L(\bar{x}, y) \quad \forall x \in C, y \in D$$

Theorem of KKT:

Thm: Let (P) fulfill (Assex). Then the following holds:

i) (\bar{x}, \bar{y}) saddle point of $L \Rightarrow \bar{x}$ optimal for (P)

and $\bar{y}_i f_i(\bar{x}) = 0$, i.e.,

$$L(\bar{x}, \bar{y}) = f_0(\bar{x})$$

ii) If \bar{x} is optimal for (P) and (Asslab), then

$\exists \bar{y} \in D$ s.t. (\bar{x}, \bar{y}) is saddle point of L .

iii) If $\alpha \in \mathbb{R}$ & (Asslab) $\exists \bar{y} \in D$ s.t.

$$\inf_{x \in C} \{f_0(x) + \sum_{i=1}^m \bar{y}_i f_i(x)\} = \alpha = \inf_{x \in C} L(x, \bar{y}) = \max_{y \in D} \inf_{x \in C} L(x, y).$$

Proof: ~~Let (\bar{x}, \bar{y}) be saddle point of $L(x, y)$,~~

~~then $\forall y \in D: L(\bar{x}, \bar{y}) \leq L(\bar{x}, y)$~~

~~$$L(\bar{x}, \bar{y}) = f_0(\bar{x}) + \sum_{i=1}^m \bar{y}_i f_i(\bar{x}) = f_0(\bar{x}) + \sum_{j=1}^m \bar{y}_j f_j(\bar{x})$$~~

~~\Rightarrow special case holds, $\bar{y} \in D$~~

~~$$\Rightarrow \forall y \in D \quad \sum_{i=1}^m \bar{y}_i f_i(\bar{x}) = \sum_{i=1}^m \bar{y}_i f_i(\bar{x})$$~~

~~$$\bar{y}_i \rightarrow +\infty \Rightarrow f_i(\bar{x}) \leq 0$$~~

~~$$\bar{y}_i \rightarrow 0 \Rightarrow \bar{y}_i f_i(\bar{x}) = 0$$~~

~~and f_j active $\Rightarrow f_j(\bar{x}) = 0 \Rightarrow \bar{x} \in S$~~

But $L(\bar{x}, \bar{y}) \leq L(x, \bar{y})$ for $x \in S$

$$f_0(\bar{x}) \leq f_0(x) + \sum_{i=1}^r \bar{y}_i f_i(x) + \sum_{i=r+1}^m \bar{y}_i f_i(x) \leq f_0(x) \quad \square$$

Proof: i) (\bar{x}, \bar{y}) saddle point of L

(4)

$$\Leftrightarrow \forall x \in C, y \in D:$$

$$L(x, \bar{y}) \stackrel{(1)}{\geq} L(\bar{x}, \bar{y}) \stackrel{(2)}{\geq} L(\bar{x}, y)$$

$$(2) \Leftrightarrow \sum_i \bar{y}_i f_i(x) \geq \sum_i \bar{y}_i f_i(\bar{x})$$

$$y_i \rightarrow \infty \Rightarrow f_i(x) \leq 0, \quad y_i \rightarrow +\infty \Rightarrow f_i(x) \leq 0$$

$$y_i \rightarrow -\infty \Rightarrow f_i(x) \geq 0, \quad y_i \rightarrow -\infty \Rightarrow f_i(x) \geq 0$$

$$y_i \rightarrow 0 \Rightarrow \sum_i \bar{y}_i f_i(x) \geq 0 \Rightarrow \bar{y}_i f_i(x) = 0$$

$$(1) \Rightarrow f_0(x) + \sum_i \bar{y}_i f_i(x) \geq f_0(\bar{x}) + \sum_i \bar{y}_i f_i(\bar{x})$$

$$\Rightarrow f_0(x) \geq f_0(\bar{x})$$

Therefore, \bar{x} is a minimum of f_0 for all $x: f_0(x) \leq f_0(\bar{x}) = 0$.

ii) \bar{x} optimal sol of (P) $\Leftrightarrow \forall x: f_0(\bar{x}) \leq f_0(x)$
 $f_i(\bar{x}) \leq 0 \quad i=1, \dots, p$
 $f_j(\bar{x}) = 0 \quad j=p+1, \dots, m$

$$(A_{convex}) + (A_{slater}) \stackrel{(4.12)}{\Rightarrow} \exists \bar{y} \in D \text{ s.t. } \forall x \in C$$

$$L(x, \bar{y}) = f_0(x) + \sum_i \bar{y}_i f_i(x) + \sum_j \bar{y}_j f_j(x) \geq f_0(x^*)$$

$$\text{for } x = \bar{x} \Rightarrow \sum_i \bar{y}_i f_i(\bar{x}) \geq 0 \text{ but } \bar{y}_i \geq 0, f_i(\bar{x}) \leq 0$$

$$\text{implies } \bar{y}_i f_i(\bar{x}) = 0$$

$$\Rightarrow L(x, \bar{y}) \geq f_0(\bar{x}) = L(\bar{x}, \bar{y}) \geq f_0(\bar{x}) + \sum_i \bar{y}_i f_i(\bar{x}) + \sum_j \bar{y}_j f_j(\bar{x})$$

$$= L(\bar{x}, \bar{y})$$

$$\text{min}) \quad \alpha := \inf_{x \in S} f(x) \in \mathbb{R}$$

$$(A_{\text{convex}}) + (A_{\text{stable}}) \stackrel{\text{LH2}}{\implies} \exists \bar{y} \in D$$

$$L(x, \bar{y}) \geq \alpha$$

$$\implies \inf_{x \in C} L(x, \bar{y}) = \inf_{x \in C} \left[f_0(x) + \sum_i \bar{y}_i f_i(x) + \sum_j \bar{y}_j g_j(x) \right]$$

$$= \inf_{x \in S} L(x, \bar{y}) = \alpha$$

$$\sup_{y \in D} L(x, y) = \begin{cases} f(x) & \text{for } x \in S \\ +\infty & \text{else} \end{cases}$$

$$\implies \inf_{x \in C} \sup_{y \in D} L(x, y) = \inf_{x \in S} f(x) = \alpha$$

$$\implies \alpha = \inf_{x \in S} \sup_{y \in D} L(x, y) \geq \sup_{y \in D} \inf_{x \in C} L(x, y)$$

$$\geq \inf_{x \in C} L(x, \bar{y}) = \alpha$$

$$\implies \alpha = \sup_{y \in D} \inf_{x \in S} L(x, y) = \max_{y \in D} \inf_{x \in S} L(x, y) \quad \square$$

Rem: note that ~~$L(x, y) \geq L(x, \bar{y})$~~ $L(x, \bar{y}) \geq L(x, y)$

$$\text{implies bc. we: } \nabla f_0(x) + \sum_i \bar{y}_i \nabla f_i(x) + \sum_j \bar{y}_j \nabla g_j(x) = 0$$

HW (x, \bar{y}) saddle point of $L \iff$ KKT conditions hold for (x, \bar{y})

PRIMAL & DUAL FORM OF SVM

(I)

Primal program

$$f_0(\omega, \xi) = \frac{1}{2} \|\omega\|^2 + C \sum_{i=1}^n \xi_i \longrightarrow \min$$

subject to constraints

$$f_i(\omega, \xi) = 1 - y_i^T \omega x^{(i)} - \xi_i \leq 0 \quad i=1, \dots, n$$

$$g_j(\omega, \xi) = -\xi_j \leq 0 \quad j=1, \dots, n$$

\Rightarrow Lagrange function

$$L(\omega, \xi, \alpha, \beta) = f_0(\omega, \xi) + \sum_{i=1}^n \alpha_i f_i(\omega, \xi) + \sum_{j=1}^n \beta_j g_j(\omega, \xi)$$

at optimum $\omega, \xi, \alpha, \beta$ we have KKT:

$$(I) \begin{pmatrix} \frac{\partial \omega}{\partial \omega} \\ \frac{\partial \omega}{\partial \xi} \\ \frac{\partial L}{\partial \alpha} \end{pmatrix} \left[f_0(\omega, \xi) + \sum_{i=1}^n \alpha_i f_i(\omega, \xi) + \sum_{j=1}^n \beta_j g_j(\omega, \xi) \right] = 0$$

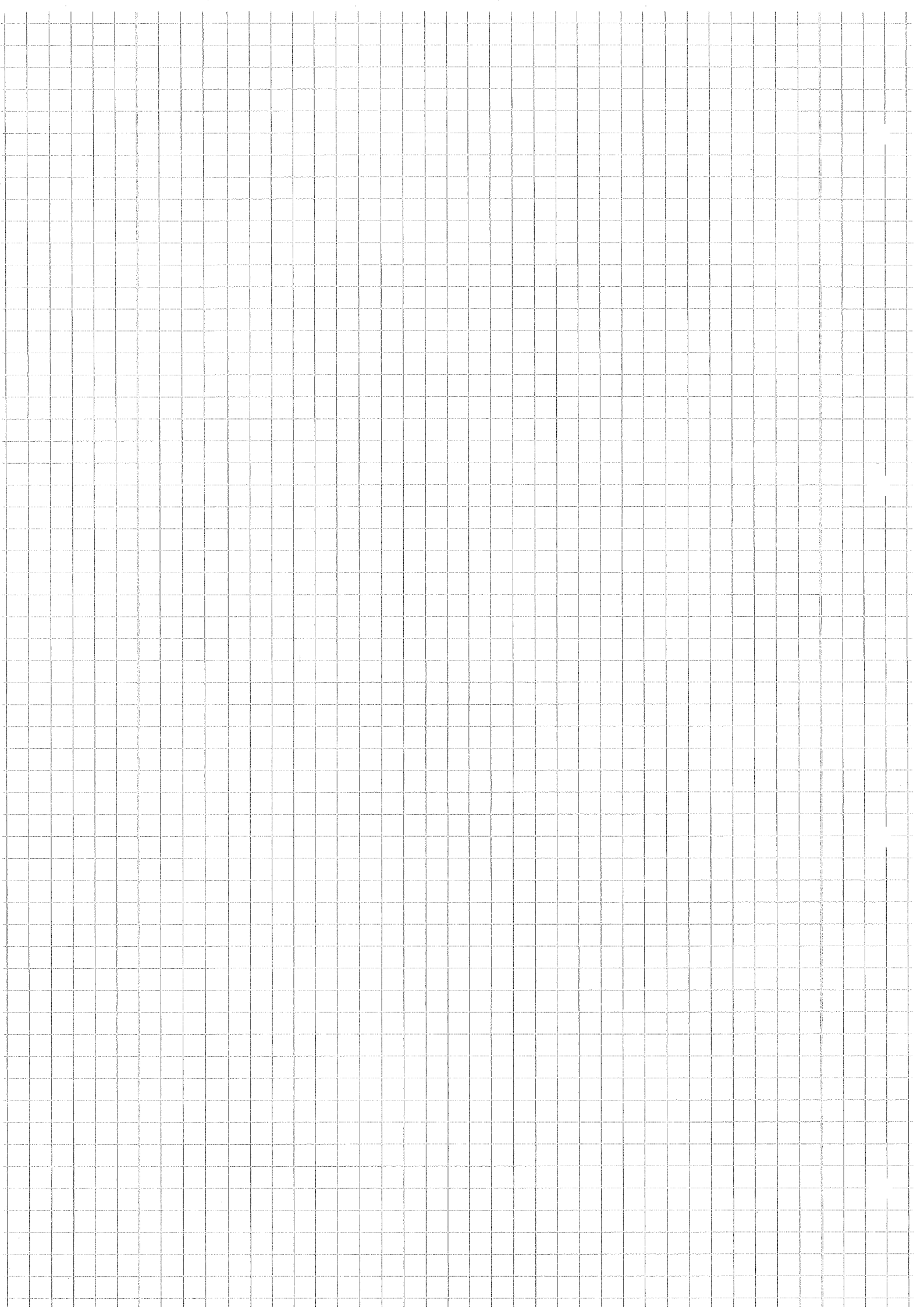
$$(II) \alpha_i f_i(\omega, \xi) = 0, \quad \beta_j g_j(\omega, \xi) = 0, \quad \alpha_i \geq 0, \quad \beta_j \geq 0$$

for $i, j=1, \dots, n$

$$\frac{\partial \omega}{\partial \omega} f_0(\omega, \xi) = 0, \quad \frac{\partial \omega}{\partial \xi} f_0(\omega, \xi) = \omega, \quad \frac{\partial \xi}{\partial \alpha} f_0(\omega, \xi) = C$$

$$\frac{\partial \omega}{\partial \omega} f_i(\omega, \xi) = -y^{(i)}, \quad \frac{\partial \omega}{\partial \xi} f_i(\omega, \xi) = -y^{(i)} x^{(i)}, \quad \frac{\partial \xi}{\partial \alpha} f_i(\omega, \xi) = -x^{(i)}$$

$$\frac{\partial \omega}{\partial \omega} g_j(\omega, \xi) = 0, \quad \frac{\partial \omega}{\partial \xi} g_j(\omega, \xi) = 0, \quad \frac{\partial \xi}{\partial \alpha} g_j(\omega, \xi) = -\xi_j$$



$$(I) \Leftrightarrow \sum_{i=1}^n \alpha_i y_i = 0$$

$$\underline{w} = \sum_{i=1}^n \alpha_i y_i \underline{x}^{(i)}$$

$$C = \alpha_i + \beta_i \quad i=1 \dots n$$

$$(II) \Leftrightarrow \begin{aligned} \alpha_i (1 - y_i) w \cdot x^{(i)} - \xi_i &= 0 & \alpha_i &\geq 0 & i=1 \dots n \\ \beta_j \xi_j &= 0 & \beta_j &\geq 0 & j=1 \dots n \end{aligned}$$

