

Machine Learning for Natural and Life Sciences

1	Introduction:.....	1
	what is machine learning, why do we need it.....	1
	features and response, notation	1
	classical solutions (hand-crafted formula with thresholding, nearest-neighbor methods) and their shortcomings.....	3
2	Loss functions and performance evaluation:	5
	squared loss, 0/1 loss	5
	training error vs. generalization error	6
	necessity of separate test sets or cross-validation	6
	lower bounding the loss via Bayes optimal rule	8
3	Linear models:.....	9
	linear decision rules and boundaries	9
	LDA, squared loss.....	10
	linear SVM, hinge loss.....	11
	logistic regression, logistic loss	12
	optimization of LR with stochastic gradient descent.....	14
4	Unbalanced risk:	15
	minimization of expected risk.....	15
	confusion matrix	16
	sensitivity/specificity vs. precision/recall.....	16
	ROC curves and PR curves	16
5	Regression (briefly):.....	17
	ordinary least squares.....	17
	regularization (ridge regression, LASSO)	18
	robust regression, robust loss functions.....	19
6	Non-linear methods (very briefly):	21
	augmented feature spaces, kernel trick.....	21
	boosting	22
	descision trees.....	22
	motivation for neural networks (they combine all of these ideas).....	22
7	Neural networks:	23
	definition of a neuron	23
	activation functions (sigmoid, tanh, ReLU, leaky ReLU, PReLU, ELU)	23
	fully connected architectures.....	24
	softmax	25
	loss functions revisited	27
	backpropagation.....	28
	training tricks (weight initialization, training rate schedule, ADAM, mini-batches, batchnorm, dropout, weight decay, data augmentation)	31
	convolutional architectures, U-net.....	35
	famous architectures (LeNet, AlexNet, VGG, ResNet) examples	

Machine Learning for Natural- and Life-Science

1

- Goal: Teach the background of ML, not just recipes.
 \Rightarrow help you understand ML literature
- Setting:
 - we are interested in properties or quantities Y , but these are ~~hard to~~ impossible or impractical to measure directly $\hat{=}$ "response"
 - properties / quantities X are related to Y and are easier to measure / observe $\hat{=}$ "features"
- \Rightarrow approximate Y^* (true values, "ground truth") by a formula $\hat{Y} = f(X) \approx Y^*$
 - \Leftarrow approximate responses of our model, "predictions"
 - problem: the formula $f(X)$ is not known
- \Rightarrow collect training data and fit the model to $f(X)$ reproduce the behavior of these data
- supervised learning: we cannot measure Y under normal field conditions, but we can obtain it in an experimental setting
 - example: crash test to determine car safety

$$TS = \{(X_i, Y_i)\}_{i=1}^N \quad N \text{ training instances}$$

X_i, Y_i : features and response of instance i

X $\hat{=}$ feature matrix is simply a table

Instance	Instance Index	Height	Weight	Gender	
Alice	$i = 1$	1.60	80	F	
Bob	$i = 2$	1.85	85	M	...
Carol	$i = 3$	1.70	70	F	
:			$x_{3,2}$		\hat{x}_j

Feature index $\rightarrow j=1, j=2, \dots, j=D$

$j \in 1, \dots, D$

D-dimensional feature space

Problem: for simplicity of the math and software, we want X to be real-valued, but many features are discrete / categorical, e.g. gender

\Rightarrow "one-hot encoding": one flag per label

modified table / matrix	gender	
	male	female
Alice	0	1
Bob	1	0
Carol	0	1

$\Rightarrow X \in \mathbb{R}^{N \times D}$, likewise Y

[if only two labels, we need only ~~one~~ encode one, because the other is implied, e.g. male = 1 - female]

- simple solutions to the problem

- threshold classifier:

- e.g. $\hat{Y} = \begin{cases} 0 & \text{patient is healthy} \\ 1 & \text{patient suffers from diabetes} \end{cases}$

cannot be measured directly \Rightarrow measure a bio-marker

~~X: sugar content in blood~~

~~X: blood sugar concentration~~

classifier: $Y = \begin{cases} 0 & \text{if } X < 100 \text{ mg/dl} \\ 1 & \text{if } X \geq 100 \text{ mg/dl} \end{cases}$

(simplified for illustration purposes)

- problem: often, there is no single feature that indicates the response

\Rightarrow combine several features in some formula returning a "score"

e.g. $Y = \begin{cases} 0 & \text{person is normal} \\ 1 & \text{person is obese / overweight} \end{cases}$

score = $\frac{\text{weight [kg]}}{\text{height}^2 [\text{m}^2]}$ "body mass index"

$Y = \begin{cases} 0 & \text{if } \text{BMI} < 25 \\ 1 & \text{if } \text{BMI} \geq 25 \end{cases}$

• problem: hand-crafted formulas have disadvantages (3)

- hard to find, expensive (may require decades)
- in complex situations, simplifications are necessary
 ⇒ loss of accuracy

⇒ overcome these difficulties using machine learning

(downside: machine learning is often less explainable than an explicit formula ⇒ explainable ML is hot research area)

- nearest(-neighbor) classifier

- classify a new instance x_{test} like the most similar instance in the training set
- define similarity by a distance function of the features

$$dist(x, x') = \begin{cases} \text{small} & \text{if } x \text{ and } x' \text{ are similar} \\ \text{large} & \text{--- if --- different} \end{cases}$$

e.g. Euclidean distance between features

$$dist(x, x') = \sqrt{\sum_{j=1}^D (x_j - x'_j)^2}$$

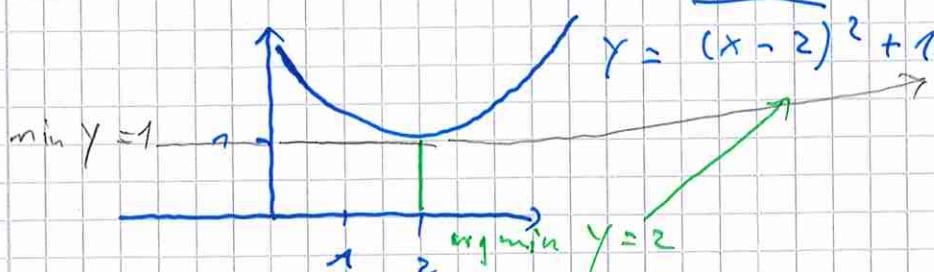
↑ sum over features / index j

- arg min - Notation:

$$i_{\text{nearest}} = \arg \min_{i=1}^N dist(x_{test}, x_i)$$

$$\min_{i=1}^N dist(x_{test}, x_i) : \begin{array}{l} \text{- iterate over all training inst.} \\ \text{- compute distance} \\ \text{- return the minimal distance} \end{array}$$

$$\arg \min_{i=1}^N \quad \begin{array}{l} \text{- iterate over all training instances} \\ \text{- compute distance} \\ \text{- return } \underline{\text{index}} \text{ of instance with minimal} \\ \text{distance} \end{array}$$



$$\Rightarrow \hat{Y}_{\text{test}} = Y^*_{\text{nearest}} = Y_{i_{\text{nearest}}} \quad \text{classifier response}$$

- often works pretty well, but has a number of shortcomings
 - need to memorize TS and search with $O(N)$
in each query \Rightarrow scales poorly to large TS
 - ~~memorize~~ only important training instances
 \Rightarrow hard to identify
 - speed-up by approximate nearest neighbor search
 - \hat{Y} is not guaranteed to converge toward Y^* even if $N \rightarrow \infty$ (error may be twice as big as the best achievable)
 - k-nearest neighbor \Rightarrow more expensive
 - requires hand-crafted distance function
 \Rightarrow same problems as hand-crafted score
 \Rightarrow metric learning (Björn Orner)
 - human similarity judgements are hard to model
 - show Voronoi diagram (animation)

Loss functions

5

- classical models:
- define a specific formula for every situation
(e.g. physics: Newtons laws, law of gravity, Maxwell's laws, etc.)
 - theoreticians hand-craft these models
- machine learning
- use generic model families that apply to many situations and have free parameters

$$\Theta : Y = f_\Theta(X)$$

e.g. linear: $Y = \alpha X + \beta$
 $\hookrightarrow \Theta = (\alpha, \beta)$

- among all possible choices Θ for Θ , pick $\hat{\Theta}$ such that $Y_i \approx f_{\hat{\Theta}}(X_i)$ works as well as possible for the TS $\{(X_i, Y_i)\}_{i=1}^N$
- "model fitting / training"

\Rightarrow need a way to quantify the "goodness of fit", i.e. to compare the performance of different choices of Θ

let $\hat{Y} = f_\Theta(X)$ the prediction for a data set X , given Θ
 $Y^* = f^*(X)$ the ~~true~~ true relationship
 (we know Y^* , but not f^*)

loss(Y^*, \hat{Y}) $\rightarrow \mathbb{R}^+$ measures the difference between truth and predictions

$$\begin{cases} = 0 & \text{if } Y^* \approx \hat{Y} \\ \text{big} & \text{otherwise ("lossy")} \end{cases}$$

if the data instances are iid (independent of each other), but generated from the same true model f^*)

This simplifies to an additive loss over all instances

$$\text{loss}(Y^*, \hat{Y}) = \frac{1}{N} \sum_{i=1}^N \text{loss}(Y_i^*, \hat{Y}_i)$$

The particular formula for the $\text{loss}(\cdot)$ determines what is considered a good fit and must be selected by the model designer / data scientist according to the application.

- most common choice : squared loss

$$\text{loss}(Y^*, \hat{Y}) = \frac{1}{N} \sum_{i=1}^N (Y_i^* - \hat{Y}_i)^2$$

$$= 0 \Leftrightarrow Y_i^* = \hat{Y}_i \text{ for all } i$$

e.g. for 2-class classification

$$Y = \begin{cases} 0 & \text{if healthy} \\ 1 & \text{if disease} \end{cases}$$

$$(Y_i^* - \hat{Y}_i)^2 = \begin{cases} 0 & (\Rightarrow Y_i^* = \hat{Y}_i \text{ correct prediction}) \\ 1 & (\Rightarrow Y_i^* \neq \hat{Y}_i \text{ wrong prediction}) \end{cases}$$

$\hat{\square}$ count the wrongly predicted instances $\hat{\square}$ one error

- $\hat{Y}_i = f_\Theta(X_i)$ depends on the choice of Θ , and so will be loss \Rightarrow choose Θ to minimize the loss of ~~the~~

on the TS

$$\hat{\Theta} = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \text{loss}(Y_i^*, \hat{Y}_i = f_\Theta(x_i))$$

Finding the optimal $\hat{\Theta}$ is in general a complicated optimization problem \Rightarrow later.

