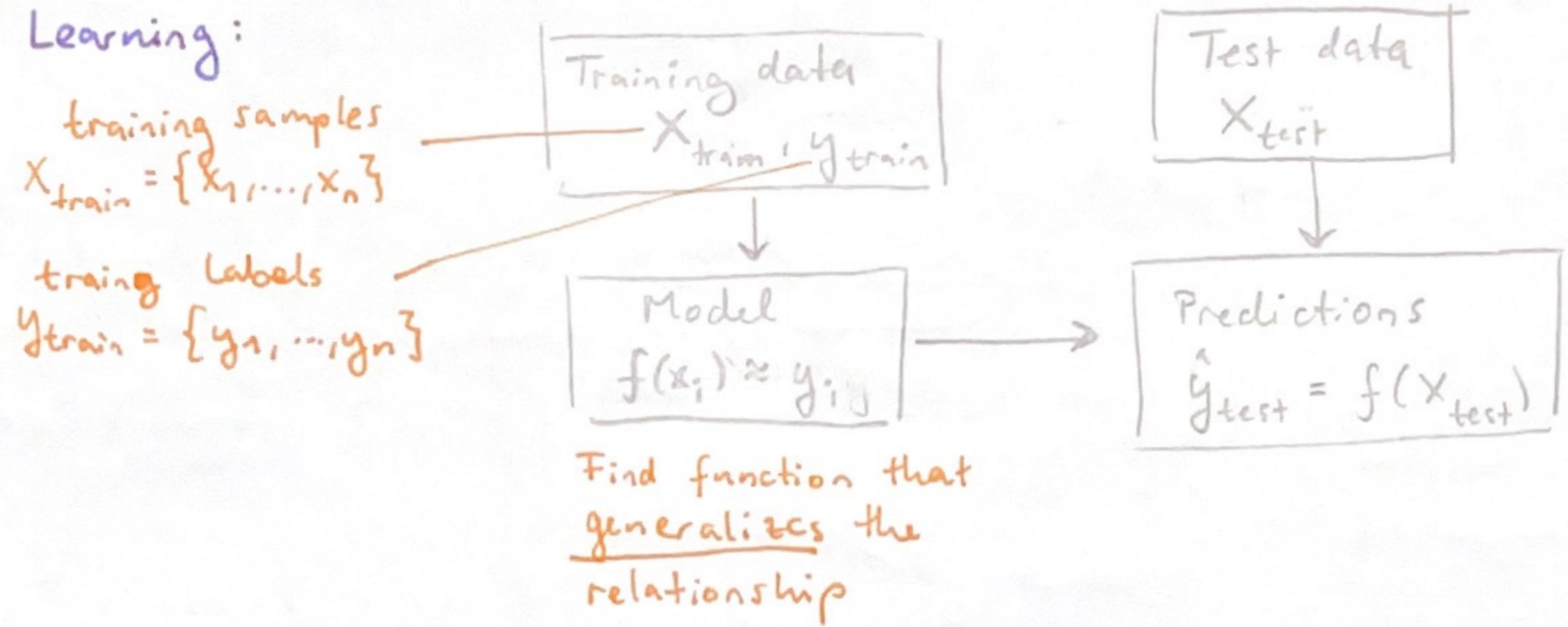


ML 01 : Intro

①

- Supervised Learning:



- Classification: If target y_i represents categories.

- Regression: If the targets y_i represent continuous numbers.

- Unsupervised Learning: Finding a structure in unlabelled data.

↳ Clustering
↳ Dimensionality Reduction

ML 02 : KNN

- 1-NN-Algorithm: Given a dataset $\mathcal{D} = \{(\vec{x}_i, y_i)\}_{i=1}^N$

→ to classify new observations

1. Define distance measure
2. Compute NN
3. Label with the NN's label

- k-NN-Algorithm: Look at multiple NNs and pick the majority label.