We note again as in the lin. sep. case that due to (I)

$$\underline{w} = \sum_{i=1}^n \alpha_i y_i \underline{x}^{(i)}$$

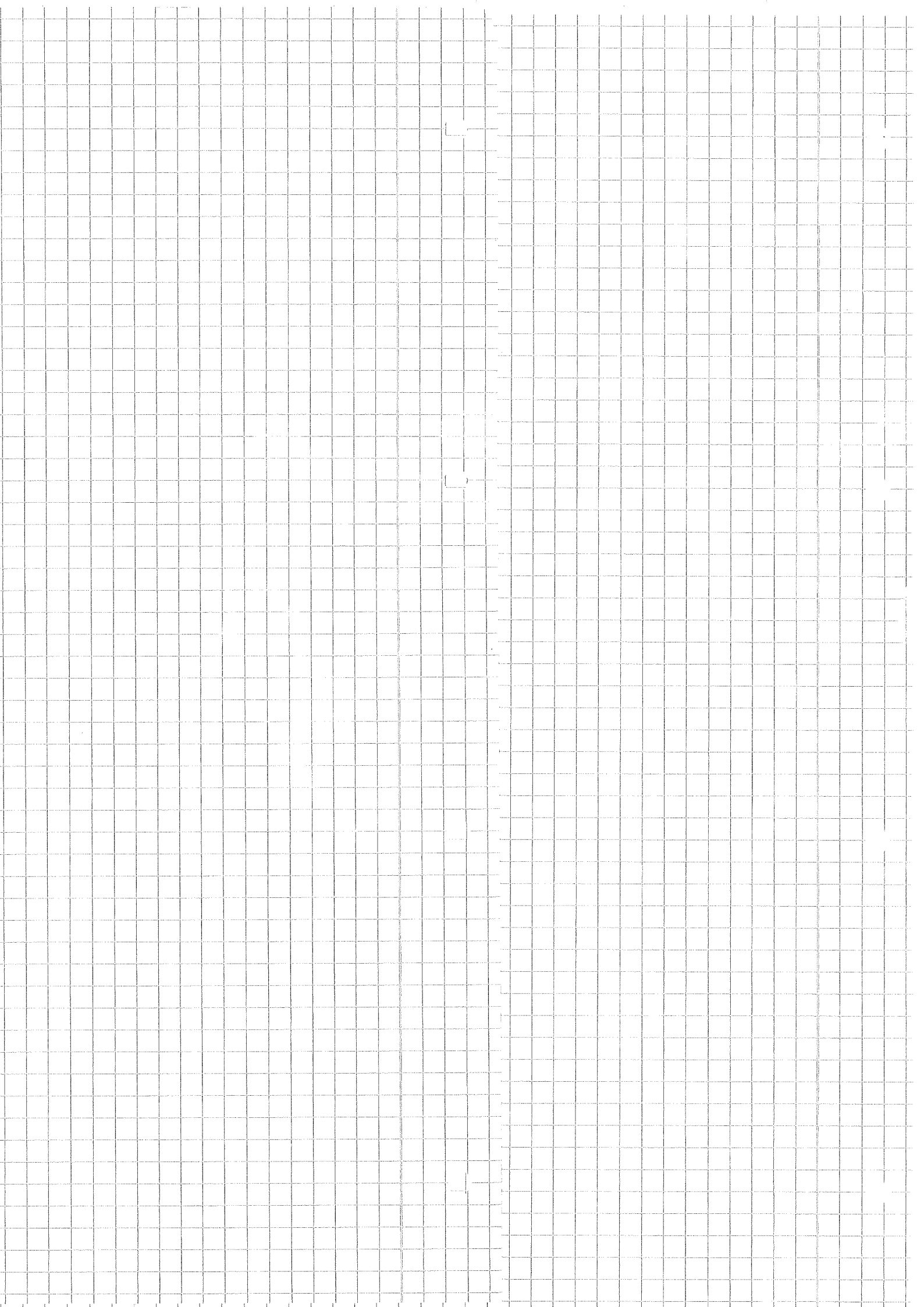
where due to (II) only vectors fulfilling $\alpha_i \neq 0$, i.e.,

$$1 - y_i w \cdot x^{(i)} - \xi_i = 0$$

contribute \Rightarrow either support vectors with $\xi_i = 0$ or other vectors for $\xi_i \neq 0$. The latter however must satisfy $\beta_i \xi_i = 0 \Rightarrow \beta_i = 0 \Rightarrow \alpha_i = C$.

Hence, vectors contributing to \underline{w} are either support vectors or outliers with $\alpha_i = C$.

Note however that even though w is unique as we have shown, the support vectors are not.



NOTE: This program has several advantages:

IV

1) $\frac{d^2}{d\alpha} \tilde{f}_0(\alpha) = -y^{(i)} x^{(i)} \cdot y^{(j)} x^{(j)}$ which is negative semi-def.
 $\Rightarrow \tilde{f}_0(\alpha)$ concave

2) Constraints are all affine

$\Rightarrow \exists!$ optimum and problem is always quadratic

We can furthermore express the activation output

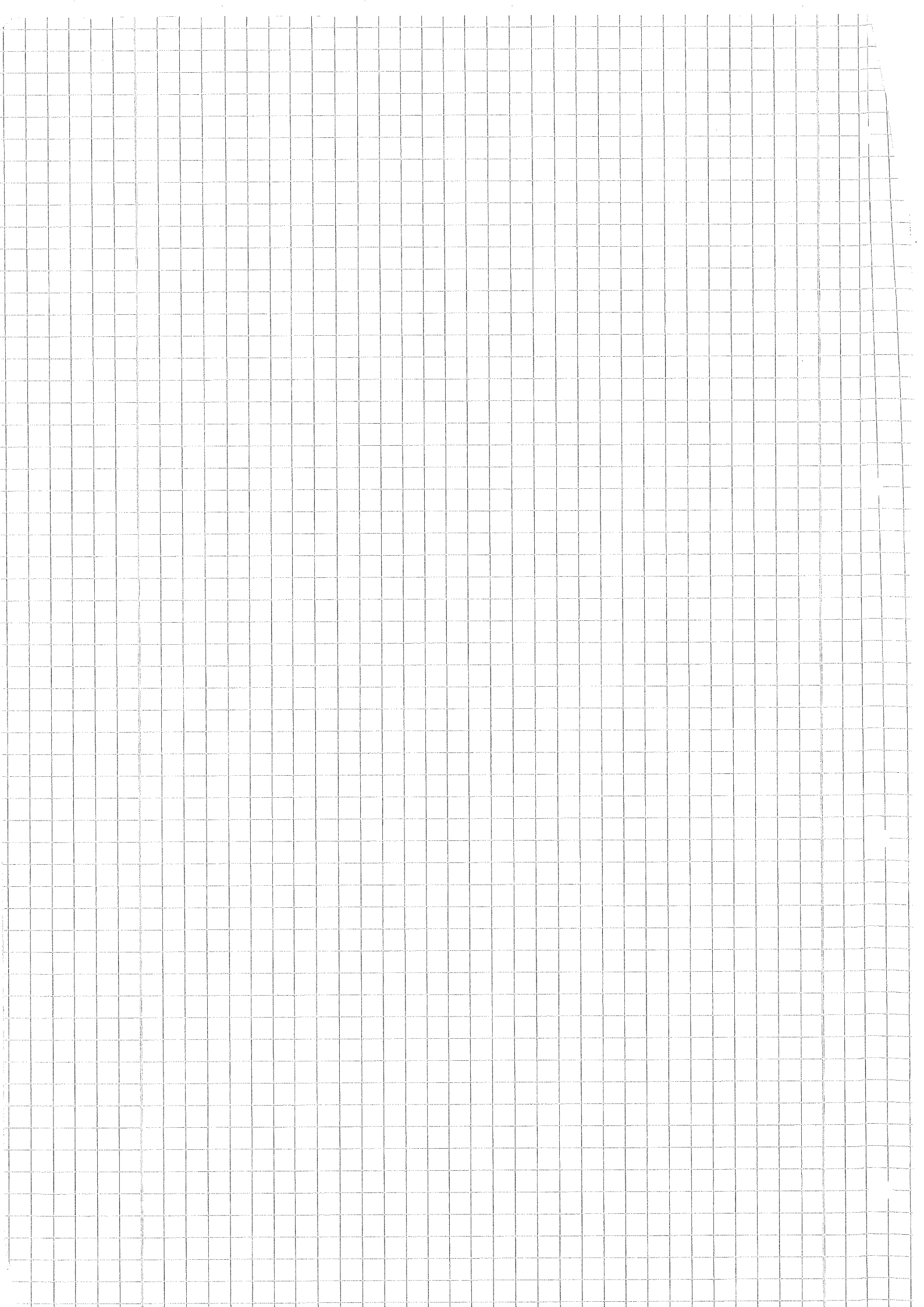
$$h(x) = w_0 x$$

by choosing one support vector $\underline{x}^{(k)}$, i.e., $w_0 \underline{x}^{(k)} = y^{(k)}$

and computing

$$\begin{aligned} w_0 &= \frac{y^{(k)}}{\underline{x}^{(k)}} - \frac{w_0 \cdot \underline{x}^{(k)}}{\underline{x}^{(k)}} \\ &= \frac{y^{(k)}}{\underline{x}^{(k)}} - \sum_{j=1}^r \alpha_j y^{(j)} \frac{\underline{x}^{(j)} \cdot \underline{x}^{(k)}}{\underline{x}^{(k)}} \end{aligned}$$

\Rightarrow The activation does not depend on the particular values of $x^{(i)}$ but only on inner products.



Duality

$$S = \{x \in C \mid \begin{array}{l} f_i(x) \leq 0 \\ f_j(x) = 0 \\ i=1, \dots, p \\ j=p+1, \dots, m \end{array} \} \quad \boxed{I}$$

minimize $f_0(x)$ $x \in C$

subject to $f_i(x) \leq 0$ $i=1, \dots, p$

$f_j(x) = 0$ $j=p+1, \dots, m$

Lagrangian function $L(x, y) = f_0(x) + \sum_{i=1}^p y_i f_i(x) + \sum_{j=p+1}^m y_j f_j(x)$

for $x \in C$, $y \in D = \{\mathbb{R}^m \mid y_i \geq 0 \ i=1, \dots, p\}$

Dual form: $\tilde{L}(y) = \inf_{x \in C} L(x, y)$

Say $\alpha = \inf_{x \in S} f_0(x)$ then

$$\tilde{L}(y) \leq \alpha$$

because

$$\tilde{L}(y) = \inf_{x \in C} L(x, y) \leq L(x^*, y) \leq f_0(x^*)$$

because the constraints are fulfilled

Hence, the optimal value of the dual

problem $S := \sup_{y \in D} \tilde{L}(y)$ satisfies

$$S \leq \alpha$$

$\alpha - S$ is called the duality gap.

Example:

$$\text{max. } f_0(x) = x^2$$

$$\text{sub. } f_1(x) = ax + b \leq 0$$

$$L(x, y) = f_0(x) + y f_1(x), y \geq 0 \quad \partial_x L(x, y) = 0$$

$$\Leftrightarrow \boxed{x = -\frac{ya}{2}}$$

$$\boxed{\tilde{L}(y) = -\frac{1}{4}y^2a^2 + yb} \quad y \geq 0$$

$$\tilde{L}'(y) = 0 \Rightarrow \boxed{y = \frac{2b}{a^2}}$$

$$\Rightarrow \sup_{y \geq 0} \tilde{L}(y) = \begin{cases} \frac{2b}{a^2} & \text{for } b \leq 0 \Rightarrow x = -\frac{b}{a} \\ 0 & \text{for } b > 0 \Rightarrow x = 0 \end{cases}$$

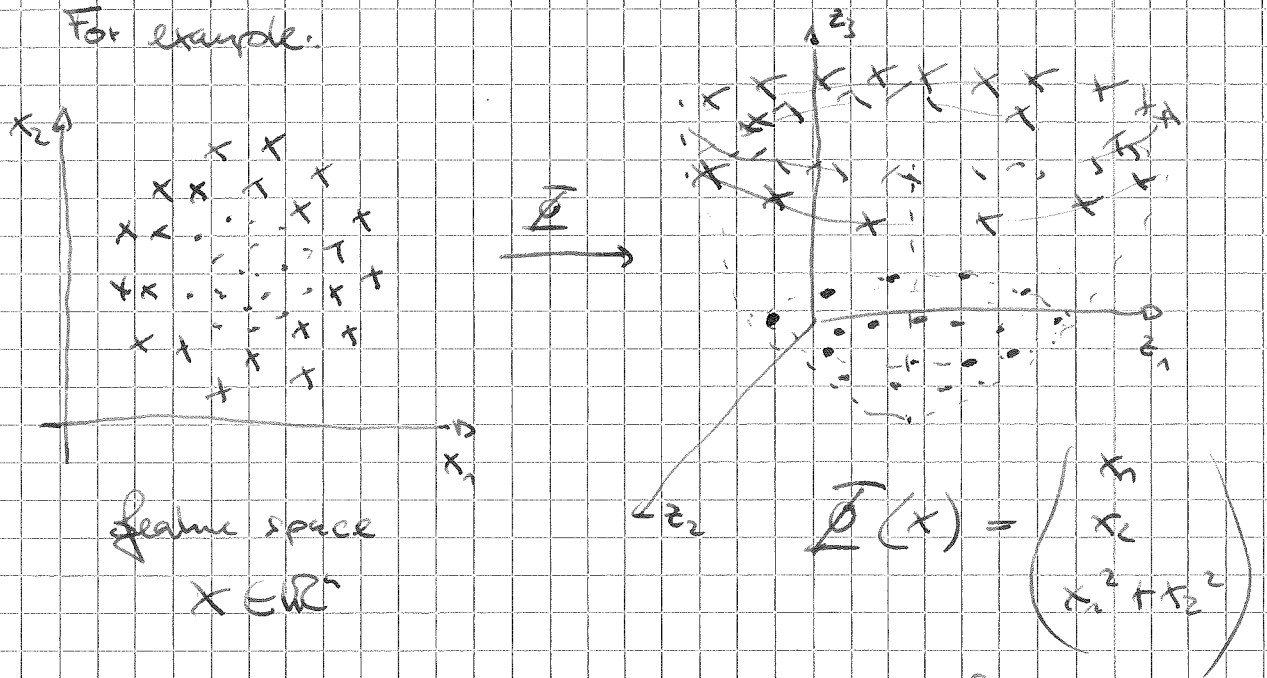
\Rightarrow

NON-LINEAR CLASSIFIERS

(1)

1) General: Oftentimes a good strategy to map the feature space $X \in \mathbb{R}^n$ into a bigger space where the training data become lin. sep.