The loss after optimising $\hat{\Theta}$ is the "training error"

$$\text{err} = \frac{1}{N} \sum_i \text{loss}(Y_i^*, \hat{Y}_i = f_{\hat{\Theta}}(x_i))$$

$$= \min_{\Theta} \frac{1}{N} \sum_i \text{loss}(Y_i^*, \hat{Y}_i = f_\Theta(x_i))$$

$\hat{\square}$ because we chose $\hat{\Theta}$ such that it achieves the minimum

- what we are actually interested in is the performance of our model in the field, i.e. on arbitrary data beyond the TS

$$\text{Err} = \mathbb{E}_{X, Y \sim p^*(X, Y)} [\text{loss}(Y^*, \hat{Y} = f_{\hat{\Theta}}(X))]$$

"generalization ~~order~~⁴ error"

- fundamental insight of ML and statistics

(7)

$$\text{Err} > \text{err}$$

generalization error is usually bigger than training error,
sometimes much bigger $\hat{\approx}$ "overfitting"

\Rightarrow take home message: fundamental insight 1

IT MAKES NO SENSE TO REPORT MODEL PERFORMANCE IN TERMS OF THE TRAINING ERROR err

- instead we must estimate and report the generalization error.

In practice, empirical estimation works best. 2 possibilities:

- independent test set: $\text{Test} = \{(x_i^{\text{test}}, y_i^{\text{test}})\}_{i=1}^{N'}$

compute the average loss of the test set

$$\text{Err} = \frac{1}{N'} \sum_{i=1}^{N'} \text{loss}(y_i^{\text{test}}, f_{\theta}(x_i^{\text{test}}))$$

Note: - Never use the test data for training.

- θ remains fixed now.

- cross-validation: if no independent test set is available, create it from the training set:

TS (randomly shuffled)



$\ell = 1 \quad 2 \quad 3 \quad 4 \quad 5$

Create k ~~sets~~ random subsets ("folds"), here $k=5$

\Rightarrow " k -fold cross-validation"

Let TS_ℓ denote fold ℓ , $TS_{\neq \ell}$ the remaining data
for $\ell = 1 \dots K$

- train on $TS_{\neq \ell}$ to get $\hat{\theta}_\ell$

- estimate Err_ℓ on TS_ℓ using $\hat{\theta}_\ell$ fixed
estimate: $\text{Err} = \frac{1}{k} \sum_{\ell} \text{Err}_\ell$

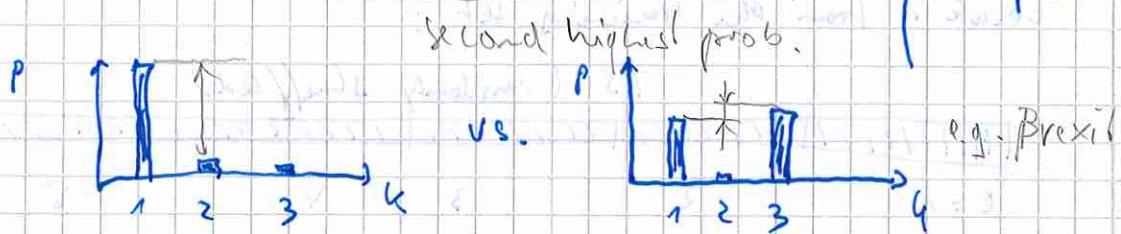
typical values for k : 10, 5 (2 if training is
very expensive); $K=N$ LOO-CV for theory

- Infrastructure for CV is available in all credible ML libraries (sk-learn, pytorch, ...)
- ⇒ don't implement it yourself (except as an exercise)

Bayes decision rule

- instead of a hard class decision, we can return a soft response: the probability of each label, given the features
 $\forall i: p(Y_i=c | X_i) \quad \text{with } \sum_{c=1}^C p(Y_i=c | X_i) = 1$
- a hard decision is easily obtained by deciding for the most probable class: $\hat{Y}_i = \arg \max_c p(Y_i=c | X_i)$
- many methods learn to approximate $p^*(Y|X)$ (the "posterior probability" \approx after measuring the features) as well as possible \Rightarrow more information than a hard decision rule, e.g. uncertainty which via margin:

$$p(\hat{Y}_i | X_i) = \arg \max_{\substack{c \\ c \neq \hat{Y}_i}} p(Y_i=c | X_i) = \begin{cases} \text{large if no doubt} \\ \approx 0 \text{ if uncertain} \end{cases}$$



- the true posterior $p^*(Y_i | X_i)$ is called the "Bayesian ~~decision~~ posterior". The Bayes decision rule

$$\hat{Y}_i = \arg \max_c p^*(Y_i=c | X_i)$$

is the theoretically best possible classifier

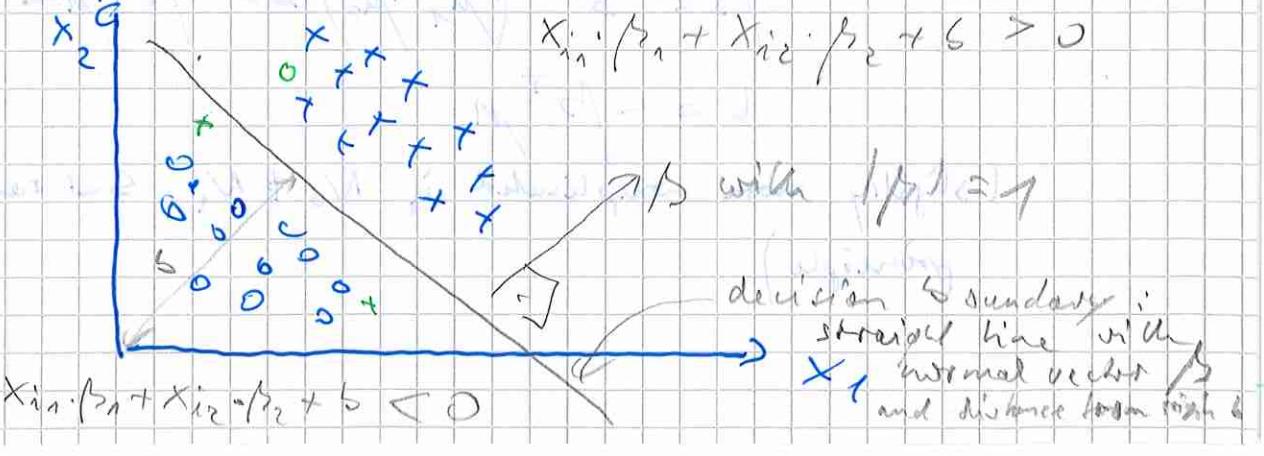
- \Rightarrow learning attempts to ensure $p^*(Y_i | X_i) \approx p^*(Y_i | X_i)$
- the true Bayesian posterior and thus the limit is generally unknown in practice, but plays an important role as an upper bound for the error in theoretical analysis

Explaining confusion matrix, difference between false positives and false negatives → see later

Linear Models

- simplest model family, applicable to
 - regression : $Y_i \in \mathbb{R}$
 - classification $Y_i \in \{1, \dots, C\}$ $Y_i \in \{0, 1\}$

(extensions to more classes possible by coupling several linear models in one-against-the-rest fashion)
- complete a weighted sum of the feature values plus offset
 - matrix notation for regression $\hat{Y}_i = \vec{x}_i \cdot \beta + b$
 - \vec{x}_i : column vector of D weights
 - b : offset
 - element notation $\hat{Y}_i = \sum_{j=1}^D x_{ij} \cdot \beta_j + b$
 - β_j : weight of feature j
 - $|\beta_j| \begin{cases} \text{big} & \text{feature } j \text{ is important} \\ \approx 0 & \text{--" -- unimportant} \end{cases}$
- matrix notation for classification $\hat{Y}_i = \Pi [x_i \cdot \beta + b > 0]$
 - $\Pi [cond] = \begin{cases} 1 & \text{if "cond" is true} \\ 0 & \text{if "cond" is false} \end{cases}$ Simply converts true/false to members 0/1
 - $\Rightarrow x_i \cdot \beta + b$ should be positive if $\hat{Y}_i^* = 1$
negative if $\hat{Y}_i^* = 0$
 - the parameters are $\Theta = \{\beta, b\}$ (D+1 numbers)
- linear classification works well if data are linearly separable, or nearly so



- many ways to find good decision planes, e.g.

optimal parameters $\hat{\beta}, \hat{b}$, controlled by the choice of loss function.

- If the data are suitable for linear classification, all solutions will be similar.

- squared loss: $\Rightarrow \text{LDA (linear discriminant analysis)}$
 define $\tilde{Y}_i = \begin{cases} 1 & \text{if } Y_i = 1 \\ -1 & \text{if } Y_i = 0 \end{cases} = 2Y_i - 1$

$$\text{optimize } \hat{\beta}, \hat{b} = \underset{\beta, b}{\operatorname{arg\min}} \frac{1}{N} \sum_{i=1}^N (\tilde{Y}_i - x_i \beta - b)^2$$

reduces learning to ordinary least squares and looks complicated, but has a simple solution

(assuming balanced classes, i.e. $N_0 = N_1 = \frac{N}{2}$)

- 1) center compute the class means

$$\mu_1 = \frac{1}{N_1} \sum_{i: Y_i=1} x_i \quad \mu_0 = \frac{1}{N_0} \sum_{i: Y_i=0} x_i$$

- 2) and total mean

$$\mu = \frac{1}{N} \sum_i x_i = \frac{\mu_0 + \mu_1}{2} \quad (N_0 = N_1)$$

- 3) compute within-class covariance matrix

$$S = \frac{1}{N} \sum_i (x_i - \mu_i)(x_i - \mu_i)^T \quad \text{outer product}$$

with $\mu_i = \begin{cases} \mu_1 & \text{if } Y_i = 1 \\ \mu_0 & \text{if } Y_i = 0 \end{cases}$ cluster mean

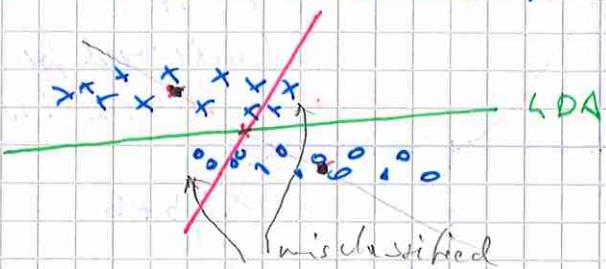
- 4) compute parameters

$$\beta = S^{-1}(\mu_1 - \mu_0)$$

$$b = -\beta^T \mu$$

(slightly more complicated if $N_0 \neq N_1$, same principle)