prob. of belonging
in class c

$$p(y=c | \vec{x}, k) = \frac{1}{k} \sum_{i \in N_k(x)} \mathbb{I}(y_i=c)$$

given k \vec{x} and
k nearest Neighbors
of x

Classification

$$\Rightarrow \hat{y} = \operatorname{argmax}_c p(y=c | \vec{x}, k)$$

Complexity per prediction:
 $\mathcal{O}(n)$

- Weighted K-NN:

$$p(y=c | \vec{x}, k) = \frac{1}{Z} \sum_{i \in N_k(x)} \frac{1}{d(x_i, x)} \mathbb{I}(y_i=c)$$

distance measure

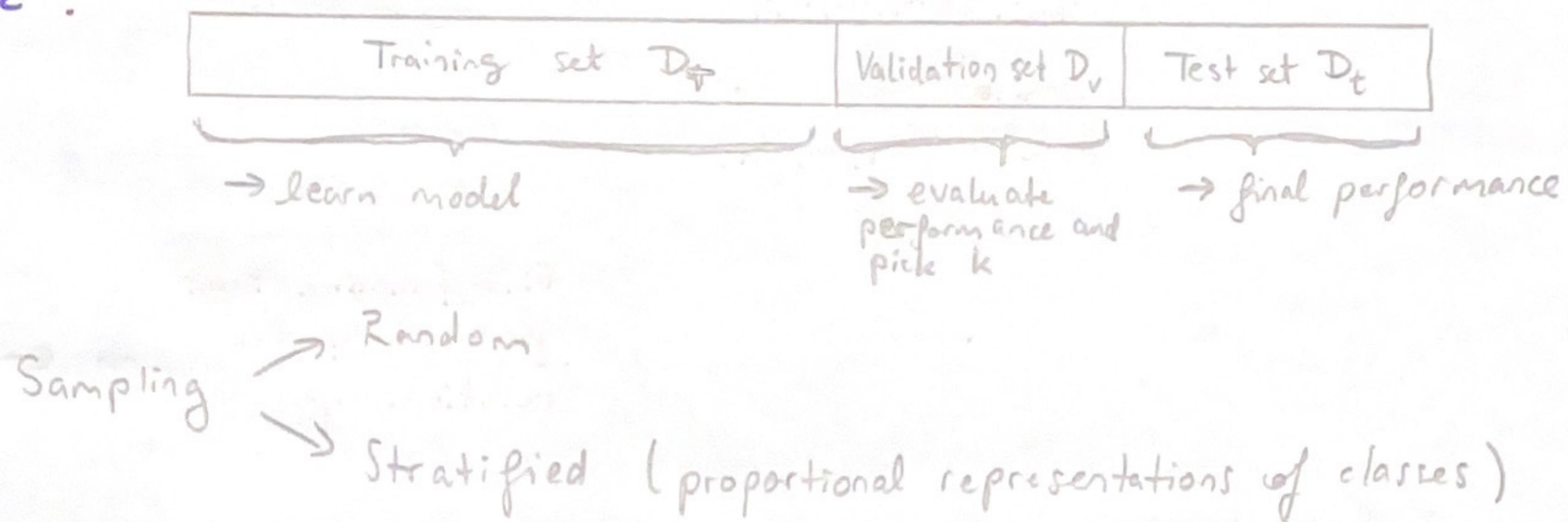
with $Z = \sum_{x \in N_k(x)} \frac{1}{d(x_i, x)}$

Normalization
weight → sum to 1

- k-NN-regression:

$$\hat{y} = \frac{1}{Z} \sum_{i \in N_k(x)} \frac{1}{d(x_i, x)} y_i \quad \text{"weighted mean"}$$

- How to choose k ?



- How to measure classification performance?

→ Confusion table

		True Label	
		$y=1$	$y=0$
Predicted	$y=1$	TP	FP
	$y=0$	FN	TN

True positives
&
True negatives

⇒ Accuracy: $acc = \frac{TP + TN}{TP + TN + FP + FN}$

precision: $prec = \frac{TP}{TP + FP}$

sensitivity / recall: $rec = \frac{TP}{TP + FN}$

...

F1 Score: $f_1 = \frac{2 \cdot prec \cdot rec}{prec + rec}$ (tries to max. prec & rec)

- Distance measures: * Euclidean L_2 : $\sqrt{\sum_i (u_i - v_i)^2}$

* L_1 norm: $\sum_i |u_i - v_i|$

* L_∞ norm: $\max_i |u_i - v_i|$

* Angle: $\cos \alpha = \frac{u^T v}{\|u\| \cdot \|v\|}$

(is equal to Euclidean distance on Normalized data) * Mahalanobis distance: $\sqrt{(u-v)^T \Sigma^{-1} (u-v)}$ Covariance matrix

- How to circumvent scaling issues?