For example:



We may thus generalize Adaboost as follows

$$h(x) = w \cdot \begin{pmatrix} \Phi(x) \end{pmatrix} = w \cdot \Phi(x)$$

and the loss function becomes

$$L(w) = \frac{1}{2} \sum_{i=1}^M (y^{(i)} - w \cdot \Phi(x)) ^2$$

Learning & update rule is always

$$w \mapsto \tilde{w} = w - \eta \frac{\partial L(w)}{\partial w}$$

(HW:) implement this with test data for non-lin. sep. features

Drawback. This can get very costly. Computation of $\mathcal{E}(x)$ can be $\mathcal{O}(n)$ and n is usually extremely large.

2) For SVMs we have seen that in the dual representation it only depends on inner products

$$\tilde{L}(x) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \langle x^{(i)}, x^{(j)} \rangle \rightarrow \max$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \sum_{i=1}^n \alpha_i y_i = 0$$

Using the map \mathcal{E} we need to compare

$$\langle \mathcal{E}(x), \mathcal{E}(x') \rangle =: K(x, x')$$

K can be chosen freely as long as the kernel of \mathcal{E} is guaranteed and, in order to assure concavity, K must be pos. def. symmetric:

DEF: $K: X \times X \rightarrow \mathbb{R}$ is pos. def. sym.

$$\Leftrightarrow \forall x_1, \dots, x_n \in X, [K(x_i, x_j)]_{i,j}$$

is a pos. def. symmetric matrix in $\mathbb{R}^{n \times n}$

THM: (Mercer's cond.) $K: X \times X \rightarrow \mathbb{R}$ cont.

and symmetric, $X \subseteq \mathbb{R}^n$ compact, then

$\exists a_n > 0, \phi_n: X \rightarrow \mathbb{R}$ cont s.t.

$$K(x, x') = \sum_{n=0}^{\infty} a_n \phi_n(x) \phi_n(x')$$

iff K is the kernel of a pos. def.

Hermitian-Schmidt operator i.e.,

$$\forall f \in L^2(X): \langle f, Kf \rangle \geq 0.$$

2) For SVM we have seen that in the dual representation the optimization problem depends entirely on inner products:

$$\text{maximize } \tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y^{(i)} y^{(j)} \langle \underline{x}^{(i)}, \underline{x}^{(j)} \rangle$$

$$\text{subject to } 0 \leq \alpha_i \leq C, \quad \sum_{j=1}^n \alpha_j y^{(j)} = 0, \quad j=1..n$$

Applying a map from feature space to some
 $\Phi: X \rightarrow \mathcal{H}$ Hilbert space \mathcal{H} with $\dim \mathcal{H} \gg n$

the scalar products turn into

$$\langle \Phi(\underline{x}^{(i)}), \Phi(\underline{x}^{(j)}) \rangle_{\mathcal{H}} = K(\underline{x}^{(i)}, \underline{x}^{(j)})$$

But since all depends on K and not directly on $\Phi(\cdot)$, there is no need to specify Φ , we only have to make sure K carries the properties that induces a scalar product.

~~Why do the SVM with kernel for soft margin work with typhoon kernel?~~

Advantages: While $\langle \Phi(x), \Phi(s) \rangle$ usually take $O(\dim \mathcal{H})$ to compute $K(x,s)$ may take much less, e.g., $O(n)$.

Thm: (Mercer's cond.)

Let $X \subset \mathbb{R}^n$ be compact and $K: X \times X \rightarrow \mathbb{R}$ is a continuous and symmetric function. Then

$$K(x, y) = \sum_{n=0}^{\infty} \lambda_n \phi_n(x) \phi_n(y)$$

~~$K(x, y) = \sum_{n=0}^{\infty} \lambda_n \phi_n(x) \phi_n(y)$~~ $T_K f := \int_X K(x, z) f(z) dz \quad \forall f \in L^2(X)$

is self-adjoint and there are eigenvalues/values $\lambda_1, \lambda_2, \dots / \phi_1, \phi_2, \dots$ ordered by their geometric multiplicity.

Furthermore, (semi-)positivity implies expansion and vice versa)

$$\forall f \in L^2: \langle f, T_K f \rangle_X \geq 0 \iff K(x, y) = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \phi_j(y)$$

~~$\langle f, T_K f \rangle_X \geq 0$~~ $\lambda_j \geq 0$

where the convergence is absolute and uniform.

RFN: Let $K(x, y)$ fulfill Mercer's cond. then the eigenvalue/vectors def. a map

$$\underline{\Phi}(x) = \left(\sqrt{\lambda_i} \phi_i(x) \right)_{i \in \mathbb{N}}$$

and a Hilbert space \mathcal{H} with a scalar product

$$\langle \underline{\Phi}(x), \underline{\Phi}(y) \rangle_{\mathcal{H}} = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \phi_j(y)$$

DEF: (Pos. def. sym. kernels) A $K: X \times X \rightarrow \mathbb{R}$ is called pos. def. sym. if $\forall \{x_1, \dots, x_n\} \subset X$ the matrix $(K(x_i, x_j))_{i, j \in \{1, \dots, n\}}$ is symmetric and pos. semidef.

Recall: A matrix $M \in \mathbb{R}^{n \times n}$ sym. and pos. semi-def.

iff it \Rightarrow eigenvalues are non-neg.

i) all eigenvalues are non-neg.

ii) for any $v \in \mathbb{R}^n$: $\langle v, Mv \rangle \geq 0$

Example: (Polynomial Kernel)

$$\forall x, y \in \mathbb{R}^n: K(x, y) = (x \cdot y + c)^d$$

for any $c \in \mathbb{R}, d \in \mathbb{N}$

For $n=2, d=2$ we find

$$\Phi(x) = \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \\ \sqrt{2c}x_1 \\ \sqrt{2c}x_2 \\ c \end{pmatrix} \Rightarrow \dim(\Phi) = 6$$

HW Show that $\dim(\Phi) = \binom{N+d}{d} = \frac{(N+d)!}{N!d!}$

~~(Gaussian Kernel)~~

~~$\forall x, y \in \mathbb{R}^n: K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$~~

~~for any $\sigma > 0$~~

Gaussian Kernel

$$K(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right) \quad \sigma > 0$$

Why pos. def. sym.

LEM Pos. def. sym. kernels are closed under
 sum, product, tensor product, point-wise limit
 and composition with power series $\sum a_n x^n$
 for $a_n \geq 0$.

Proof: sum) clear

product) as $(K_{ij})_{1 \leq i, j \leq n} = (K(x_i, x_j))_{1 \leq i, j \leq n}$ is pos. def. sym.

$$\exists R \text{ s.t. } K = R^* R$$

the kernel matrix associated with $K(x, y) = K(x, y)$

$\Rightarrow (K_{ij} K'_{ij})_{1 \leq i, j \leq n}$ hence for any $v \in \mathbb{R}^n$

$$\sum_{ij} K_{ij} K'_{ij} v_i v_j = \sum_{ij} \sum_k R_{ik} R_{jk} K'_{ij} v_i v_j$$

$$= \sum_k w_k K' w_k \geq 0 \quad \text{for } w_k = \text{row} \begin{pmatrix} v_1 R_{1k} \\ \vdots \\ v_n R_{nk} \end{pmatrix}$$

and furthermore sym.

tensor product) $\tilde{K}: (x_1, x'_1, x_2, x'_2) \mapsto K(x_1, x_2)$

$\tilde{K}': (x_1, x'_1, x_1, x'_1) \mapsto K'(x_1, x_1)$

$\Rightarrow K \otimes K' (x_1, x'_1, x_2, x'_2) = \tilde{K} \tilde{K}'$ pos. def. sym.
 as product

point-wise limit) say $(K_n)_{n \in \mathbb{N}}$ pos. def. sym.
 kernels and $K_n(x, y) \rightarrow K(x, y)$ for
 for all $v \in \mathbb{R}^n$

$$\langle v, K_n v \rangle \geq 0 \Rightarrow \lim_{n \rightarrow \infty} \langle v, K_n v \rangle \geq 0 \text{ and sym.}$$

power series) $\sum a_n x^n$ converges with
 conv. radius R . Suppose $|K(x,y)| < R$
 $\forall x,y \in X$, then $a_n K^n$ pos. def. sym. by product)
 and $\sum a_n K^n$ pos. def. sym. by limit)

HW:

~~For any pos. def. sym. kernel~~
 $K'(x,y) = \begin{cases} 0 & \text{for } K(x,x) = 0 = K(y,y) \\ K(x,y) & \text{else} \end{cases}$
~~is also a kernel pos. def. sym.~~



HW:

Show that $K(x,y) = \tanh(a(x+y) + b)$
 is a pos. def. sym. kernel for $a \geq 0, b \geq 0$.

Thm:

(Reproducing kernel Hilbert space)

Let $K: X \times X \rightarrow \mathbb{R}$ be a pos. def. sym. kernel
 of real-valued functions
 Then: i) \exists a Hilbert space \mathcal{H} and a map
 $\Phi: X \rightarrow \mathcal{H}$ s.t.

$$K(x,y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{H}}$$

ii) $\forall f \in \mathcal{H}, x \in X:$

$$f(x) = \langle f, K(x, \cdot) \rangle_{\mathcal{H}}$$

Proof: For any $x, y \in X$ def.

$$\bar{\Phi}(x)(y) = k(x, y)$$

$$\text{Def. } \mathcal{R}^0 := \left\{ \sum_{i \in I} a_i \bar{\Phi}_i(x_i) \mid a_i \in \mathbb{R}, x_i \in X, |I| < \infty \right\}$$

Def: $\langle \cdot, \cdot \rangle : \mathcal{D}^0 \times \mathcal{R}^0 \rightarrow \mathbb{R}$ st for $f, g \in \mathcal{R}^0$,
i.e.,

$$f = \sum_{i \in I} a_i \bar{\Phi}_i(x_i)$$

$$g = \sum_{j \in J} b_j \bar{\Phi}_j(y_j),$$

$$\begin{aligned} \langle f, g \rangle &:= \sum_{\substack{i \in I \\ j \in J}} a_i b_j \underbrace{\bar{\Phi}_i(x_i) \bar{\Phi}_j(y_j)}_{k(x_i, y_j)} = \left[\sum_{i \in I} a_i g(x_i) \right] \\ &= \sum_{j \in J} b_j f(y_j) \end{aligned}$$

Note: i) $\langle f, g \rangle = \langle g, f \rangle$

ii) Representer indep. \Rightarrow well-def.

$$\begin{aligned} \text{iii) } \langle f, f \rangle &= \sum_{i, j \in I} a_i a_j k(x_i, x_j) \geq 0 \\ &= \langle \cdot, \cdot \rangle \text{ pos. semi-def. bilinear form} \end{aligned}$$

$\Rightarrow \langle \cdot, \cdot \rangle$ pos. semi-def. on \mathcal{D}^0 but $\forall f \in \mathcal{D}^0$

$$\begin{aligned} \langle f, \bar{\Phi}(x) \rangle^2 &\leq \langle f, f \rangle \langle \bar{\Phi}(x), \bar{\Phi}(x) \rangle \\ &= \langle f, f \rangle \underbrace{k(x, x)}_{\geq 0} \end{aligned}$$

$$\text{and } f(x) = \sum_{i \in I} a_i k(x_i, x) = \langle f, \bar{\Phi}(x) \rangle$$

$$\Rightarrow |f(x)|^2 \leq \langle f, f \rangle \langle \bar{\Phi}(x), \bar{\Phi}(x) \rangle$$

and therefore

$$f=0 \Leftrightarrow \langle f, f \rangle = 0$$

Given $C: \mathcal{X} \rightarrow \mathbb{R}$ definite \Rightarrow inner product on \mathcal{X}^0

Def $\mathcal{H} := \overline{\mathcal{X}^0} \subset \mathcal{X}$

Furthermore, $L_x: f \mapsto \langle f, \Phi(x) \rangle$ is bounded, hence, reproducing property holds on \mathcal{H} . \square

Kernel Normalization (*)

Conclusion: SVM for non-linear problems

1) $K(x, y)$ pos. def. sym.

2) maximize $\tilde{L}(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{j=1}^n \alpha_j \sum_{k=1}^n \alpha_k y_j^{(1)} y_k^{(1)} K(x_j^{(1)}, x_k^{(1)})$

subject to $0 \leq \alpha_i \leq C, \sum_i \alpha_i y_i^{(1)} = 0$

3) activation function (hypothesis)

$$h(x) = \sigma \left(\sum_i \alpha_i K(x_i^{(1)}, x) + w_0 \right)$$

$$w_0 = y^{(1)} - \sum_j \alpha_j y_j^{(1)} K(y_j^{(1)}, x_j^{(1)})$$

for any $x^{(1)}$ with $\alpha_i \in [0, C]$

Property 3) is quite general.

Module effect the hypothesis function is a lin. comb of $k(x_i, \cdot)$:

Theorem: (Reproducing Theorem)

Let $k: X \times X \rightarrow \mathbb{R}$ pos. def. sym.,

$G: \mathbb{R} \rightarrow \mathbb{R}$ non-dec, $L: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$

$$h^* = \underset{h \in \mathcal{H}}{\operatorname{argmin}} \left[\underbrace{G(\|h\|_X) + L(h(x_i), i=1..n)}_{F(h)} \right]$$

admits sol. of the form

$$h^* = \sum_{i=1}^n a_i k(x_i, \cdot)$$

and if G increasing only such one.