- Comparison between LDA and a naive rule: choose 11
vector between class means



normal vector of decision boundary

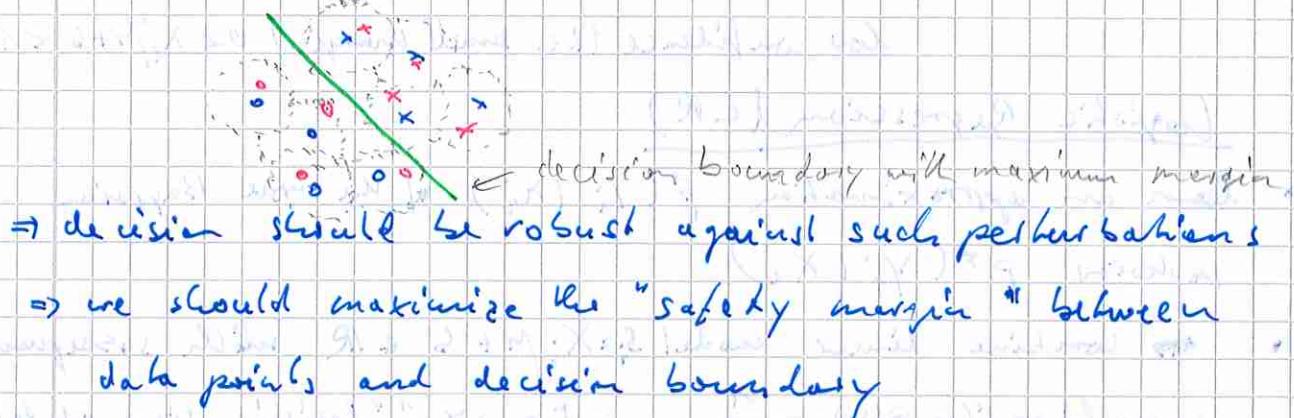
$$\beta = S^{-1}(\mu_1 - \mu_0)$$

rotates decision boundary according to cluster shape
(no rotation if clusters are round)

- strong assumptions:
 - cluster shapes are close to Gaussian bell (elliptic)
 - both classes have same cluster shapes
 - somewhat robust if not fulfilled exactly

Linear Support Vector Machine (SVM)

- insight: features are not exact, but noisy \Rightarrow slightly perturbed TS is equally plausible



$$\hat{\beta}, \hat{b} = \arg \min_{\beta, b} \frac{\beta^T \beta}{2} + \lambda \sum_i \text{Hinge Loss}(Y_i^*, X_i \beta + b)$$

maximizes margin for possible TS perturbations
"regularization term"

measures quality of prediction on the TS
"data term"

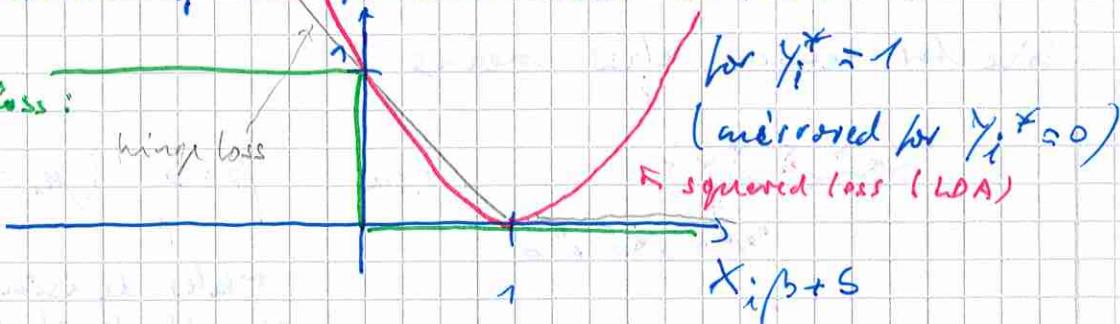
regularization parameter: adjusts trade-off between training error and robustness
(hyperparameters, must be chosen by user)

$$\text{Hinge Loss}(Y_i^*, X_i \beta + b) = \begin{cases} \max(0, 1 - (X_i \beta + b)) & \text{if } Y_i^* = 1 \\ \max(0, 1 + (X_i \beta + b)) & \text{if } Y_i^* = 0 \end{cases}$$

no analytic solution, but efficient iterative algorithms

- comparison of loss functions seen so far!

0/1 loss:



- 0/1 loss: no penalty if correct sign (i.e. correct decision)
constant penalty if wrong sign

- squared loss: quadratic penalty if \$x_i\beta + \varsigma \neq 1\$
(i.e. over confident responses are also penalized)

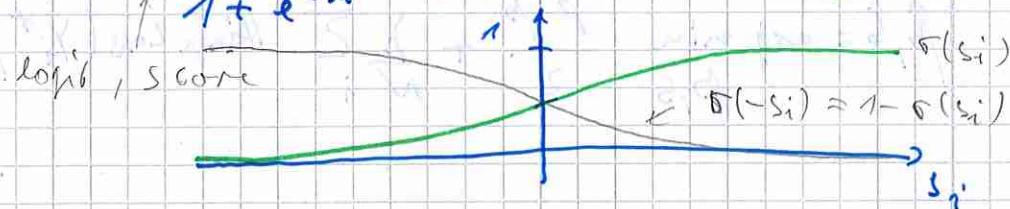
- hinge loss: no penalty for sufficiently confident correct answers: \$x_i\beta + \varsigma \geq 1\$

linearly increasing penalty for incorrect answers: \$x_i\beta + \varsigma < 0\$

small penalty for correct answers with low confidence (i.e. small margin) \$0 < x_i\beta + \varsigma < 1\$

Logistic Regression (LR)

- learn an approximation \$\hat{p}(Y_i=1|x_i)\$ of the true Bayesian posterior \$p^*(Y_i=1|x_i)\$
- combine linear model \$s_i = x_i\beta + \varsigma \in \mathbb{R}\$ with non-linearity \$\sigma(s_i) = \frac{1}{1+e^{-s_i}} \in [0, 1] \Rightarrow\$ "logistic sigmoid fn."



$$\hat{p}(Y_i=1|x_i) = \sigma(x_i\beta + \varsigma)$$

$$\begin{aligned} \hat{p}(Y_i=0|x_i) &= 1 - \hat{p}(Y_i=1|x_i) = 1 - \sigma(x_i\beta + \varsigma) \\ &= \sigma(-(x_i\beta + \varsigma)) \end{aligned}$$

- optimal decision: \$Y_i = \begin{cases} 1 & \text{if } \hat{p}(Y_i=1|x_i) > \hat{p}(Y_i=0|x_i) \Leftrightarrow x_i\beta + \varsigma > 0 \\ 0 & \text{otherwise} \end{cases}\$

\$\Rightarrow\$ posterior is non-linear, but decision rule is linear

- optimization problem: define $\hat{\beta}$, i.e. $\hat{\beta}$ and \hat{b} such that
the combined posterior of the training labels is maximized
⇒ maximum likelihood principle
due to i.i.d.

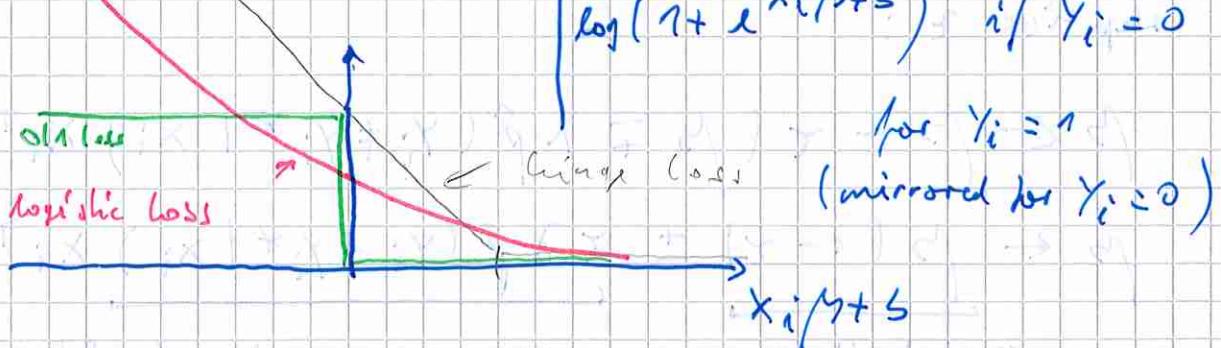
$$\begin{aligned}\hat{\beta}, \hat{b} &= \underset{\beta, b}{\operatorname{argmax}} \prod_{i: Y_i=1} p^*(Y_i=1|x_i) \cdot \prod_{i: Y_i=0} p^*(Y_i=0|x_i) \\ &= \underset{\beta, b}{\operatorname{argmax}} \prod_{i: Y_i=1} \sigma(x_i \beta + b) \cdot \prod_{i: Y_i=0} \sigma(-(x_i \beta + b))\end{aligned}$$

- after some manipulations and introduction of regularization (i.e. maximization of safety margin), this results in

$$\hat{\beta}, \hat{b} = \underset{\beta, b}{\operatorname{argmin}} \frac{\beta^T \beta}{2} + \frac{1}{n} \sum_i \text{logistic Loss}(Y_i, x_i \beta + b)$$

solve iteratively, similar to SVM

$$\text{logistic Loss}(Y_i, x_i \beta + b) = \begin{cases} \log(1 + e^{-(x_i \beta + b)}) & \text{if } Y_i = 1 \\ \log(1 + e^{x_i \beta + b}) & \text{if } Y_i = 0 \end{cases}$$



⇒ logistic loss is a smooth version of the hinge loss,
similar behavior

- for large dataset ($N > 10^4 \dots 10^5$), stochastic gradient descent (SGD), or one of its improved variants, is the fastest optimization algorithm

(for small datasets, $N \leq 10^3 \dots 10^4$, Newton or quasi-Newton algorithms are faster)

- SGD selects a single random instance or a mini-batch to estimate the loss in every iteration
 $i \sim \text{uniform}(1 \dots N)$

loss in current iteration:

$$\text{loss}_i = \frac{\beta^T \beta}{2} + \lambda \log \left(1 + e^{-\hat{y}(x_i \beta + s)} \right)$$

with - if $y_i = 1$ and + if $y_i = 0$

$$\begin{aligned}\frac{\partial \text{loss}_i}{\partial \beta} &= \frac{\beta}{2} + \frac{\lambda}{1 + e^{-\hat{y}(x_i \beta + s)}} \cdot e^{-\hat{y}(x_i \beta + s)} \cdot (-x_i^T) \\ &= \frac{\beta}{2} + \lambda \cdot \underbrace{\delta(-\hat{y}(x_i \beta + s))}_{= 1 - \hat{p}(y_i = y_i^* | x_i)} x_i^T \\ &= \frac{\beta}{2} + \lambda \cdot \hat{p}(y_i \neq y_i^* | x_i) x_i^T \\ &= \beta + \lambda \cdot \hat{p}(y_i \neq y_i^* | x_i) x_i^T\end{aligned}$$

Since we want to minimize the loss, we perform gradient descent, i.e. the update goes in the opposite direction with step size τ :

$$\begin{aligned}\beta &\leftarrow \beta - \tau (\beta + \lambda \hat{p}(y_i \neq y_i^* | x_i) x_i^T) \\ \beta &\leftarrow \underbrace{\beta (1 - \tau)}_{\text{regularization}} \pm \underbrace{\tau \lambda \hat{p}(y_i \neq y_i^* | x_i) x_i^T}_{\text{update weight}}\end{aligned}$$

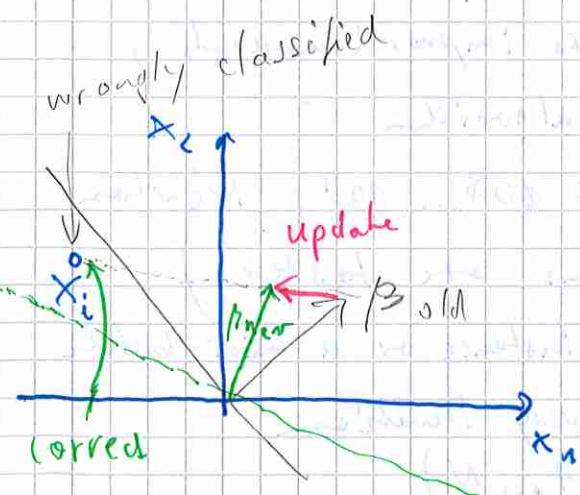
regularization:
reduce magnitude
of β

update weight:
if classification
is already correct:
 $\hat{p}(y_i \neq y_i^* | x_i) \approx 0$
 \Rightarrow almost no update

- if classification is wrong:

$$\hat{p}(y_i \neq y_i^* | x_i) \approx 1$$

\Rightarrow move β towards x_i (if $y_i^* = 1$)
or away from x_i (if $y_i^* = 0$)



Unbalanced loss functions / classes

- in two-class problems, there are two kinds of errors:
 - false positive $\hat{Y}_i = 1, Y_i^* = 0$
 - false negative $\hat{Y}_i = 0, Y_i^* = 1$
 - often, consequences of these errors are very different,
e.g. access control to a bank safe
 - false negative: authorized person is rejected:
 \Rightarrow annoying, but simply retry (cf. fingerprint sensor on cell phones)
 - false positive: burglar is admitted
 \Rightarrow very expensive
- \Rightarrow weight each loss term with its cost

$$\hat{\beta}, \hat{\gamma} = \underset{\beta, \gamma}{\operatorname{arg\,min}} \frac{1}{N} \sum_i w_i \text{loss}(Y_i^*, \hat{Y}_i | \beta, \gamma)$$

for access control:

$$w_i = \begin{cases} \text{big} & \text{if } Y_i^* = 0 \text{ (false positive)} \\ \text{small} & \text{if } Y_i^* = 1 \end{cases}$$

choose weights according to application

- unbalanced classes: one class (e.g. positive = disease) is much more rare than the other (e.g. healthy)
 - \Rightarrow errors may be misleading
 - consider test for a disease with 1% false positives and 1% false negatives - sounds pretty good
 - positive to negative ratio is $1/100$
 - \sim simple calculation shows that if
 - $\hat{Y}_i = 1$ (suspected disease) $\rightarrow p(Y_i^* = 1 | x_i) \approx 0.5$
 (actual disease)
 - \Rightarrow many patients are unnecessarily scared
 - \Rightarrow report precision and recall instead of error rates

- even worse: if only true positive rate (e.g. 99%)
false negative rate (e.g. 1%) is reported, but not the false positive rate,
the test quality cannot be judged at all!
(cf. blood breast cancer blood test scandal)

- confusion matrix for $C=2$ ($C > 2$ analogously)

\hat{Y}_i^*	0	1	should be diagonal
0	#TN	#FN	
1	#FP	#TP	

true positive rate \hat{s} sensitivity:
 \hat{s} recall

$$\frac{\# TP}{\# P} = \frac{\# TP}{\# TP + \# FN}$$

true negative rate \hat{s} specificity:

$$\frac{\# TN}{\# N} = \frac{\# TN}{\# TN + \# FP}$$

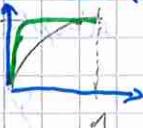
precision:

$$\frac{\# TP}{\# TP + \# FP}$$

- If we have a score or posterior, we can freely choose the threshold, e.g. $\hat{p}(Y_i=1|X_i) > 0.5$ (balanced)
 $\hat{p}(Y_i=1|X_i) > 0.9$ (avoid false positives)
 \Rightarrow precision, recall, specificity change according to threshold

\Rightarrow ROC curve and precision/recall curve depict this dependency \Rightarrow judge quality of results by "area under curve"
(AUC)

fundamental insight 2:



DEFINING APPROPRIATE LOSS FCT. FOR AN APPLICATION IS A MAJOR PART OF ▲ MACHINE LEARNING DESIGN.

\Rightarrow easier than hand-crafted model, because only a critic, not a constructive solution

Regression

- classification: Y is discrete (e.g. C classes)
- regression: Y is a real number
- closely related: $s_i = x_i \beta + \epsilon$ is a real-valued score
 $\hat{y}_i = \text{sign}(x_i \beta + s)$ a discrete class
 \Rightarrow regression is often an intermediate part of classification
and often a goal of its own
(examples:)

- linear regression is most basic approach
- \Rightarrow ordinary least squares (OLS) method

$$\hat{\beta}, \hat{s} = \underset{\beta, s}{\arg \min} \frac{1}{N} \sum_{i=1}^N (y_i^* - (x_i \beta + s))^2$$

y_i^*
arbitrary real number

~~$\hat{\beta}, \hat{s}$~~

- we can eliminate s by centralizing the data

$$\tilde{x}_i = x_i - \mu \quad \mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \text{mean}$$

(\tilde{x}_i has now zero mean)

$$\hat{\beta} = \underset{\beta}{\arg \min} \frac{1}{N} (y^* - x \beta)^T (y^* - x \beta)$$

(matrix notation)

- find the optimum by setting the derivative to zero

$$\frac{\partial}{\partial \beta} \frac{1}{N} (y^* - x \beta)^T (y^* - x \beta) = \frac{2}{N} (-y^* x + x^T x \beta) \stackrel{!}{=} 0$$

$$\frac{\partial}{\partial \beta} \frac{1}{N} (y^* - x \beta)^T (y^* - x \beta) = \frac{2}{N} (-x^T y^* + x^T x \beta) \stackrel{!}{=} 0$$

$$x^T x \beta = x^T y^*$$

$$\hat{\beta} = (X^T X)^{-1} X^T y^*$$

$\approx X^+$ pseudo-inverse

Moore-Penrose inverse
generalization of inverse for rectangular matrices

- in order to give all features equal influence in the loss, it makes sense to standardize features \Rightarrow unit-free quantities

$$x_i \approx \frac{x_i - \mu}{\sigma} \leftarrow \text{element-wise division}$$

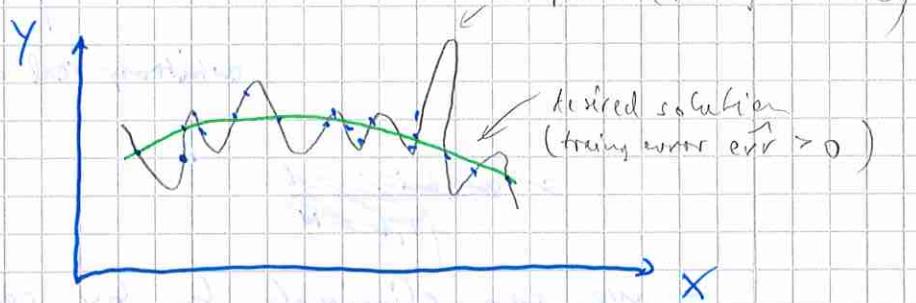
σ : vector of standard deviations (σ_j : std dev of feature j)

(this is also helpful in classification, standard data preprocessing)

- numeric solution of OLS: reusable algorithms in all ML libraries, using QR decomposition or singular value decomposition (SVD)

overfitting

- non-linear example



- effective counter-measure: regularization of parameters, e.g. penalize large values of β -coefficients

most common: L2 regularization: restrict Euclidean norm squared

$$\|\beta\|_2^2 = \sum_j \beta_j^2 = \beta^T \beta$$

\Rightarrow ridge regression

$$\hat{\beta} = \underset{\beta}{\operatorname{arg\,min}} \frac{\beta^T \beta}{2} + \frac{\lambda}{n} \sum_i (y_i^* - x_i^T \beta)^2$$

identical to support vector machine, but with squared loss \Rightarrow robustness against noise

- L1 regularization: restrict L1 norm

$$\|\beta\|_1 = \sum_j |\beta_j|$$

$\Rightarrow \text{LASSO regression}$

$$\hat{\beta} = \underset{\beta}{\operatorname{arg\,min}} \| \beta \|_1 + \frac{\lambda}{n} \sum_i (y_i^* - x_i \beta)^2$$