Proof: def. $\mathcal{X}_n := \operatorname{span} \langle k(x_i, \cdot), i=1..n \rangle$

and split $\mathcal{H} = \mathcal{X}_n \oplus \mathcal{X}_n^\perp$ s.t. $\forall h \in \mathcal{H}$

$\exists h_n \in \mathcal{X}_n, h_n^\perp \in \mathcal{X}_n^\perp$ with $h = h_n + h_n^\perp$

Clearly $G(\|h_n\|_X) \leq G(\|h\|_X)$

and by reproducing property

$$h(x_i) = \langle h, k(x_i, \cdot) \rangle = h_n(x_i)$$

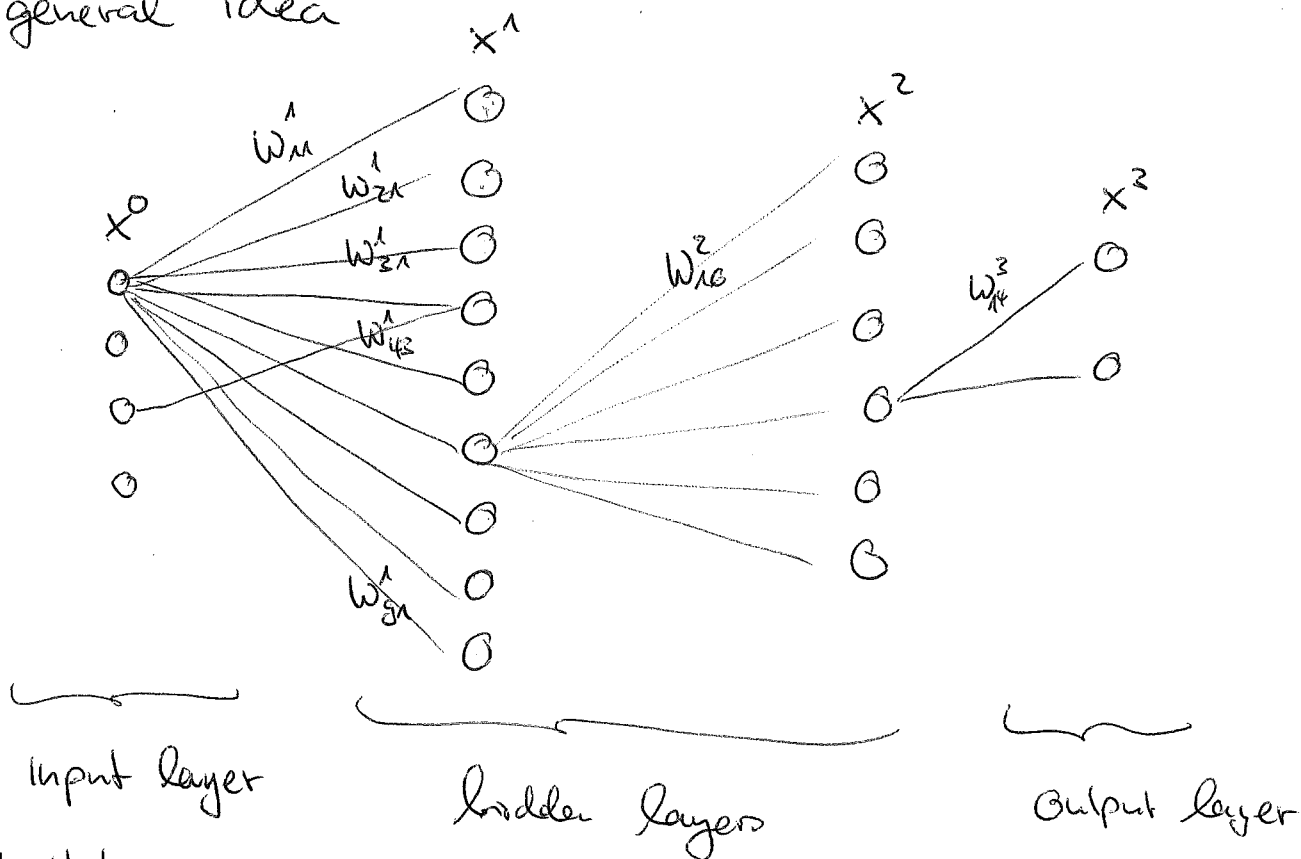
$$\Rightarrow L(h(x_i), i=1..n) = L(h_n(x_i), i=1..n)$$

$$\Rightarrow F(h_n) \leq F(h)$$

The last assertion follows from the corresp. strict neg. \square

MULTILAYER NETWORKS

general idea



network's forward pass:

$$\mathbb{R}^{n^0} \ni x^0 \mapsto \underbrace{W^1 \cdot x^1 + b^1}_{\mathbb{R}^{n^1}} = z^1$$

$$x^1 := \sigma^1(z^1) \mapsto \underbrace{W^2 \cdot x^1 + b^2}_{\mathbb{R}^{n^2}} = z^2$$

$$x^2 := \sigma^2(z^2) \mapsto \underbrace{W^3 \cdot x^2 + b^3}_{\mathbb{R}^{n^3}} = z^3$$

$$x^3 := \sigma^3(z^3)$$

- each layer takes as input x^{k-1} and applies

the maps $x^{k-1} \mapsto z^k := W^k x^{k-1} + b^k \rightarrow \sigma^k(z^k)$

where $W^k \in \mathbb{R}^{n^k \times n^{k-1}}$, $b^k \in \mathbb{R}^{n^k}$, $\sigma^k: \mathbb{R}^{n^k} \rightarrow \mathbb{R}^{n^k}$

$$(z_i) \mapsto (\sigma^k(z_i))_{1 \leq i \leq n^k}$$

for some $\sigma^k: \mathbb{R} \rightarrow \mathbb{R}$

• the output of an N -layer neural network is

$$x^N = z(z^N) = z(W^N x^{N-1} + b^N) = \dots = x^N(x^0)$$

this is what we usually called the hypothesis

function $h(x^0) \equiv x^N(x^0)$

• given the training data $(x^{(i)}, y^{(i)})_{i=1 \dots M}$
with $x^{(i)} \in \mathbb{R}^{n^0}$ being the features and
 $y^{(i)} \in \mathbb{R}^{n^N}$ the "labels".

• for classification we use the convention

class $i \leftrightarrow \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{n^N}$ i-th position

• for regression $y^{(i)} \in \mathbb{R}^{n^N}$ are the
 n^N continuous values of the output
that is desired

After the forward pass we want to compare
the output of the network to the labels of
our training data and update the weights ^{we} and
biases b^l accordingly:

$$L = \frac{1}{M} \sum_{i=1}^M l_i \quad \text{as loss function}$$

for some function l_i depending on the i -th pair $(x^{(i)}, y^{(i)})$
for example $l_i = \frac{1}{2} (y^{(i)} - x^N(x^{(i)}))^2$

How should the weights and biases be updated? (2)

Gradient descent: let p^k be a placeholder for some parameter in the k -th layer, i.e., W_{ij}^k, b_i^k for $1 \leq i \leq n^k, 1 \leq j \leq n^{k-1}$

$$\frac{\partial L}{\partial p^k} = \frac{1}{M} \sum_{i=1}^M \frac{\partial \ell_i}{\partial p^k} \quad \text{recall } \ell_i = \ell(y^{(i)}, x^{(i)})$$

$$\frac{\partial \ell}{\partial p^k} = \sum_{i=1}^{n^N} \underbrace{\frac{\partial}{\partial x_i^N} \ell(y, x^N)}_{=: \Delta(x^N)} \cdot \frac{\partial x_i^N}{\partial p^k}$$

depends on all the weights W^l and biases b^l

$$\frac{\partial x_i^N}{\partial p^k} = \frac{\partial [\sigma^N(z_i^N)]}{\partial p^k} = \sum_{i=1}^{n^N} \sigma^{N'}(z_i^N) \frac{\partial z_i^N}{\partial p^k}$$

$$= \sum_{i=1}^{n^N} \sigma^{N'}(z_i^N) \frac{\partial \left[\sum_{j=1}^{n^{N-1}} W_{ij}^N x_j^{N-1} + b_i^N \right]}{\partial p^k}$$

$$= \sum_{i=1}^{n^N} \sigma^{N'}(z_i^N) \left[\sum_{j=1}^{n^{N-1}} W_{ij}^N \frac{\partial x_j^{N-1}}{\partial p^k} + \sum_{j=1}^{n^{N-1}} \frac{\partial W_{ij}^N}{\partial p^k} x_j^{N-1} + \frac{\partial b_i^N}{\partial p^k} \right]$$

We note that

$$\frac{\partial x^l}{\partial p^k} = 0 \quad \text{for } l \leq k$$

$$\frac{\partial W^l}{\partial p^k} = 0 = \frac{\partial b^l}{\partial p^k} \quad \text{for } l \neq k$$

Hence, the last expression simplifies depending on layer number k :

if $k \leq N-1$ this has to be expanded

$$\frac{\partial x^N}{\partial p^k} = \sigma^{N'}(z^N) \odot W^N \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$+ \sigma^{N'}(z^N) \odot \left[\frac{\partial W^N}{\partial p^k} x^{N-1} + \frac{\partial b^N}{\partial p^k} \right]$$

if $k \neq N$ this term is zero

for $k < N$:

$$\frac{\partial x^N}{\partial p^k} = \sigma^{N'}(z^N) \odot W^N \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$= \sigma^{N'}(z^N) \odot W^N \cdot \left[\sigma^{N-1'}(z^{N-1}) \odot W^{N-1} \cdot \frac{\partial x^{N-1}}{\partial p^k} \right.$$

$$\left. + \sigma^{N-1'}(z^{N-1}) \odot \left(\frac{\partial W^{N-1}}{\partial p^k} \cdot x^{N-2} + \frac{\partial b^{N-1}}{\partial p^k} \right) \right]$$

and, hence, we get the following structure:

$$\frac{\partial L}{\partial p^k} = \Delta(x^N) \cdot \left[\sigma^{N'}(z^N) \odot W^N \right] \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$= \Delta(x^N) \cdot \left[\sigma^{N'}(z^N) \odot W^N \right] \cdot \left[\sigma^{N-1'}(z^{N-1}) \odot W^{N-1} \right]$$

$$\cdot \left[\dots \right] \cdot \left[\sigma^{k+1'}(z^{k+1}) \odot W^{k+1} \right]$$

$$\cdot \left[\sigma^{k'}(z^k) \odot \left(\frac{\partial W^k}{\partial p^k} \cdot x^{k-1} + \frac{\partial b^k}{\partial p^k} \right) \right]$$

An efficient update rule: Backpropagation

③

Algorithm: 1) Input $x^0 \leftarrow x^{(i)}$

2) Forward pass: $x^0 \mapsto z^1 \mapsto x^1 \mapsto \dots \mapsto z^N \mapsto x^N$

3) Compute $\Delta(x^N) = l(y^{(i)}, x^N(x^{(i)}))$

4) Compute backward pass

$$\Delta^N := \Delta(x^N)$$

$$\Delta^{k-1} := \Delta^k \cdot [z^{k'}(z^k) \odot W^k]$$

5) Compute

$$\frac{\partial l_i}{\partial p^k} = \Delta^k \frac{\partial x^k}{\partial p^k}$$

$$= \Delta^k \left[z^{k'}(z^k) \odot \left(\frac{\partial W^k}{\partial p^k} \cdot x^{k-1} + \frac{\partial b^k}{\partial p^k} \right) \right]$$

6) Compute l_i for $i \in I$ (mini-batches $I \subseteq \{1, \dots, n\}$)
and average

$$\frac{\partial L_I}{\partial p^k} = \frac{1}{|I|} \sum_{i \in I} \frac{\partial l_i}{\partial p^k}$$

7) Update weights accordingly

$$W_{ij}^k \mapsto W_{ij}^k - \eta \frac{\partial L_I}{\partial W_{ij}^k}$$

$$b_i^k \mapsto b_i^k - \eta \frac{\partial L_I}{\partial b_i^k}$$

8) repeat for all mini-batches and epochs.