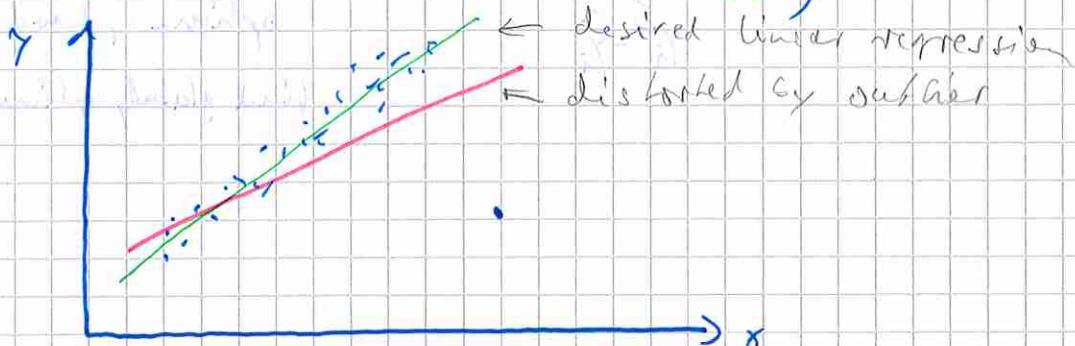
\Rightarrow sparsity enforcement: many coefficients $\beta_j = 0$
non-zero coefficients denote important features

\Rightarrow feature selection

- regularization parameter λ must be manually chosen for each application as a hyper-parameter to optimize the bias-variance trade-off

robust regression

- results can be severely distorted when training data is contaminated by outliers \Leftrightarrow grossly wrong data points
(e.g. satellite data: transmission errors as opposed to measurement noise)



- solution 1: identify outliers and ignore outliers

\Rightarrow RANSAC algorithm

- solution 2: modify the loss function to become robust against outliers

- squared loss: outliers are highly weighted, because error term is squared

- absolute loss: $\text{loss}(y_i^*, \hat{y}_i) = |y_i^* - \hat{y}_i|$

- much more robust, influence of outliers greatly reduced

- disadvantages:

- often no unique solution for $\hat{\beta}$

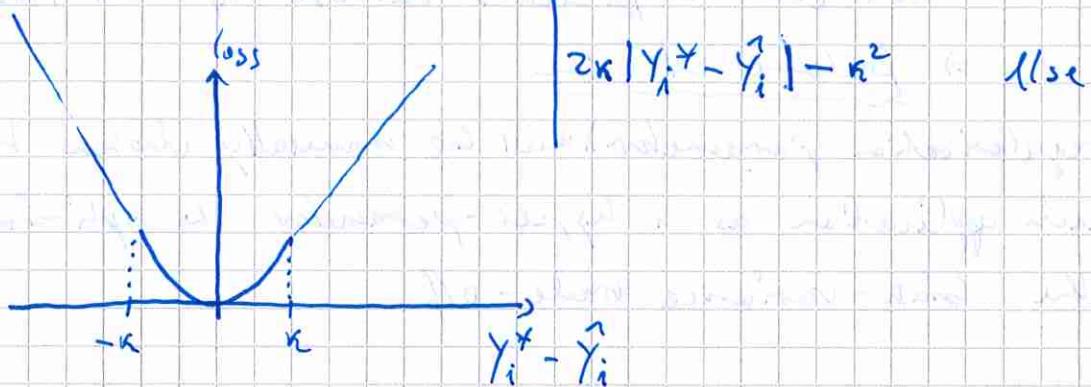
- higher generalization error on new data

- Huber Loss: best of both worlds

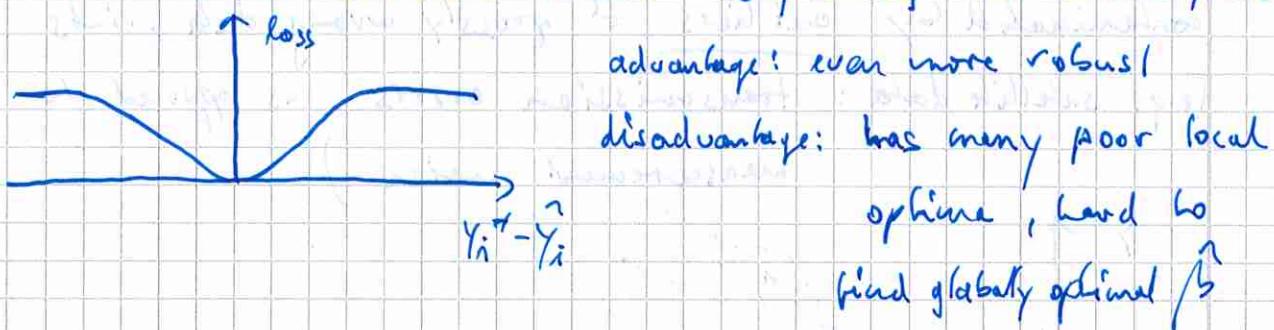
- squared loss for small errors \Rightarrow high accuracy, uniqueness

- absolute loss for high errors \Rightarrow robustness

$$\text{Huber Loss}(Y_i^*, \hat{Y}_i) = \begin{cases} (Y_i^* - \hat{Y}_i)^2 & \text{if } |Y_i^* - \hat{Y}_i| \leq \kappa \\ 2\kappa|Y_i^* - \hat{Y}_i| - \kappa^2 & \text{else} \end{cases}$$



- biweight functions: loss for high errors saturates,
there is no distinction between "very wrong" and "extremely wrong"



advantage: even more robust!

disadvantage: has many poor local optima, hard to find globally optimal β

- outlier handling $\hat{\approx}$ out-of-distribution detection \approx novelty detection
 $\hat{\equiv}$ robustness against adversarial attacks / samples
 is still a hot research topic \Rightarrow show adversarial examples

Non-linear Methods

- most data are not adequately represented by linear models
- approach 1: manually design non-linear model and fit its parameters, e.g. non-linear least squares, Levenberg - Marquardt algorithm.
problems:
 - has many poor local optima, good initial guess required to find best $\hat{\theta}$
 - manual model design is difficult, esp. for highly complex problems (e.g. medicine, biology, social sciences)

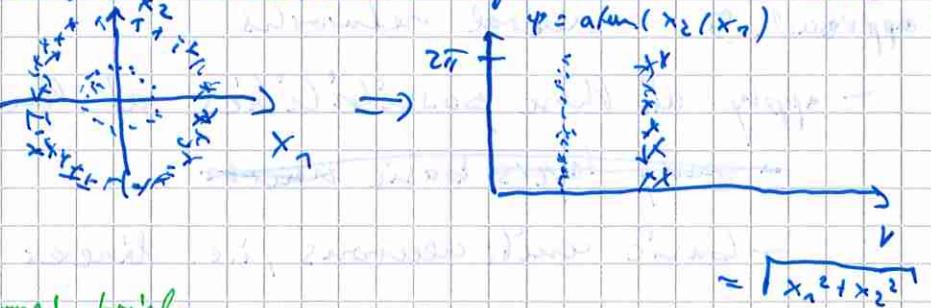
- approach 2: ~~clearly combine several linear models into a non-linear model~~
- approach 2: augmented feature spaces

(1) calculate new features by non-linear functions of the given original features

(2) apply a linear model in this augmented feature space

example: body mass index = ~~height~~ weight / height²

example: Cartesian coordinates \Rightarrow polar coordinates



Kernel trick

variant: rewrite loss and regularization in terms of a

non-linear kernel function \Rightarrow augmented feature

space needs not be constructed explicitly, can have

high and even infinite dimension

advantage: - training (i.e. finding optimal $\hat{\theta}$) is easier than in approach 1, good optimization algorithms

disadvantage: - manual design of good non-linear transformations or

- kernels is difficult

• approach 3: boosting

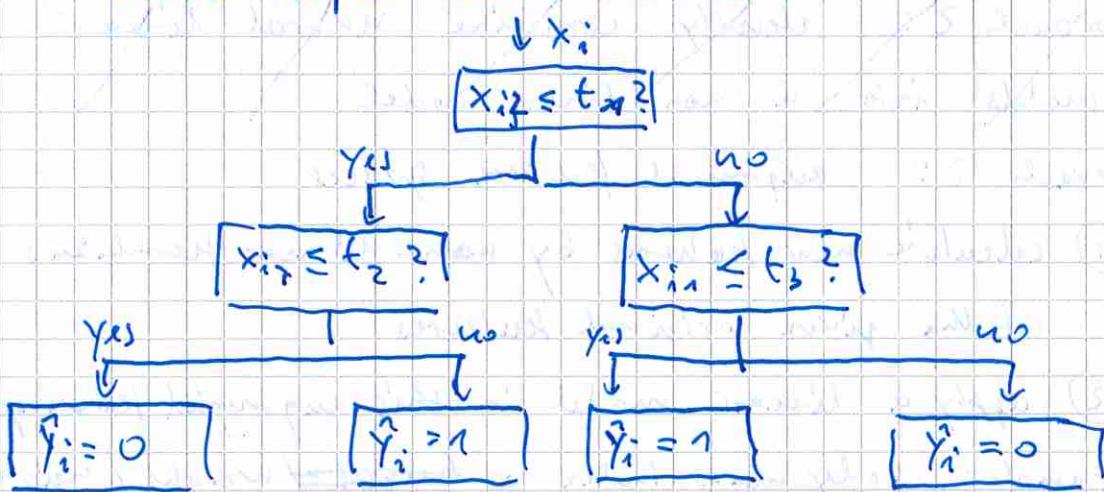
- run several linear models in parallel, then take a weighted vote

$$\hat{Y}_i = \text{sign} \left(\frac{1}{L} \sum_{e=1}^L w_e \cdot \text{sign}(X_i \beta_e + b_e) \right)$$

$$\Theta = \{ (\beta_1, b_1, w_1), \dots, (\beta_L, b_L, w_L) \}$$

• approach 4: decision trees

- run several linear models in series \Rightarrow split feature space into subregions, where each subregion behaves simpler than the whole.



• approach 5: neural networks

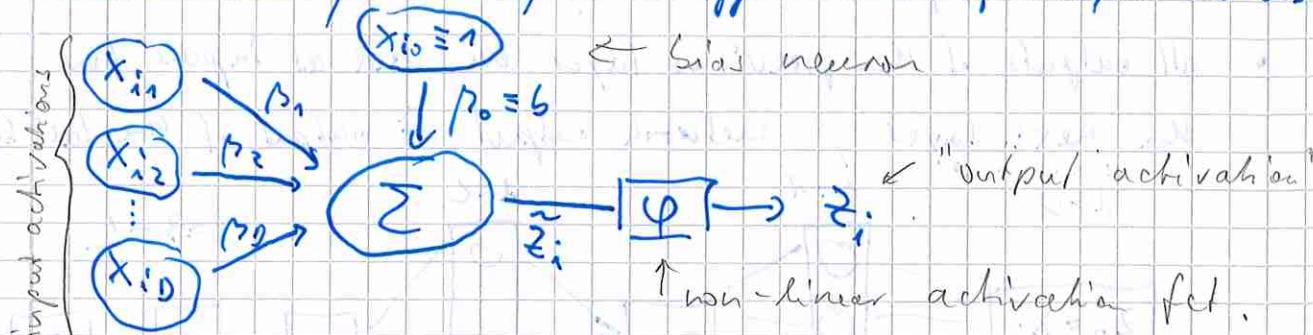
- apply all those possibilities at the same time
 - ~~many layers basic neurons~~
- basic units: neurons, i.e. linear models, followed by simple non-linearity
- sequential processing: arrange neurons in layers
- parallel processing: each layer contains many independent neurons

why does this solve the problems?