HW Derive this backpropagation algorithm from the KKT condition for the Lagrange formula of:

$$\min_{\mathbf{x}} \frac{1}{N} \sum_{i=1}^M \ell(y_i^0, x^N)$$

Under constraints $x^k = z(w^k x^{k-1} + b^k)$

1) formula Lagrange

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^M \ell(y_i^0, x^N) + \sum_{k=1}^M \sum_{i=1}^N \beta_i^k \cdot [x^k - z(w^k x^{k-1} + b^k)]$$

2) $\frac{\partial \mathcal{L}}{\partial \beta_i}$ gives forward pass

3) $\frac{\partial \mathcal{L}}{\partial x_i}$ gives backward pass

4) $\frac{\partial \mathcal{L}}{\partial w^k}, \frac{\partial \mathcal{L}}{\partial b^k}$ give update rule

Python implementation

- MNIST data download
- Number browser
- example :
 - randomly initialized weights
 - $h = 3, n^0 = 28 \times 28 = 784,$
 - $n^1 = 25, n^2 = 10, \text{quadratic loss}$

\Rightarrow efficiency $\approx 34\%$

BACK PROPAGATION



Network forward pass

$$x^{(i)} \xrightarrow{\text{id}} x^0 \in \mathbb{R}^{n^0}$$



$$z^1 = W^1 x^0 + b^1 \in \mathbb{R}^{n^1}, \text{ i.e., } W^k \in \mathbb{R}^{n^k \times n^{k-1}}, b^k \in \mathbb{R}^{n^k}$$



$$x^1 = \sigma^1(z^1), \text{ i.e., } \sigma^k: \mathbb{R}^k \rightarrow \mathbb{R}^k$$

$$(z_i)_i \mapsto (\sigma^k(z_i))_i$$



where $\tilde{\sigma}^k: \mathbb{R} \rightarrow \mathbb{R}$

⋮

$$x^{k-1} = \sigma^{k-1}(z^{k-1})$$



k -th layer

$$\left\{ \begin{array}{l} z^k = W^k z^{k-1} + b^k \\ \downarrow \\ x^k = \sigma^k(z^k) \end{array} \right.$$



⋮



$$h(x^{(i)}) \xrightarrow{\text{id}} x^N = \sigma(z^N)$$

Compare \leftarrow

$\sigma(y^{(i)})$, where

$$\sigma: \{0, 1, \dots, M\} \rightarrow \mathbb{R}^{n^N}$$
$$i \mapsto \begin{pmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{pmatrix} - \text{ith}$$

Loss
function

$$L = \frac{1}{M} \sum_{i=1}^M \ell(v(y^{(i)}), h(x^{(i)}))$$

where $\ell(v, x)$ is a "distance" or a function of it, e.g.

$$\ell(v, x) = \frac{1}{2} (v - x)^2$$

Optimization $\circ = \frac{\partial L}{\partial p^k}$ where p^k stands for $\begin{pmatrix} W_{ij}^k \\ b_i^k \end{pmatrix}$

$$\frac{\partial L}{\partial p^k} = \frac{1}{M} \sum_i \frac{\partial \ell^{(i)}}{\partial p^k}$$

$$\frac{\partial \ell}{\partial p^k} = \left\{ \ell'_i(v(y^{(i)}), h(x^{(i)})) \frac{\partial h_j}{\partial p^k} \right\} = \overbrace{\ell'(v(y^{(i)}), x^{(i)})}^{\Delta(x^{(i)})} \cdot \frac{\partial x^N}{\partial p^k}$$

~~$$\frac{\partial \ell}{\partial p^k} = \sum_i \delta^{(i)}(z_i^N) \frac{\partial z_i^N}{\partial p^k}$$~~

$$\frac{\partial x^N}{\partial p^k} = \sum_i \delta^{(i)}(z_i^N) \frac{\partial z_i^N}{\partial p^k}$$

$$= \sum_i \delta^{(i)}(z_i^N) \left[\sum_j \frac{\partial W_{ij}^N}{\partial p^k} x_j^{N-1} + \sum_j W_{ij}^N \frac{\partial x_j^{N-1}}{\partial p^k} + \frac{\partial b_i^N}{\partial p^k} \right]$$

$$= \delta^{(i)}(z_i^N) \odot W^N \cdot \frac{\partial x^{N-1}}{\partial p^k}$$

$$+ \delta^{(i)}(z_i^N) \odot W^N \cdot \left[\frac{\partial W^N}{\partial p^k} \cdot x^{N-1} + \frac{\partial b^{N-1}}{\partial p^k} \right]$$

~~Again, the same~~

$$\frac{\partial x^e}{\partial p^k} = 0 \text{ for } e \neq k$$

$$\frac{\partial W^e}{\partial p^k} = 0 \text{ for } e \neq k$$

$$\frac{\partial b^e}{\partial p^k} = 0$$

Thus, we get the following structure by induction

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial p^k} &= \Delta(x^N) \cdot [z^{N'}(z^N) \odot W^N] \cdot \frac{\partial x^{N-1}}{\partial p^k} \\ &= \Delta(x^N) \cdot [z^{N'}(z^N) \odot W^N] \\ &\quad \cdot [z^{(N-1)'}(z^{N-1}) \odot W^{N-1}] \\ &\quad \cdot \dots \\ &\quad \cdot [z^{(k+1)'}(z^{k+1}) \odot W^{k+1}] \\ &\quad \cdot [z^{k'}(z^k) \odot \left(\frac{dw^k}{dp^k} x^{k-1} + \frac{db^k}{dp^k} \right)] \end{aligned}$$

Backpropagation algorithm

① Compute forward pass

$$\begin{aligned} x^{(1)} = x^0 &\mapsto z^1 = W^1 x^0 + b^1 \mapsto x^1 = z^1 \\ &\mapsto \dots \mapsto x^{k-1} \mapsto z^k = W^k x^{k-1} + b^k \\ &\mapsto x^k = z^k \mapsto \dots \mapsto x^N = z^N \end{aligned}$$

② Compute $\Delta(x^1) = \ell'(v(y^1), h(x^1))$

③ Compute backward pass

$$\Delta^N = \Delta(x^1)$$

$$\Delta^{k+1} = \Delta^k \cdot [z^{(k+1)'}(z^{k+1}) \odot W^{k+1}]$$

and with that

$$\frac{dl}{dp^k} = \Delta^{k\alpha} \cdot \frac{dx^\alpha}{dp^k}$$

$$= \Delta^k \cdot \left[z^{k'}(z^k) \odot \left(\frac{\partial w^k}{\partial p^k} \cdot x^{k'} + \frac{\partial b^k}{\partial p^k} \right) \right]$$

LEARNING BEHAVIOR

1) Consider one neuron with quadratic loss

$$\mathbb{R} \ni x^0 \rightarrow \text{neuron} \rightarrow \text{output} \quad x^1 = z(z^1) = z(Wx^0 + b^1) \in \mathbb{R}$$

dropping the indices

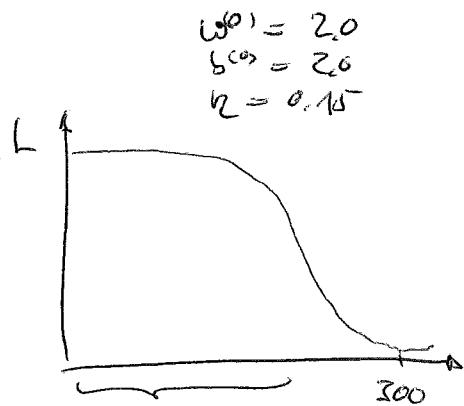
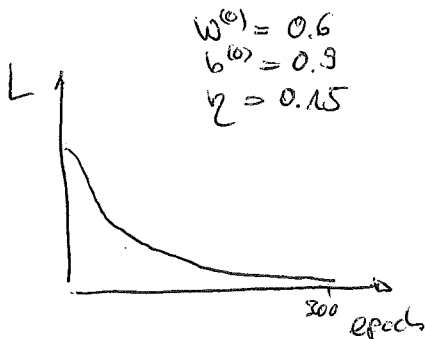
$$h(x) = z(Wx + b) \quad \text{for } W, b \in \mathbb{R}$$

Cost function for one training data point $(x^{(1)}, y^{(1)}) = (1, 0)$

$$L(W, b) = \frac{1}{2} (0 - h(1))^2 = \frac{1}{2} h(1)^2$$

$$\frac{\partial L}{\partial W} = h(1) \frac{\partial}{\partial W} z(Wx + b) = h(1) z'(Wx + b) \cdot 1$$

$$\frac{\partial L}{\partial b} = h(1) z'(Wx + b) \cdot 1$$



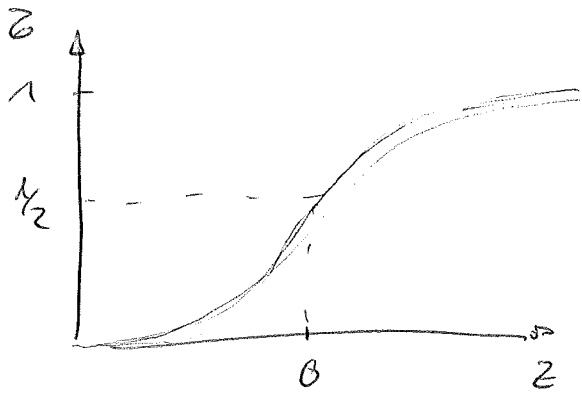
updates per epoch

$$W \mapsto W - \eta \frac{\partial L}{\partial W}$$
$$b \mapsto b - \eta \frac{\partial L}{\partial b}$$

even though we are predicting "very" wrongly we only take little notice and only adapt the learning in the last 100 epochs!

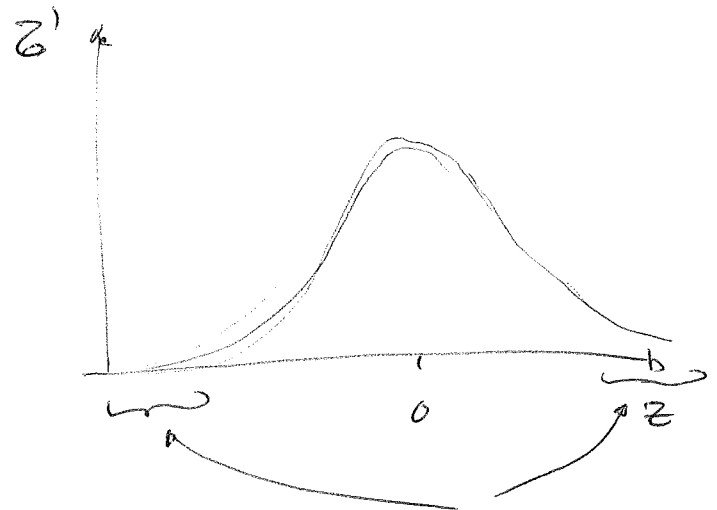
→ Show plots for weights, loss, output

Recall the sigmoid



$$z(z) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow z'(z) = z(z)(1 - z(z))$$



region where
the learning
is really slow
because $z'(z)$
is very small

Conclusion: In order to avoid this slow learning curves we have to change either 1) $z(z)$ or 2) L

$z(z)$ is however desirable
so let's try to find a better L .

ENTROPIE

$$H(P) = - \sum_{w \in \Omega} p(w) \log_2 p(w)$$

je überraschender
das Ereignis umso
weniger Information

→ hier Anzahl unabhängiger bits

Bsp: • Wurf einer Münze $P_{\text{fair}}(\text{Kopf}) = \frac{1}{2} = P_{\text{fair}}(\text{Zahl})$

$$H(P_{\text{fair}}) = - \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = 1$$

• Wurf eines Würfels $P_{\text{fair}}(\text{Augen} = a) = \frac{1}{6} \quad \forall a \in \{1, \dots, 6\}$

$$H(P_{\text{fair}}) = - \log_2 \frac{1}{6} = 2,585 \dots$$

Bem: 8-seitiger Würfel liefert $H(P_{\text{fair}}) = 3$

So ist $H(P)$ ein Maß für die durchschnittliche Überraschung, oder die Anzahl der Bits, die im Durchschnitt benötigt werden.

Bsp: • Unfairer Würfel

w	1	2	3	4	5	6
P	$\frac{1}{4}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{2}$

$$H(P) = 2$$

Die Augenzahl 6 kommt im Durchschnitt häufiger vor als bei einem fairen Würfel, deshalb braucht man im \emptyset weniger Bits.

Entropie ist also ein Maß für die Unschärfe über Wahrscheinlichkeitsverteilung P .

DEF: Für zwei Ω -Räume (Ω, P) und (Ω, Q) in diesem Fall def. wir die Kreuzentropie

$$H[P, Q] = - \sum_{\omega \in \Omega} P(\omega) \log_2 Q(\omega)$$

Interpretation: Q ist eine Schätzung von P .
Der Informationsgehalt für $\omega \in \Omega$ ist nach unserer Schätzung $-\log_2 Q(\omega)$ dieser tritt aber $P(\omega)$ verteilt auf, also ist $H[P, Q]$ der erwartete Informationsgehalt nach Q .

THEM: $H[P, Q] \geq H(P) = H(P, P)$

$$\begin{aligned} \forall x & \quad l(x) = -y \ln x - (1-y) \ln(1-x) \\ l'(x) &= -\frac{y}{x} + \frac{1-y}{1-x} = 0 \Rightarrow x=y \\ l''(x) &= \frac{y}{x^2} + \frac{1-y}{(1-x)^2} \geq 0 \end{aligned}$$

Interpretation von Neuron

Output $h(x) = \sigma(wx + b) \in \{0, 1\}$

bei input $x^{(i)}$ wird output $y^{(i)} \in \{0, 1\}$

erwartet, also ist $P_i(x^{(i)}) = y^{(i)}$ und $P_i(\text{not } x^{(i)}) = 1 - y^{(i)}$
 $Q_i(x^{(i)}) = h(x^{(i)})$ und $Q_i(\text{not } x^{(i)}) = 1 - h(x^{(i)})$

denn

$$H[P, Q] = -\frac{1}{N} \sum_{i=1}^N y^{(i)} \log_2 h(x^{(i)}) + (1 - y^{(i)}) \log_2 (1 - h(x^{(i)}))$$

② One neuron with cross-entropy



$$L(w, b) = -y \ln h(x) - (1-y) \ln(1-h(x)), \quad z = z(w, x + b)$$

(for training data $x=1, y=0$)

$$\frac{\partial L}{\partial w} = - \left[\frac{y}{h(x)} - \frac{1-y}{1-h(x)} \right] \frac{dh'(x)}{dw} = - \left[\frac{y}{z(z)} - \frac{1-y}{1-z(z)} \right] z'(z) \cdot x$$

$$\frac{\partial L}{\partial b} = -y$$

$$\text{but } z'(z) = z(z)(1-z(z))$$

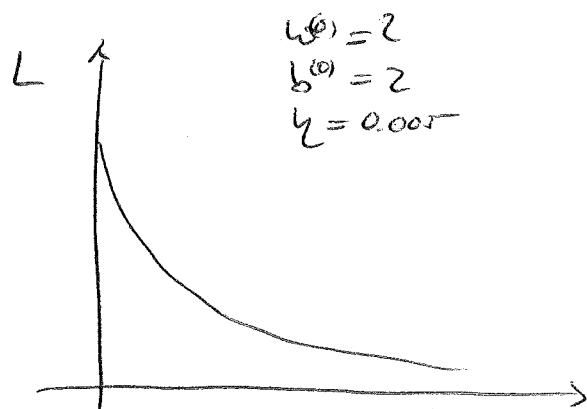
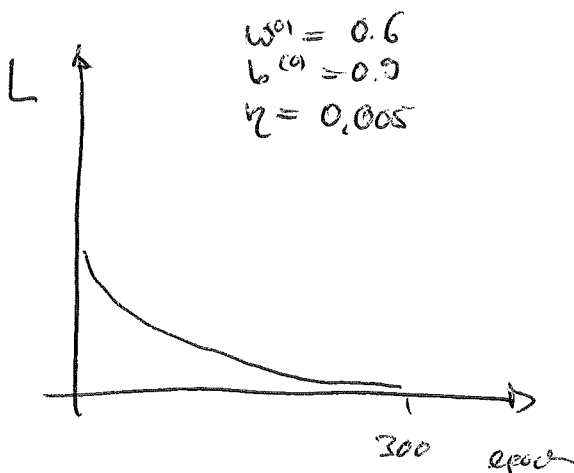
$$\begin{aligned} \Rightarrow \frac{\partial L}{\partial w} &= - \left[y(1-z(z)) - (1-y)z(z) \right] \cdot x \\ &= [z(z) - y] \cdot x \end{aligned}$$

and hence
for $x=1, y=0$

$$\frac{\partial L}{\partial w} = z(z)$$

$$\frac{\partial L}{\partial b} = z(z)$$

Hence, the neuron learns the weights proportional to the error made.



For the general case of $x^{(i)} \in \mathbb{R}^n$, $y^{(i)} \in \mathbb{R}^c$
 We may define

$$L(W^{(1)} \dots W^{(N)}, b^{(1)} \dots b^{(N)}) \\
= -\frac{1}{M} \sum_{i=1}^M \sum_{j=1}^c \left(y_j^{(i)} \ln h_j(x^{(i)}) + (1 - y_j^{(i)}) \ln h_j(x^{(i)}) \right)$$

REF: • Cross-entropy nearly always a better choice for sigmoid neurons

• however for $\sigma(z) = z$, i.e., in the linear neuron case the quadratic loss performs equally well

$$L(W, b) = \frac{1}{2} (Wx + b - y)^2 \\
\frac{\partial L}{\partial W} = (Wx + b - y)x \stackrel{x=y=0}{=} W + b \\
\frac{\partial L}{\partial b} = (Wx + b - y) = W + b$$

General case for cross entropy is:

$$\frac{\partial L}{\partial p^k} = \Delta(x^{(N)}) \cdot [z^{(N)}(z^{(N)}) \odot W^N] \dots [z^{(k)}(z^{(k)}) \odot \left(\frac{\partial W^k}{\partial p^k} \cdot x^k + \frac{\partial b^k}{\partial p^k} \right)] \\
= z(z^{(N)})(1 - z(z^{(N)}))$$

$$\rightarrow L'_i(x) = - \left[\frac{y_i}{z(z_i^{(N)})} + \frac{1 - y_i}{1 - z(z_i^{(N)})} \right] z'(z_i^{(N)}) \\
= (z(z_i^{(N)}) - y_i) = \text{Chances} \cdot (x^N - y_i)$$

Back to MNIST

VI

30 hidden neurons, 10 fixed cross-labels,
 $\eta = 0.5$, epochs 30

gives results close to quadratic loss $\sim 95.42\%$

100 hidden neurons and everything else the same $\sim 96.59\%$

REGULARIZATION

$$L' = L + \frac{\lambda}{N} \sum_k \|W^k\|_2^2$$

(for all k generally)

$$R(W^1 \dots W^N, b^1 \dots b^N)$$

$$\frac{\partial R}{\partial W^k} = \frac{\lambda}{N} \frac{\partial}{\partial W^k} \|W^k\|_2^2$$

new update rule:

$$W^k \mapsto W^k - \eta \frac{\partial L'}{\partial W^k} = \left(1 - \eta \frac{\lambda}{N}\right) W^k - \eta \frac{\partial L}{\partial W^k}$$

Example: $\lambda = 5$, same parameters as for
100 hidden neurons and cross-entropy

for $\eta = 0.1$ we break the 99% barrier

Other methods of regularization:

1) L_1 regularization $R = \frac{\lambda}{N} \sum_{l=1}^N \|w^l\|_1$

$$w \mapsto w - \eta \frac{\lambda}{N} \text{sign}(w) - \eta \frac{\partial L}{\partial w}$$

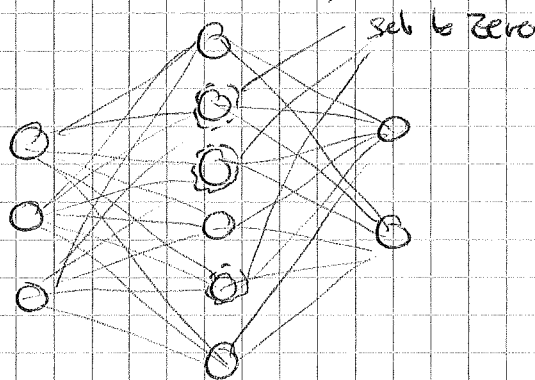
except for $w=0$

which shrinks the weights by a constant towards zero;

opposed to L_2 which is proportional to w itself.

$\Rightarrow L_1$ tends to concentrate on high number of high-importance connections

2) Dropout (2012)



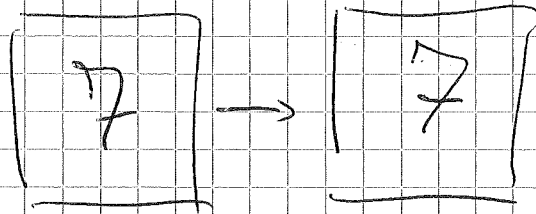
- 1) choose half of the hidden neurons and switch them off; learn weights
- 2) do this process again for other randomly chosen neurons deleted
- 3) average the weights of both networks

Idea: superposition of networks, average value

• minimizes co-adaptation

(2012) Dropout + L2 regularization gave 98.7%

3) artificially expanding the training set



- took a bit
- translate ...
- skew ... 98.9%
- elastic distortions 99.3%

Other improvements

Weight initialization

Ex: say we have 1000 neurons n in layer

say for 500 ^{its} ~~of~~ $W_{ij} = 0$

$$z_i = \sum_j W_{ij} x_j + b_i$$

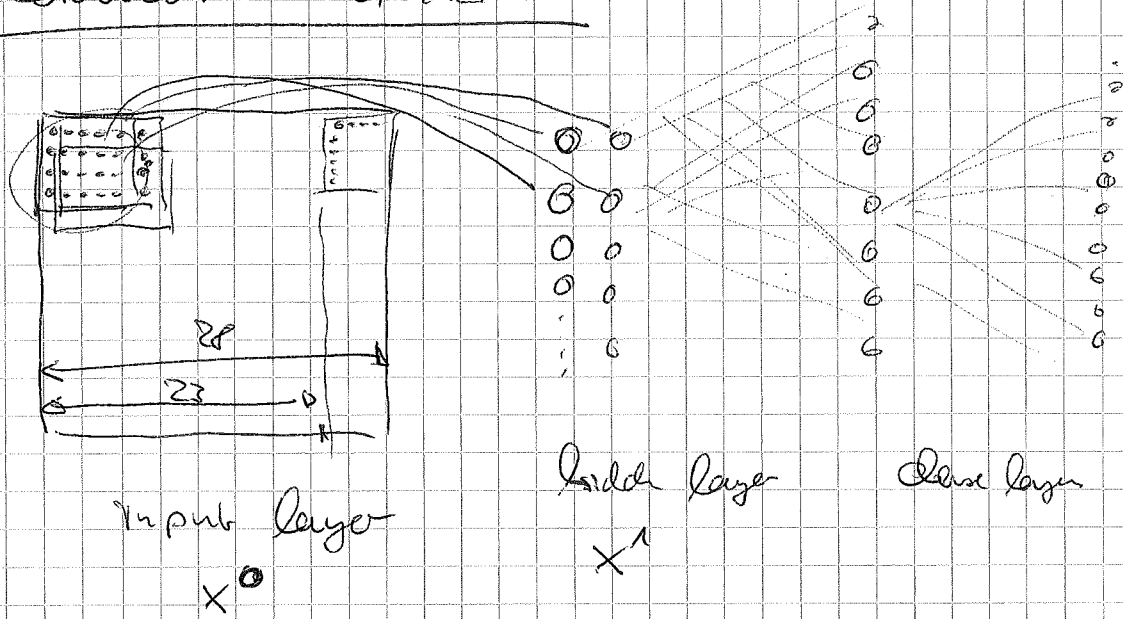
\sim random variable gaussian with mean zero and $\sigma \sim \sqrt{1/n} \approx 22\%$

\Rightarrow signals will be saturated

Solution: initialize weights by

$$R^{n_{k+1}, n_k} \Rightarrow W_{ij}^k \sim \mathcal{N}(0, \frac{1}{f(n_k)})$$

CONVOLUTION LAYERS



convolution layer

$$(X^0)_{ij} \rightarrow (z^1)_{ij} \quad \text{for} \quad z^1_{ij} = \sum_{k=0}^4 \sum_{l=0}^4 w_{kl} X^0_{i+k, j+l} + b$$

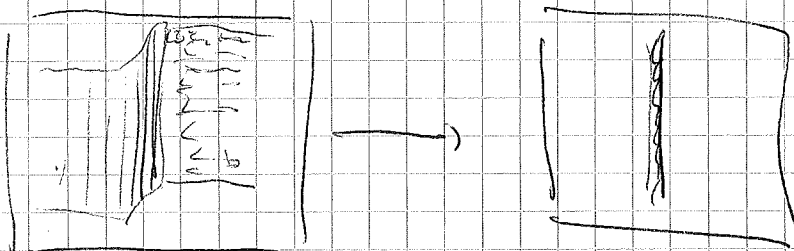
$$\downarrow$$

$$(z(z^1))_{ij}$$

What does a convolution layer do?

$$w = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{pmatrix}$$

Laplace filter



TUNING

① LEARNING BEHAVIOR

Let us consider one neuron with quadratic loss:

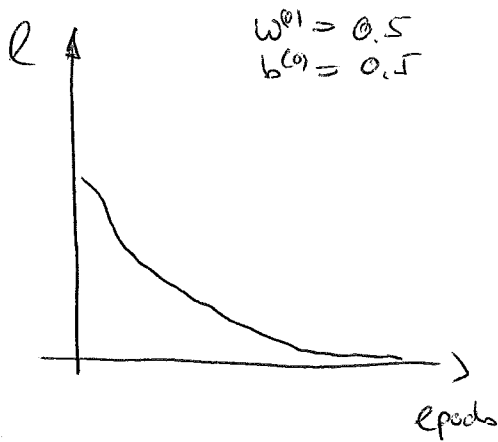
$$\mathbb{R} \ni x^0 \mapsto z^1 := \underbrace{w}_{\mathbb{R}} x^0 + \underbrace{b}_{\mathbb{R}} \mapsto x^1 = \sigma(z^1) \in \mathbb{R}$$

loss $l(y, x) = \frac{1}{2} (y - x)^2$

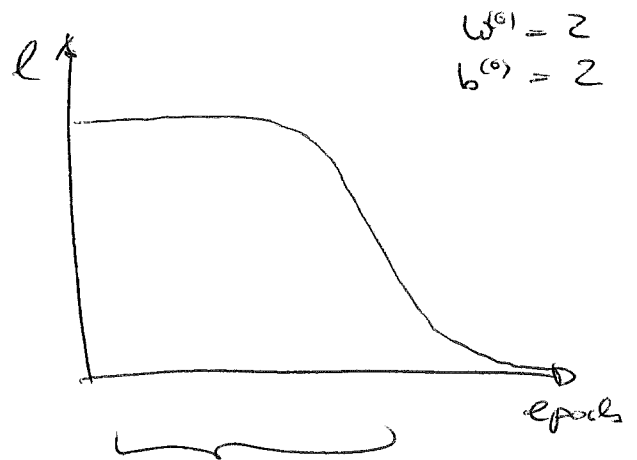
$$\Rightarrow \frac{\partial l}{\partial w} = [y - x^1(x)] (-1) \sigma'(z^1) \times$$

$$\frac{\partial l}{\partial b} = [y - x^1(x)] (-1) \sigma'(z^1)$$

} let's use
back to learn
 $x^0 = 1 \mapsto x^1 = 0$



same
 η

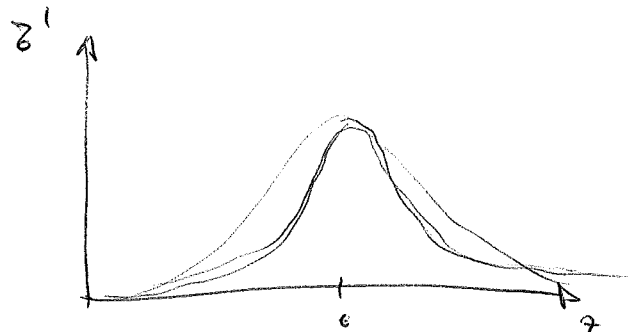
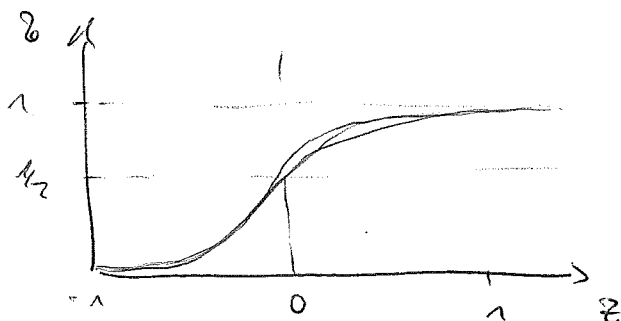


steads learning
using $w \mapsto w - \eta \frac{\partial l}{\partial w}$
 $b \mapsto b - \eta \frac{\partial l}{\partial b}$

long phase where
not much is happening

recall the sigmoid function

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \quad \sigma'(z) = \sigma(z)(1 - \sigma(z))$$



Condition: slow learning, due to saturation of the neuron.