- scales to very complex models ("deep" networks, "wide" networks)
- outputs of interior layers are "learned orthonormal feature spaces"
- most local minima are reasonably good \Rightarrow robust optimization
- hand-crafting good architecture and non-linearities is still hard, but performance is relatively robust against suboptimal solutions

Neural Networks

- neurons are simple linear models, followed by a 1-dimensional non-linearity — very crude approximation of biological neurons



bias neuron: add feature 0, which is constant = 1
 \Rightarrow absorb bias parameter b into β as β_0
 $(\beta \in \mathbb{R}^{D+1}, j = 0 \dots D)$

output of the neuron:

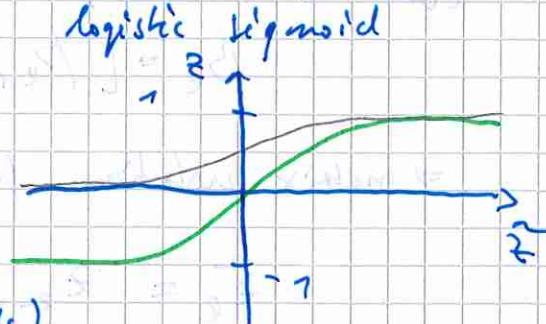
$$z_i = \varphi \left(\sum_{j=0}^D x_{ij} \cdot \beta_j \right) \quad \text{row vector}$$

- activation fct. should be almost everywhere differentiable,
 so that we can train by gradient descent,
 $(\text{sign}(z) \text{ is not suitable})$

- traditional choices:

$$z = \varphi(\tilde{z}) = \frac{1}{1 + e^{-\tilde{z}}} \quad \text{logistic sigmoid}$$

$$z = \varphi(\tilde{z}) = \tanh(\tilde{z})$$



$z = \tilde{z}$ (for linear layers,
 e.g. regression outputs)

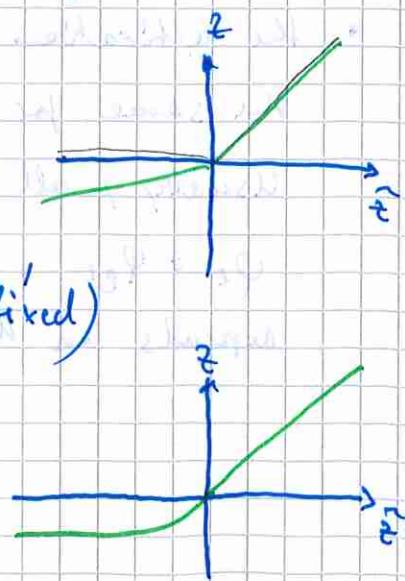
- modern choices

$$z = \varphi(\tilde{z}) = \text{ReLU}(\tilde{z}) = \max(\tilde{z}, 0)$$

$$z = \text{LeakyReLU}(\tilde{z}) = \begin{cases} \tilde{z} & \text{if } \tilde{z} \geq 0 \\ a \cdot \tilde{z} & \text{if } \tilde{z} < 0, \\ (\text{a} \ll 1 \text{ fixed}) \end{cases}$$

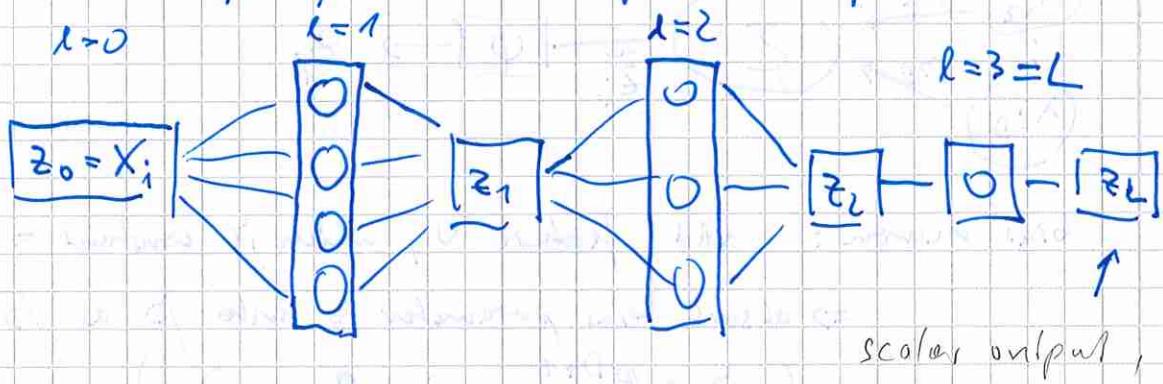
$z = \text{PReLU}(\tilde{z})$: like leaky ReLU,
 but a is learned

$$z = \text{ELU}(\tilde{z}) = \begin{cases} \tilde{z} & \text{if } \tilde{z} \geq 0 \\ a(e^{\tilde{z}} - 1) & \text{otherwise} \\ (\text{a fixed}) \end{cases}$$



Fully connected networks

- consist of the input layer ($\hat{=}$ the original feature vector x_0) and L ~~compu-~~ neuron layers
 - all outputs of the previous layer are used as inputs for the next layer ; network output $\hat{=}$ output of the last layer



e.g. regression in 1-D
 • posterior probab. $\hat{p}(\vec{Y}_i=1 | \vec{X}_i)$

- layer l contains H_l neurons (+ bias neuron - not shown)
 (here $H_1 = 4$, $H_2 = 3$, $H_3 = 1$)
 \Rightarrow each layer needs H_l different weight vectors / β :
 $\beta_{lm} \quad l = 1 \dots L \quad m = 1 \dots H_e^{M_l}$

(combine the weight vectors of each layer into matrix)

$$\beta_e = [\beta_{e1}, \dots, \beta_{eM_e}]$$

\Rightarrow matrix notation for pre-activations of layer l :

$$\tilde{z}_\ell = z_{\ell-1} \cdot \beta_\ell$$

- the activation functions (σ (non-linearities)) are usually the same for all neurons in a given layer: φ_e .
Usually, all interior layers have the same activation
 $\varphi_e = \varphi_{e1}$, except for last layer φ_L , whose activation depends on the application

$$\hat{y}_i = \hat{z}_3 = \varphi_3 \left(\varphi_2 \left(\varphi_1 \left(z_0 \cdot \beta_1 \right) \cdot \beta_2 \right) \cdot \beta_3 \right)$$

↓ ↑ ↓ ↓ ↓
 non-linearities parallel processing
 serial processing

[Note: the bias channel $z_{l=0} \equiv 1$ is always implicitly present, even if not shown explicitly]

- By increasing L and M_L , we can adapt model complexity to the application
- output activation functions depend on the task :
 - regression : y_i is real or real-valued vector
 $\varphi_L(\tilde{z}_L) = \tilde{r}_L$ linear activation
 \Rightarrow last layer performs linear regression, using penultimate activations \tilde{z}_{L-1} as augmented feature space
 - classification with $C=2$: estimate posterior for positive class
 $z_L \stackrel{\text{def}}{=} \hat{p}(Y_i = 1 | x_i) = \sigma(\tilde{z}_L)$ sigmoid fct.
 $[\hat{p}(Y_i = 0 | x_i) = 1 - z_L \text{ is trivial, doesn't have to be learned}]$

[class label easily obtained from posterior via

$$\hat{Y}_i = \begin{cases} 1 & \text{if } z_L \geq 0.5 \\ 0 & \text{if } z_L < 0.5 \end{cases}$$

(or other threshold for asymmetric costs)

- classification with $C \geq 2$: z_L is a probability vector of length C , i.e.

$$z_{L,k} = \hat{p}(Y_i = k | x_i)$$

with normalization $\sum_{k=1}^C z_{L,k} = 1$

softmax activation generalizes sigmoid for multi-class problem (reduces to $z_{L1} = \sigma(\tilde{z}_{L1}, 2)$ if $C=2$)

$$\tilde{z}_{Lk} = \frac{e^{\tilde{z}_{Lk}}}{\sum_{k=1}^C e^{\tilde{z}_{Lk}}}$$

- illustration of how it works, using a 2-layer network (with sigmoid activation (not suitable for gradient training, but easy to understand))

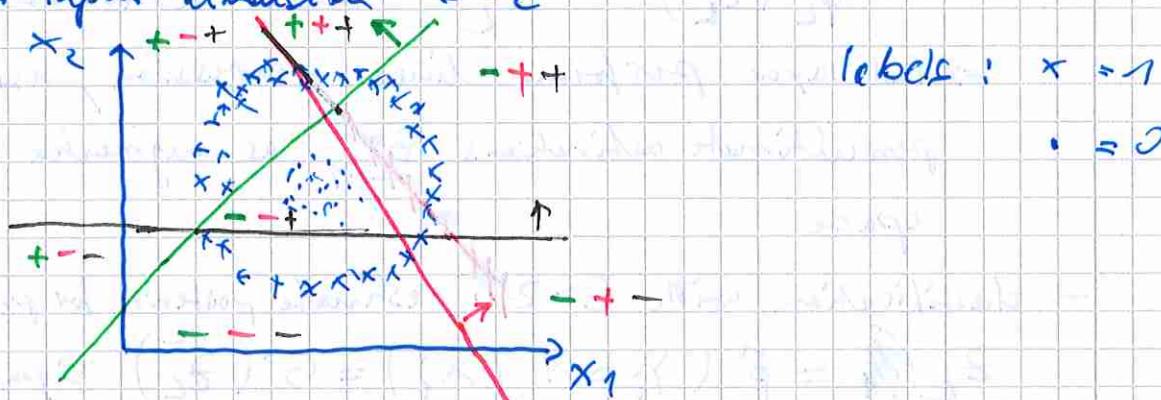
We have

$$z_2 = \varphi_2 (\varphi_1 (x_i \cdot B_1) \cdot B_2)$$

with $\varphi_{1,2}(\tilde{z}) = \text{sign}(\tilde{z}) = \{-, +\}$

- we use $M_1 = 3$, $M_2 = 1$ (4 neurons in total)

for input dimension $D = 2$



first layer neurons define 3 decision planes

with arrows indicating the positive side

\Rightarrow each subregion is assigned a particular sign combination

(7 subregions in total), activations \tilde{z}_1 are corners



Subregions are learned such

that each contains a unique label.



second layer neuron must learn a decision plane that cuts off the corner $(-1, -1, +1)$ from the remaining corners of the cube,

$$\text{e.g. } \beta_2 = (2, 1, 1, 1, -1) \\ \text{bias weight}$$

$$\text{with } \hat{Y}_i = \begin{cases} 1 & \text{if } z_{11i} + z_{12i} + z_{13i} + 2 > 0 \\ 0 & \text{else} \end{cases}$$

- loss functions are essentially the same as for linear models
 - regression: squared loss $\text{Loss}(Y_i^*, \hat{Y}_i) = \|Y_i^* - \hat{Y}_i\|_2^2$
 - classification with $C=2$:

$$z_{Li} = \hat{p}(Y_i = 1 | x_i) \quad \begin{matrix} \text{scalar posterior prob. value} \\ (\text{sigmoid output}) \end{matrix} \\ (\hat{p}(Y_i = 0 | x_i) = 1 - z_{Li} \text{ trivial})$$

\Rightarrow logistic loss

$$\text{Loss}(Y_i^*, z_{Li}) = \begin{cases} -\log z_{Li} & \text{if } Y_i^* = 1 \\ -\log(1 - z_{Li}) & \text{if } Y_i^* = 0 \end{cases}$$

- classification with $C > 2$

$$z_{Lik} = \hat{p}(Y_i = k | x_i) \quad \begin{matrix} \text{full probability vector} \\ (\text{softmax output}) \end{matrix}$$

$$\text{with } \sum_k z_{Lik} = 1$$

\Rightarrow cross-entropy loss

$$\text{Loss}(Y_i^*, z_{Lik}) = -\log z_{Lik} \quad \text{with } k = Y_i^*$$

(reduces to logistic loss for $C=2$)

$\hat{p}(Y_i = Y_i^* | x_i)$ should be close to 1 (and the others close to 0)

$$-\log z_{Lik, k=Y_i^*} = \begin{cases} \approx 0 & \text{if } z_{Lik, k=Y_i^*} \approx 1 \\ \infty & \text{else} \end{cases}$$

training by stochastic gradient descent

- (1) pick a random instance i
- (2) compute $\text{Loss}(Y_i^*, \tilde{z}_{li})$ for current parameters $\Theta^{(t)}$
- (3) compute $\Delta\Theta = \frac{\partial \text{Loss}(Y_i^*, \tilde{z}_{li})}{\partial \Theta}$ (derivative of loss w.r.t. parameters)

$$(4) \text{ update } \Theta^{(t+1)} = \Theta^{(t)} - \tau \Delta\Theta$$

↑ learning rate

Subtraction, because we need descent to minimize loss

- Step (3) is the critical part. The derivative is computed by back-propagation, i.e. by passing errors from output towards input.

Introduce auxilliary variable $\tilde{\delta}_e$ (row vector)

$$\tilde{\delta}_e = \frac{\partial \text{Loss}}{\partial \tilde{z}_e} \quad \text{with } \tilde{z}_e = z_{e:} \cdot B_e \quad \text{the}$$

pre-activations of layer e
(before applying non-linearity $\varphi_e(\cdot)$)

- The chain rule gives

$$\tilde{\delta}_e = \frac{\partial \text{Loss}}{\partial \tilde{z}_e} = \frac{\partial \text{Loss}}{\partial \tilde{z}_{e+1}} \cdot \frac{\partial \tilde{z}_{e+1}}{\partial \tilde{z}_e} \cdot \frac{\partial \tilde{z}_e}{\partial \tilde{z}_e}$$

The RHS is simple:

$$\frac{\partial \text{Loss}}{\partial \tilde{z}_{e+1}} =: \tilde{\delta}_{e+1} \quad \begin{array}{l} \text{by definition of } \tilde{\delta}. \text{ We already know this} \\ \text{term because we work back-to-front.} \end{array}$$

$$\frac{\partial \tilde{z}_{e+1}}{\partial \tilde{z}_e} = \frac{\partial (z_{e:} \cdot B_{e+1})}{\partial \tilde{z}_e} = B_{e+1}^T$$

$$\frac{\partial \tilde{z}_e}{\partial \tilde{z}_e} = \frac{\partial (\varphi_e(\tilde{z}_e))}{\partial \tilde{z}_e} = \text{diag}(\varphi'_e(\tilde{z}_e))$$

$$\tilde{\delta}_e = \tilde{\delta}_{e+1} \cdot B_{e+1}^T \cdot \text{diag}(\varphi'_e(\tilde{z}_e))$$

- The start of backpropagation is determined by the output activation

• regression: $\varphi_L(\hat{z}_L) = \hat{z}_L$

$$\text{Loss}(y_i^*, \hat{z}_{L,i}) = \frac{1}{2} (y_i^* - \hat{z}_{L,i})^2$$

$$= \frac{1}{2} (y_i^* - \hat{z}_{L,i})^2$$

$$\Rightarrow \frac{\partial \text{Loss}}{\partial \hat{z}_L} = \hat{f}_L = \hat{z}_{L,i} - y_i^*$$

• classification: $\varphi_L(\hat{z}_L) = \text{softmax}(\hat{z}_L)$

$$\text{Loss}(y_i^*, \hat{z}_L) = \text{cross-entropy}(y_i^*, \hat{z}_L)$$

$$\Rightarrow \left(\frac{\partial \text{Loss}}{\partial \hat{z}_L} \right)_k = \begin{cases} \hat{z}_{L,k} - 1 & \text{if } y_i^* = k \\ \hat{z}_{L,k} & \text{if } y_i^* \neq k \end{cases}$$

- The derivative w.r.t. the parameters B_e is easily computed from \hat{f}_e

$$\frac{\partial \text{Loss}}{\partial B_e^T} = \left(\frac{\partial \text{Loss}}{\partial \hat{z}_e} \right)^T \cdot \left(\frac{\partial \hat{z}_e}{\partial B_e} \right)^T = \hat{f}_e^T \cdot \hat{z}_{e-1}^T$$

$$\left[\frac{\partial \hat{z}_e}{\partial B_e^T} = \frac{\partial (\hat{z}_{e-1} \cdot B_e)}{\partial B_e^T} = \hat{z}_{e-1} \right]$$

⇒ gradient update with learning rate τ

$$B_e^{(t+1)} = B_e^{(t)} - \tau \cdot (\hat{z}_{e-1}^T \cdot \hat{f}_e^T)$$

outer product

- Fortunately, modern neural network libraries (pyTorch, tensorflow) contain a module "autograd", which computes these (and more complex) derivatives automatically ⇒ you don't have to do these calculations by hand and implement the resulting algorithm

⇒ eliminated big source of bugs ~~and allows for very complex and~~

Fundamental insight 3

NEURAL NETWORKS SUPPORT ~~MODERATE~~ NON-LINEAR MODELS OF

ALMOST ARBITRARY COMPLEXITY, BECAUSE:

- we can construct highly complex architectures by replicating a single basic unit (neuron \Rightarrow linear model plus non-linear activation) in parallel (layers) and in series (deep)
- auto grad allows to compute the gradients needed for training automatically, without additional programming
- modern GPUs make the forward and backward computation sufficiently fast.

Why do neural networks work so well?

- Theoretically not yet understood. NNs have so many parameters that they should overfit terribly (i.e. have small training error, but large generalization error). But this does not happen in practice.
- Theoretical guarantee: universal approximation theorem
 - 1-D regression: $x_i \in \mathbb{R}$
 - 2-class classification $\hat{y}_i = p^*(y_i=1|x_i) \in [0,1]$

can be shown, that a 2-layer network ($L=2$) can approximate arbitrarily complicated true mappings ($y_i^* = f^*(x_i)$ or $p^*(y_i=1|x_i)$) if the hidden layer is sufficiently wide ($M_1 \rightarrow \infty$) and the weights B_1, B_2 are suitably chosen. However, this is a pure existence proof, gives no algorithm to determine M_1, B_1, B_2 .
- Finding the globally optimal B_e for any network is NP-hard (i.e. virtually impossible)

IF THE NETWORK IS SUFFICIENTLY LARGE (WIDE AND DEEP), AND WE HAVE ENOUGH TRAINING DATA, TYPICAL LOCAL OPTIMA TEND TO BE PRETTY GOOD (E.G. CLOSE TO THE GLOBAL ONE, I.E. HAVE ~~HIGH~~ LOW GENERALIZATION ERROR).

This is an empirical finding, not yet theoretically understood, why this is the case.

Training tricks

- early stopping: regularly check generalization error during training on an independent validation set. Stop if this error estimate starts to increase (even if the training error would keep decreasing)
- training rate schedule: when training error plateaus, decrease training rate $\tau \leftarrow \frac{\tau}{10}$.
(repeat 2 or 3 times)
- mini-batch SGD: instead of minimizing the loss from a single random instance, as in pure SGD, use the average over a random set of size k ("mini-batch")
 $k = 32 \dots 10^4$, as GPU-RAM allows
- ADAM optimizer (and similar relatives): Plain SGD uses uniform learning rate τ for all parameters. ADAM adjusts learning rate for each parameter \Rightarrow much faster convergence (or divergence, if unlucky).