Hence, either we change b or l ?

b is, however, desirable, therefore, let's change l .

New loss function chosen that b is compensated:

$$l(y, x^1) = -y \ln x^1 - (1-y) \ln (1-x^1)$$

is called the cross entropy, for which we find

$$\frac{\partial l}{\partial w} = - \left[\frac{y}{x^1} - \frac{1-y}{1-x^1} \right] \frac{\partial x^1}{\partial w}$$

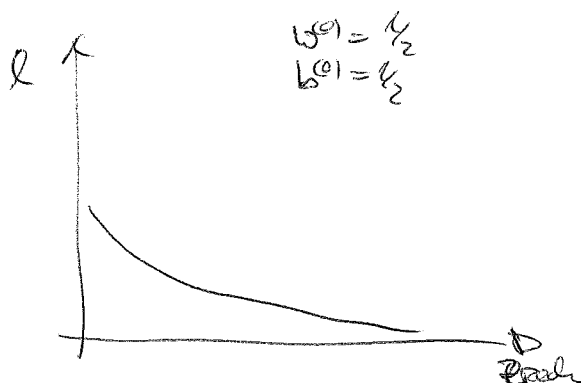
$$= - \left[\frac{y}{\sigma(z^1)} - \frac{1-y}{1-\sigma(z^1)} \right] \sigma'(z^1) x^0$$

$$\frac{\partial l}{\partial b} = - \left[\frac{y}{\sigma(z^1)} - \frac{1-y}{1-\sigma(z^1)} \right] \sigma'(z^1)$$

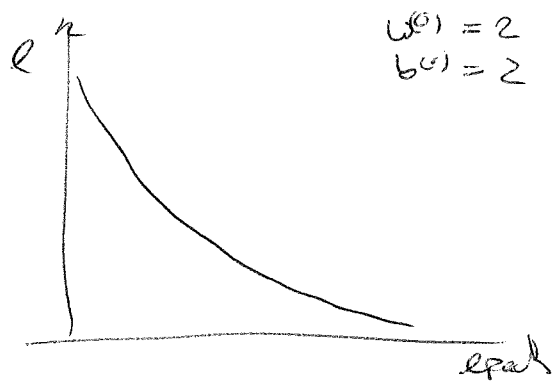
but since $\sigma'(z) = \sigma(z)(1-\sigma(z))$ we find

$$\frac{\partial l}{\partial w} = - [y - x^1] x^0$$

$$\frac{\partial l}{\partial b} = - [y - x^1]$$



same η



When does the entropy come from?

Let $(P, \Omega, \mathcal{F}(\Omega))$ be a probability space for a discrete Ω ,

DEF: (Entropy) $H(P) = - \sum_{\omega \in \Omega} P(\omega) \log_2 P(\omega)$

of independent bits needed to store the results

On average with respect to measure P

Example: • Tossing of a coin $P(\text{head}) = \frac{1}{2} = P(\text{tail})$

$H(P) = - \left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right] = 1$

• Casting a dice $P(\text{shows } a) = \frac{1}{6} \quad \forall a=1..6$

$H(P) = - \log_2 \frac{1}{6} \approx 2.58$

• Unfair dice

ω	1	2	3	4	5	6
P_{unfair}	$\frac{1}{4}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{2}$

$H(P_{\text{unfair}}) = 2$

The number 6 shows up more often, hence on average we need less bits to store the outcome

\Rightarrow Entropy is a measure of the uniformity of the probabilities

DEF. Let $(\Omega, P, \mathcal{P}(\Omega))$ and $(\Omega, Q, \mathcal{P}(\Omega))$ be two probability spaces for a discrete Ω we define the cross-entropy

$$H[P, Q] = - \sum_{\omega \in \Omega} P(\omega) \log_2 Q(\omega)$$

Interpretation: If Q is an estimate of the measure P then $-\log_2 Q(\omega)$ is the # of bits needed to store the information, $H[P, Q]$ is the (according to the estimate) the average v.r.b. P of this estimate.

THM: $H[P, Q] \geq H[P, P] = H[P, \#]$

Back to our δ -neurons:

We may interpret the input $x^{(i)}$ as probability event and the output $y^{(i)}$ as probability. The problem of our neural $x^1(x^{(i)})$ is the estimate to the real probability. The resulting cross-entropy is

$$L(y^{(i)}, x^1(x^{(i)})) = -y^{(i)} \log_2 x^1(x^{(i)}) - (1-y^{(i)}) \log_2 (1-x^1(x^{(i)}))$$

- HW
- 1) Does it play a role whether we take \log_2 or \ln ?
 - 2) Prove $L(y, x)$ has min at x and is convex.

Hence, for the general case we have

$$L = \frac{1}{n} \sum_{i=1}^M \sum_{j=1}^{n^N} (-1) \left[y_{ij}^{(i)} \ln(x_{ij}^{(i)}) + (1 - y_{ij}^{(i)}) \ln(1 - x_{ij}^{(i)}) \right]$$

HW Is the cross entropy always better concerning the slow learning? Consider for example the last layer to be linear-neurons.

② Initializing the weights and biases

- it makes sense to initialize the weights at random
- consider the output
$$z_i^k = \sum_j (w_{ij}^k x_j^{k-1}) + b_i^k$$

$\underbrace{\hspace{10em}}$
 n_{k-1} summands

If these were i.i.d. with variance = 1 we would get variance = n_{k-1} ⇒ *2-neuron saturates easily*

- hence ~~it~~ makes sense to normalize the initial random weights

$$w_{ij}^k \sim \mathcal{N}(0, 1/\sqrt{n^{k-1}})$$

HW: Compare the cases:

- all weights & biases $\sim \mathcal{N}(0, 1)$
- $\text{---} \text{---} \text{---} = 0$

③ Regularization

- In order to prevent that some weights become very dominant while others are seldomly used we may change our loss function

$$L' := L + \frac{\lambda}{N} \sum_{l=1}^N \frac{1}{2} \|W^l\|_2^2$$

which changes the update rule accordingly to

$$W^k \mapsto W^k - \eta \frac{\partial L'}{\partial W} = \left(1 - \eta \frac{\lambda}{N}\right) W^k - \eta \frac{\partial L}{\partial W^k}$$

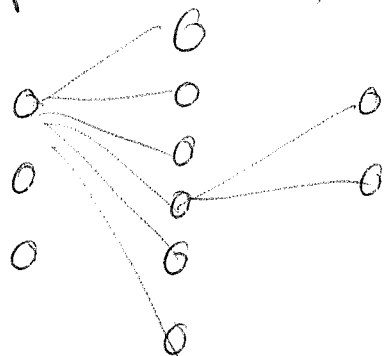
Python implementation: $\lambda = 5$, $\eta = 0.1$, cross-entropy, L_2 regularization breaks 9%

- HW Consider $\|W^l\|_1$ instead of L_2 norm

$$W^k \mapsto W^k \left(1 - \eta \frac{\lambda}{N} \text{sign}(W^k)\right) - \eta \frac{\partial L}{\partial W^k}$$

What is the new distance?

- Dropouts (2012)



1) choose half neurons of hidden layer and delete them

2) learn weights

3) choose other half, delete the learn weights

4) average over the two networks

→ Superposition, minimization of coadaptation

• automatically expanding the training data

hidden layer

APPROXIMATION BY NEURAL NETWORKS

1

Representation of boolean functions

DEF:

The goal is to take an arbitrary boolean function

$$f: \{0, 1\}^d \rightarrow \{0, 1\}$$

neural net,
vocabulary,
feed forward

and find a neural network, i.e.,
activation function σ^i , weights w^i and
biases b^i for $i=1 \dots N$ layers s.t.

$$\begin{aligned} f(x) &= x^N(x) = \sigma^N(w^N x^{N-1}(x) + b^N) \\ &= \sigma^N(w^N \sigma^{N-1}(w^{N-2} x^{N-3} + b^{N-2}) + b^N) \\ &\dots \\ &= \sigma^N(\dots \dots \dots (\sigma(b^1 x + b^1))) \end{aligned}$$

THEM: Every boolean function can be exactly
represented by a neural network with
 $N \geq 2$ and at most 2^d neurons in the
hidden layer and $\sigma(z) = \begin{cases} 1 & \text{for } z \geq 0 \\ 0 & \text{for } z < 0 \end{cases}$.

Proof: Let us denote $f^{-1}(\{1\}) = \{a_1, a_2, \dots, a_m\}$

Then we find a representation

$$f(x) = \sigma\left(-1 + \sum_{i=1}^m \delta_{x, a_i}\right)$$

Next, we find a representation for $\delta_{x, a} = \begin{cases} 1 & \text{for } x=a \\ 0 & \text{for } x \neq a \end{cases}$
for $x, a \in \{0, 1\}^d$.

For this we consider for $d=1$

$$\begin{array}{ccc} a & b & 2ab - a - b \\ 0 & 0 & 0 \\ 0 & 1 & -1 \\ 1 & 0 & -1 \\ 1 & 1 & 0 \end{array}$$

Hence, $\delta_{x,a} = \sigma(2xa - x - a)$ and for $d \geq 1$ we get

$$\delta_{x,a} = \sigma\left(\sum_{j=1}^d 2x_j a_j - x_j - a_j\right)$$

which allows to recast f into

$$f(x) = \sigma\left(-1 + \sum_{i=1}^m \sigma\left(\sum_{j=1}^d 2x_j a_j^i - x_j - a_j^i\right)\right)$$

Here we can read off the weights and biases

$$\bullet \ y_i^1 = w^1 x + b^1 \stackrel{!}{=} \sum_{j=1}^d 2x_j a_j^i - x_j - a_j^i$$

$$\Rightarrow w_{ij}^1 = 2a_j^i - 1 \quad 1 \leq i \leq m, 1 \leq j \leq d$$

$$b_i^1 = -\sum_{j=1}^d a_j^i \quad 1 \leq i \leq m$$

$$x_i^1 = \sigma(y_i^1)$$

$$\bullet \ y^2 = w^2 x^1 + b^2 \stackrel{!}{=} \sum_{i=1}^m x_i^1 - 1$$

$$\Rightarrow w_i^2 = 1$$

$$b^2 = -1$$

and finally

$$x^2 = \delta(y^2)$$

Hence, we need weights and biases

$$W^1 \in \mathbb{R}^{m \times d}, \quad b^1 \in \mathbb{R}^m$$

$$W^2 \in \mathbb{R}^m, \quad b^2 \in \mathbb{R}, \quad \text{and } N=2$$

where d is the number of input neurons
 m is the number of hidden neurons

and for m we get the estimate

$$0 \leq m \leq 2^d$$

for the worst case where the function
is highly complex $f(x) = 1 \quad \forall x$. \square

HW: Implement XOR gate by hand.

REM: • This construction implies an exponential increase in hidden layer neurons.

- Any network with $N > 2$ is of course also able to represent the binary function, too.

- Bound can be sharpened to

$$0 \leq m \leq \min(|f^{-1}(\{0\})|, |f^{-1}(\{1\})|)$$

Representation of binary classification on \mathbb{R}^d

The goal is to take an arbitrary ~~test~~-data

set $\mathcal{D} = (x^{(i)}, y^{(i)})_{1 \leq i \leq M}$ and find a

function $f: \mathbb{R}^d \rightarrow \{-1, +1\}$ s.t.

$$f(x^{(i)}) = y^{(i)} \quad \forall 1 \leq i \leq M \quad (*)$$

that can be represented by a neural network

THEM: Given test data $(x^{(i)}, y^{(i)})_{1 \leq i \leq M}$, there is an $f: \mathbb{R}^d \rightarrow \{-1, +1\}$ that fulfills (*) and can be represented by a neural network with $N=2$ and not more than $2M$ hidden neurons and $\sigma(z) = \text{sgn}(z)$.

Proof: Def. $\{a_1, \dots, a_m\} = \{x^{(i)} \mid y^{(i)} = 1, 1 \leq i \leq M\}$

Since this set is finite, for each $1 \leq i \leq m$ we find $w_i, b_i \in \mathbb{R}^d$ s.t.

$$h_i: z \mapsto w_i z + b_i \quad \text{with} \quad \text{range } h_i \cap \{a_1, \dots, a_m\} = \{a_i\}$$

furthermore, also due to finiteness, there is an $\varepsilon_i > 0$ (e.g. half of min distance to nearest $a_j, j \neq i$) s.t.

$$I_i: z \mapsto \sigma(\varepsilon_i + (w_i z + b_i)) + \sigma(\varepsilon_i - (w_i z + b_i))$$

fulfills

$$I_i(a_i) = 1$$

$$I_i(a_j) = 0 \text{ for } j \neq i$$

(iii)

With this we construct

$$f(z) = \sigma\left(-1 + \sum_{i=1}^m I_i(z)\right)$$

which by construction fulfills

$$f(a_i) = 1, \quad f(z) = 0 \quad \forall z \neq a_i, 1 \leq i \leq m$$

$$\Rightarrow f(x^{(1)}) = y^{(1)}$$

From f we read off the appropriate weights and biases of the corresponding neural network with $N=2$

$$\mathbb{R}^d \ni x^0 = z$$

\downarrow

$$\mathbb{R}^m \ni y^1 = (y_i^0)_{1 \leq i \leq 2m} \text{ with}$$

$$= W^1 x^0 + b^1 \text{ for}$$

for $1 \leq l \leq m$

$$y_{2l-1} = \sigma_+(W_l x^0 + b_l)$$

$$y_{2l} = \sigma_-(W_l x^0 + b_l)$$

$$W_{2l-1, j}^1 = W_{l, j}, \quad b_{2l-1} = b_l$$

$$-W_{2l, j}^1 = W_{l, j}, \quad b_{2l} = -b_l$$

$$x^1 = \sigma(y^1)$$

\downarrow

$$y^2 = W^2 x^1 + b^2 \text{ for } W_i^2 = 1, b^2 = -1$$

\downarrow

$$x^2 = \sigma(y^2)$$

Clearly, the number of hidden neurons
 \Rightarrow bounded by

$$0 \leq 2m \leq 2M.$$

□

REF: • The bound on the number of hidden neurons can easily be improved if test data is distributed in general position, e.s.,

Approximation of real-valued functions

Goal \Rightarrow to show that a feedforward neural network can approximate any function $f \in C^0([a, b])$.

The set of hypothesis of a $N=2$ feed-forward neural network with m hidden neurons, $\sigma^1(z) = \begin{cases} 1 & \text{for } z \rightarrow +\infty \\ 0 & \text{for } z \rightarrow -\infty \end{cases}$ and bounded and $\sigma^2(z) = z$ is given by

$$\mathcal{H}_m^z = \left\{ g: \mathbb{R} \rightarrow \mathbb{R} \mid g(x) = \sum_{i=1}^m w_i^2 \sigma(w_i^1 x + b_i^1) \right. \\ \left. \text{for } w^1 \in \mathbb{R}^{m \times 1}, w^2 \in \mathbb{R}^{1 \times m} \text{ and } b^1 \in \mathbb{R}^m \right\}$$

THM: Let $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ bounded s.t. $\lim_{z \rightarrow +\infty} \sigma(z) = 1$ and $\lim_{z \rightarrow -\infty} \sigma(z) = 0$. Then $\exists C(\sigma) < \infty$ s.t.
 $\forall m \in \mathbb{N}: \inf_{g_m \in \mathcal{H}_m^z} \|f - g_m\|_\infty \leq C(\sigma) \sup_{|x-y| < \frac{1}{m}} |f(x) - f(y)|.$

REM: $C(\mathcal{Z}) = \|\delta\|_{\infty}$

(IV)

- Distance of best and
- for $\mathcal{Z}(\mathcal{Z}) \neq \mathcal{B}_m \mathcal{Z}(\mathcal{Z})$ suffices
 $\mathcal{Z} \rightarrow +\infty$ $\mathcal{Z} \rightarrow -\infty$
- f Lipschitz $\Rightarrow \sup_{|x-y| \leq \frac{1}{m}} |f(x) - f(y)| \leq \frac{L}{m}$

Proof: It suffices to construct a best

$$g_m \in \mathcal{D}_m$$

to estimate the minimum

$$\inf_{g \in \mathcal{D}_m} \|f - g\|_{\infty} \leq \|f - g_m\|_{\infty}$$

wlog $f \in C^0([0, 1])$, i.e., $a=0, b=1$

$$\text{def. } x_i = \frac{i}{m} \quad i = 0, \dots, m+1$$

$$\text{and } g_m(x) = f(x_i) \text{ for } x \in [x_i, x_{i+1})$$

which can be written as

$$g_m(x) = f(x_0) + \sum_{i=1}^m [f(x_{i+1}) - f(x_i)] \mathbb{1}_{x \in [x_i, x_{i+1})}$$

which we need to estimate with an element $g_m \in \mathcal{D}_m$, which for $\alpha > 0$ we take

$$g_m(x) = f(x_0) + \sum_{i=1}^m [f(x_{i+\alpha}) - f(x_i)] \mathbb{1}_{x \in [x_i, x_{i+1})}$$

\cap
 \mathcal{D}_m

REM: (PPP)

- \mathcal{Z} may not be dense polynomial simply because $(\mathcal{H}_m^{\mathcal{Z}})$ would still be only polynomials of the order d and therefore cannot be dense in $\mathcal{C}^0([a,b])$
- also increasing the number of hidden layers does not help
- on the contrary $(\mathcal{H}_m^{\mathcal{Z}})$ is dense in $\mathcal{C}^0([a,b])$ iff \mathcal{Z} is not a polynomial - we however proved this for a certain subset of non-polynomial \mathcal{Z} only.

LEM: For \mathcal{Z} as before, $\bigcup_m \mathcal{H}_m^{\mathcal{Z}}$ is dense in $\mathcal{C}^0(\mathbb{R}^d, \mathbb{R}^c)$

Proof: WLOG we may assume $c=1$ because of we can approximate each $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$ of $f = (f_i)_{i=1..c}$ to accuracy $\|f_i - \hat{f}_i\|_{\infty} \leq \frac{\epsilon}{c}$ with m neurons, we can approximate f with $\|f - \hat{f}\|_{\infty} < \epsilon$ with $m \cdot c$ neurons.

In order to use our result for $\mathcal{C}^0([a,b], \mathbb{R})$ we use that

$$\mathcal{D} := \text{span} \{ e^{w \cdot x} \mid w \in \mathbb{R}^d \}$$

is dense in $\mathcal{C}^0([a,b]^d, \mathbb{R})$ according to Stone-Weierstrass, since

- i) \mathcal{D} forms an algebra under multiplication and linear combination (Vector space closed under multiplication)
- ii) $1 \in \mathcal{D}$ (non-zero constant function) (non-zero constant function)
- iii) $\forall x, y \in [a,b]^d \exists w = x - y$ s.t. $e^{w \cdot x} \neq e^{w \cdot y}$ (separation of points)

$\Rightarrow \forall \varepsilon > 0 \exists m \in \mathbb{N}, v_1, \dots, v_m \in \mathbb{R}^n, c_1, \dots, c_m \in \mathbb{R}$
s.t.

$$\sup_{x \in [a, b]^n} \left| f(x) - \sum_{i=1}^m c_i e^{v_i \cdot x} \right| < \frac{\varepsilon}{2}$$

Now we only need to approximate the exponentials, which can be done using our last theorem: $\exists \tilde{a} < \tilde{b}$ s.t.

$$[\tilde{a}, \tilde{b}] \supseteq \bigcup_{i=1}^m \{v_i \cdot x \mid x \in [a, b]^n\}$$

Then, for all $\varepsilon > 0 \exists$ $f_{\varepsilon}(y) = \sum_{j=1}^l w_j^2 \delta(w_j^1 y + b_j^1)$
with

$$\|e^y - f_{\varepsilon}\|_{\infty} < \frac{\varepsilon}{2 \cdot m \max_{k=1, \dots, l} |c_k|} \in \mathbb{R}_m^2$$

$$\Rightarrow \sup_{x \in [a, b]^n} \left| f(x) - \sum_{i=1}^m c_i \sum_{j=1}^l w_j^2 \delta(w_j^1 v_i \cdot x + b_j^1) \right|$$

$$\leq \sup_{x \in [a, b]^n} \left| f(x) - \underbrace{\sum_{i=1}^m c_i e^{v_i \cdot x}}_{= f_{\varepsilon}} \right| < \frac{\varepsilon}{2}$$

$$+ \sup_{x \in [a, b]^n} \left| \sum_{i=1}^m c_i (e^{v_i \cdot x} - f_{\varepsilon}(v_i \cdot x)) \right|$$

$$\leq m \max_{k=1, \dots, l} |c_k| \sup_{y \in [a, b]} |e^y - f_{\varepsilon}(y)|$$

$$< \frac{\varepsilon}{2 \cdot m \max_{k=1, \dots, l} |c_k|}$$

$< \varepsilon \quad \square$

- REN:
- each time step looks like a usual feed forward network with two layers
 - however, information about the last step may "flow through time" by dependence on h_{t-1}
 - memory certainly depends on size of hidden layer

Learning and update rule

We define a loss function for each sequence $(x_1, \bar{y}_1), \dots, (x_N, \bar{y}_N)$

$$L = \sum_{t=1}^N \ell(y_t(x_t), \bar{y}_t) \quad , \quad \ell(y, \bar{y}) = \frac{1}{2} (y - \bar{y})^2$$

and perform a mini-batch gradient descent update over all sequences of length N :

The free parameters are W, b, V :

$$W \in \mathbb{R}^{n \times d} \quad , \quad b \in \mathbb{R}^n$$

$$V \in \mathbb{R}^{c \times n}$$

$$\begin{pmatrix} W \\ b \\ V \end{pmatrix} \mapsto \begin{pmatrix} W \\ b \\ V \end{pmatrix} - \eta \begin{pmatrix} \frac{\partial L}{\partial W} \\ \frac{\partial L}{\partial b} \\ \frac{\partial L}{\partial V} \end{pmatrix}$$

for learning rate η .

Algorithm: For length $L = 1 \dots N$:

mini-batch grad. desc. for all seq. of length L .

Backpropagation through time

11

In order to compute the update efficiently we regard an unrolling of the recurrent network

$$z_t := W_t x_t + U_t h_{t-1} + b_t$$

$$h_t := \sigma(z_t)$$

$$y_t := V h_t$$

This is the same network except for the fact that the weights are not shared. Hence, the loss function is a function

of

$$\tilde{L} = \sum_{t=1}^N \tilde{\ell}(w_1, \dots, w_t, u_1, \dots, u_t, b_1, \dots, b_t, V, x_1, \dots, x_t, y_t)$$

However, the loss function of the recurrent network is given by

$$L = \tilde{L} \Big|_{w_t=w, u_t=u, b_t=b \text{ for } t=1, \dots, N}$$

This implies that

$$\begin{aligned} \frac{\partial L}{\partial w} &= \sum_{t=1}^N \frac{\partial \tilde{\ell}(y_t, \tilde{z}_t) \Big|_{w:=w}}{\partial w} \\ &= \sum_{t=1}^N \sum_{i=1}^t \frac{\partial \tilde{\ell}(y_t, \tilde{z}_t)}{\partial w_i} \Big|_{w_i=w \text{ } i=1, \dots, t} \end{aligned}$$

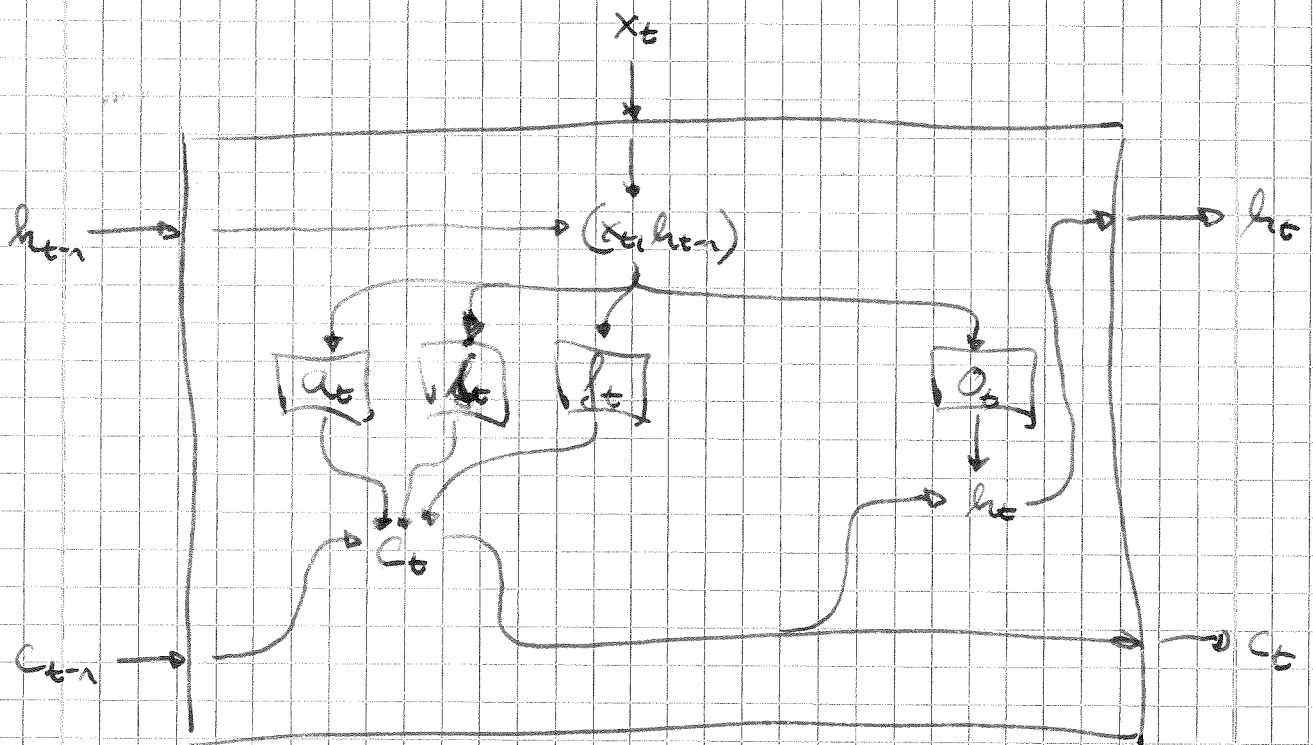
and we can use ordinary backpropagation for the unrolled network summing the gradients.

RECI: Since the unrolled networks become as deep as the length of the sequence, problems in the training become significant.

$\frac{\partial L}{\partial W}$ expands to a ^{long} product of matrices:

- regions of slow learning are extremely slow
- regions of fast learning explode.

LONG SHORT TERM MEMORY NETWORKS



add $a_t = \tanh(W_a x_t + U_a h_{t-1})$

input $i_t = \sigma(W_i x_t + U_i h_{t-1})$

forget $f_t = \sigma(W_f x_t + U_f h_{t-1})$

output $o_t = \sigma(W_o x_t + U_o h_{t-1})$

Control how much of the input to write to c_t , how much to forget

Control how much to put out for h_t

memory operation

$$c_t = i_t \odot a_t + f_t \odot c_{t-1}$$

$$h_t = o_t \odot \tanh(c_t)$$