Let $g^{(t)} = \frac{\partial \text{loss}^{(t)}}{\partial \theta^{(t)}}$ the loss derivative in iteration t

$$\text{Plain SGD: } \theta^{(t+1)} = \theta^{(t)} - \tau g^{(t)}$$

ADAM maintains a running average of g and g^2 :

$$g^{(t)} = \mu_1 \tilde{g}^{(t-1)} + (1-\mu_1) g^{(t)}$$

$$(g^2)^{(t)} = \mu_2 (g^2)^{(t-1)} + (1-\mu_2) (g^{(t)})^2$$

$$(\mu_1 \approx 0.9, \mu_2 \approx 0.999)$$

$$\tilde{g}^{(t)} = \frac{g^{(t)}}{1-\mu_1^t} \quad (\tilde{g}^2)^{(t)} = \frac{(g^2)^{(t)}}{1-\mu_2^t}$$

"burn-in phase correction"

$$B^{(t+1)} = B^{(t)} - \gamma \frac{\tilde{g}^{(t)}}{\sqrt{(\tilde{g}^2)^{(t)}} + \epsilon} \quad \epsilon = 10^{-8} \quad (\text{avoid div-by-zero})$$

- batch normalization: normalizing pre-activations \tilde{z}_e

by their mean and std-dev within current mini-batch

- standard network $\tilde{z}_e = z_{e-1} \cdot B_e, z_e = \psi_e(\tilde{z}_e)$

- batch norm: $z_{e,i}^1 = z_{e-1,i} \cdot B_e$ for i in mini-batch

$$\mu = \frac{1}{k} \sum_i z_{e,i}^1 \quad \sigma = \sqrt{\frac{1}{M} \sum_i (z_{e,i}^1 - \mu)^2 + \epsilon}$$

$$z_{e,i}'' = \frac{z_{e,i}^1 - \mu}{\sigma}$$

$$\tilde{z}_{e,i} = a_e \cdot z_{e,i}'' + b_e$$

$$z_{e,i} = \psi_e(\tilde{z}_{e,i})$$

learnable parameters

often significant improvements of training and generalization errors

- dropout: deactivate part (e.g. 50%) of the network in every mini-batch:

sample dropout masks $\nu_e \in \{0, 1\}^{M_e}$ and replace activations $z_{e,i}$ with $z_{e,i} = \nu_{e,i} \cdot \psi_e(\tilde{z}_{e,i})$

\Rightarrow reduces overfitting, because subtle interactions between neurons are made impossible

during prediction, downscale weights B_e by dropout factor (e.g. $B_e \leftarrow B_e \cdot 0.1$) and predict without dropout

- weight initialization:

$$B_e \sim N(0, \sigma^2 = \frac{1}{M_{e-1} + 1}) \quad (\text{for tanh(.)})$$

$$\sim N(0, \sigma^2 = \frac{2}{M_{e-1} + 1}) \quad (\text{for ReLU(.)})$$

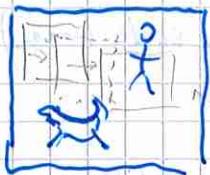
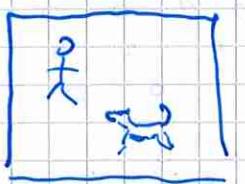
\Rightarrow ensures that norm of the input weights of each neuron expected (incl. bias weight) is 1.

- piecewise linear activation functions (ReLU, PReLU, Leaky ReLU) tend to work better than \tanh and \tanh' , because they do not suffer from vanishing gradients when $\tilde{x}_e \gg 0$.
- data augmentation: artificially increase training set size by applying random perturbations to the existing instances, e.g.:
 - mirror or rotate images
 - add slight noise
 - random non-rigid morphing
- \Rightarrow make training network robust against those perturbations (because y_i^* remains unchanged or is perturbed accordingly)
- except for data augmentation, preditioned modules for all these tricks are already available in NN libraries
- weight decay: add regularizer to loss that penalizes very large weights:
 - L_2 : $\lambda \sum_e \|B_e\|_F^2$
 - L_1 : $\lambda \sum_e \|B_e\|_1$ (elementwise abs)

1. *Leucosia* *leucostoma* (Fabricius) (Lepidoptera: Geometridae)

Convolutional Neural Networks

- in many applications, i.e. analysis of audio, images, videos, data are (approximately) translation invariant



human and dog

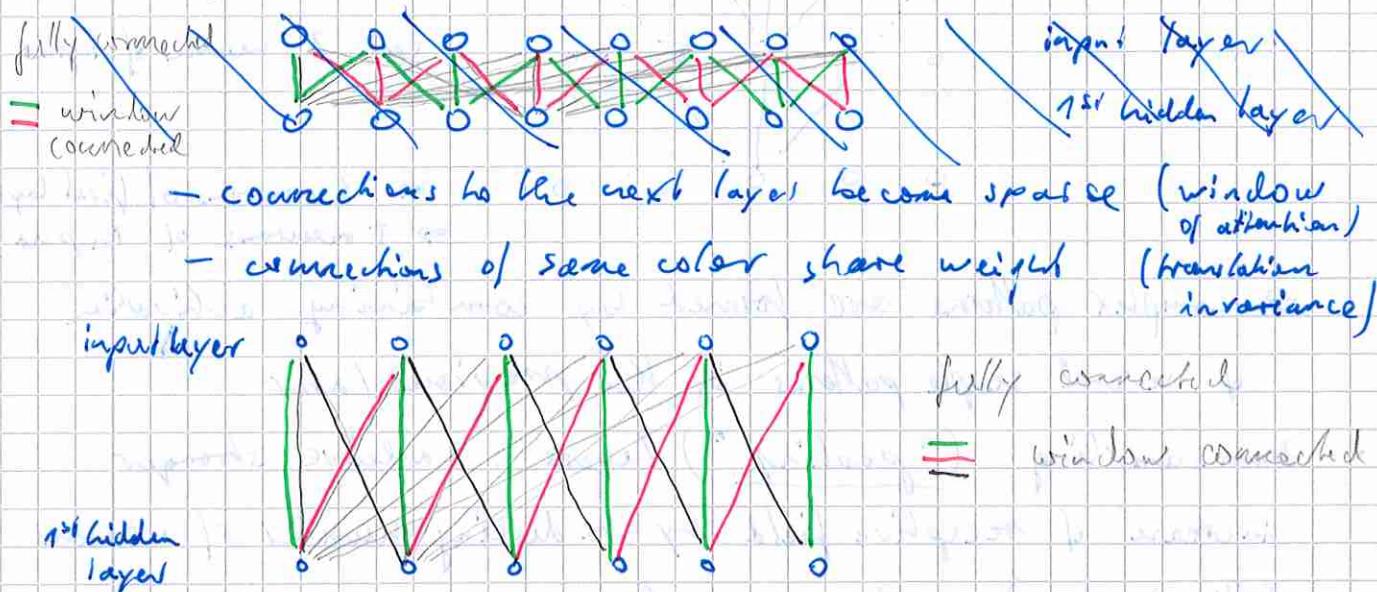
sliding window

\Rightarrow important receptive fields should not

\Rightarrow should not only look at image as a whole, but also at every sub-region of a certain (or various) size
"receptive field"

- these receptive fields should cover the entire image
in a sliding window fashion

- consequences for network architecture (1D illustration)



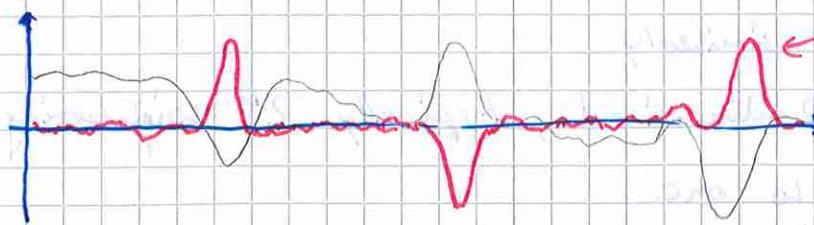
\Rightarrow We do not need to repeat the redundant network architecture, but instead use a loop to actually move the window around. The weights in a window act as a filter that selects patterns of interest.

Moving filter with weights

(+1) -2 (+1)

high response if local shape matches filter shape

ReLU: set negative response to zero

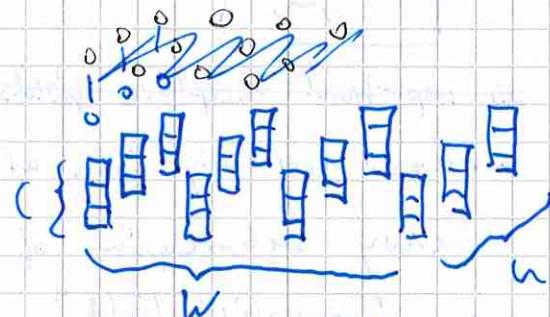


- Images contain many different patterns of interest \Rightarrow run many filters in parallel. The output of each filter forms a channel of the next neuron layer.

input layer (gray scale img)
 $w \times h$



1st hidden layer after C filters:
 $C \times w \times h$



- deep networks: stacking layers with small receptive fields yields a larger receptive field



\Rightarrow 3 neurons of first layer
 \Rightarrow 5 neurons of input

\Rightarrow complex patterns are formed by combining activation of several simple patterns in the previous layer

- downsampling ("pooling") layers: achieve stronger increase of receptive field by reducing number of neurons (\approx pixels). 10 example

if 4th layer has

$z_{4,1}$	$z_{4,2}$	$z_{4,3}$	$z_{4,4}$	$z_{4,5}$	$z_{4,6}$
0	0	0	0	0	0

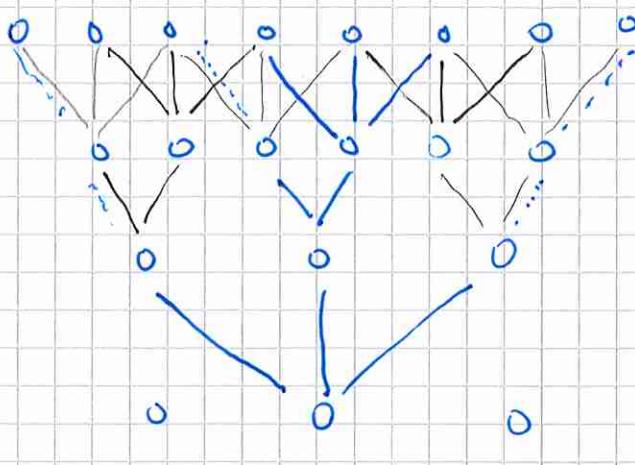
max pooling

$$z_{(k+1),1} = \max(z_{4,1}, z_{4,2}) \quad z_{(k+1),3} = \max(z_{4,5}, z_{4,6})$$

$$z_{(k+1),2} = \max(z_{4,3}, z_{4,4})$$

average pooling similarly

- if signal is 0-dimensional, typically 2^D neighboring elements are reduced to one.



neuron sees 8 neurons in the input layer

- strided convolution has a similar effect: don't move window to the next pixel, but skip one or more = "stride"
⇒ downsampling according to how many pixels are skipped
- convolutional layers are specified by window size and number of filters:
 - $5 \times 5 \times 8$: 8 filters of size 5×5
 - $3 \times 3 \times 32$: 32 - - - 3×3
 - $1 \times 1 \times 128$: 128 projections of the input channels (just scalar products of the center pixel)
- show famous networks on slides
 - classification networks (detect objects in image)
Le Net, Alex Net, GoogLe Net, VGG, Res Net
 - segmentation networks: U-net, fully convolutional networks
- show applications on slides