→ Data standardization: $x_{i,\text{std}} = \frac{x_i - \mu_i}{\sigma_i}$ (Scale each feature to zero mean and unit variance)

- Curse of dimensionality: With an increase in dims the space becomes more empty.

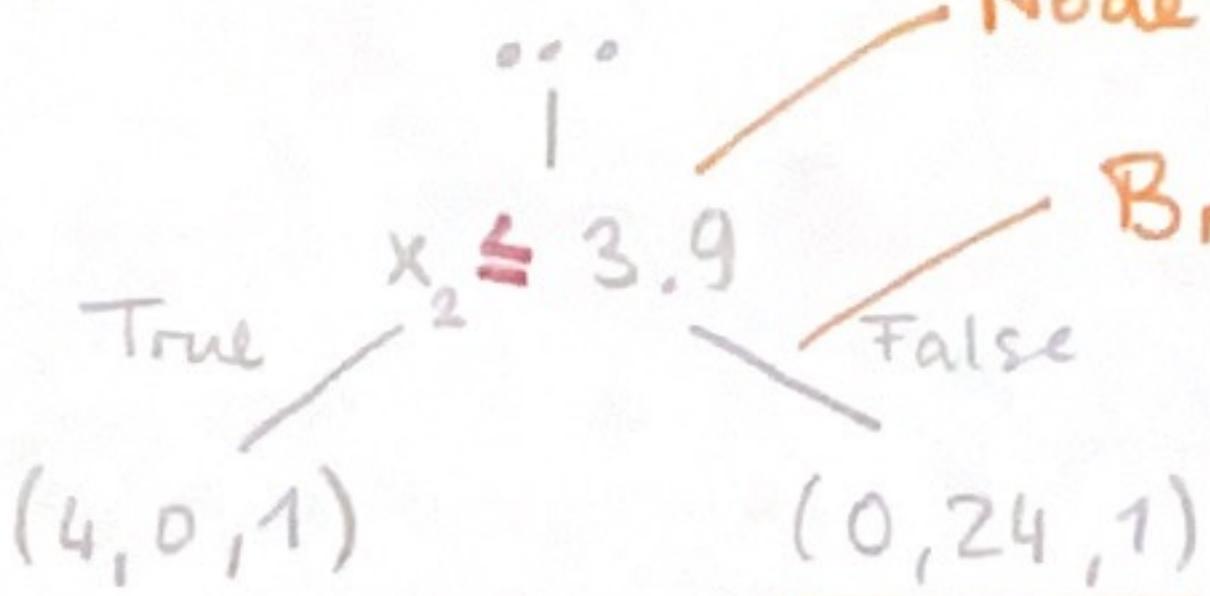
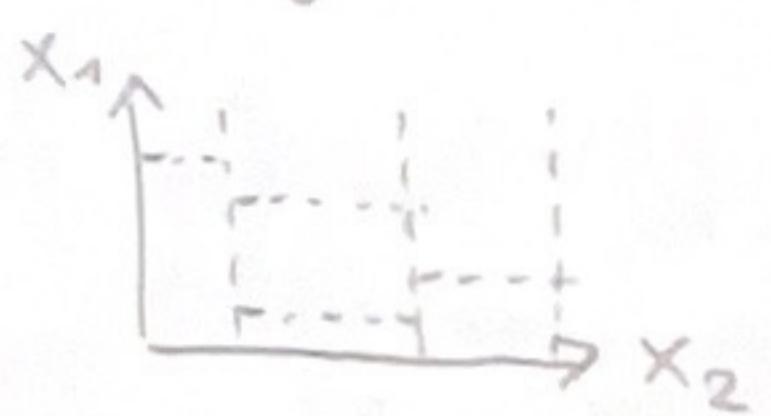
(in kNN the NN is far away → N would have to grow exponentially with the num of features)

ML 02: Decision Trees (DT)

(2)

- Decision tree:

- * partition the input space into cuboid regions.



Inference:
 Predicting a new sample.

Branch $\hat{=}$ different outcome of the feature test

Leaf $\hat{=}$ final nodes and the represented region in the input space.

- Probability for classifying:

$$p(y=c|R) = \frac{n_{c,R}}{\sum_{c_i \in C} n_{c_i,R}}$$

} points of all classes in this region

Probability that a point should be in a class given the region of the input space

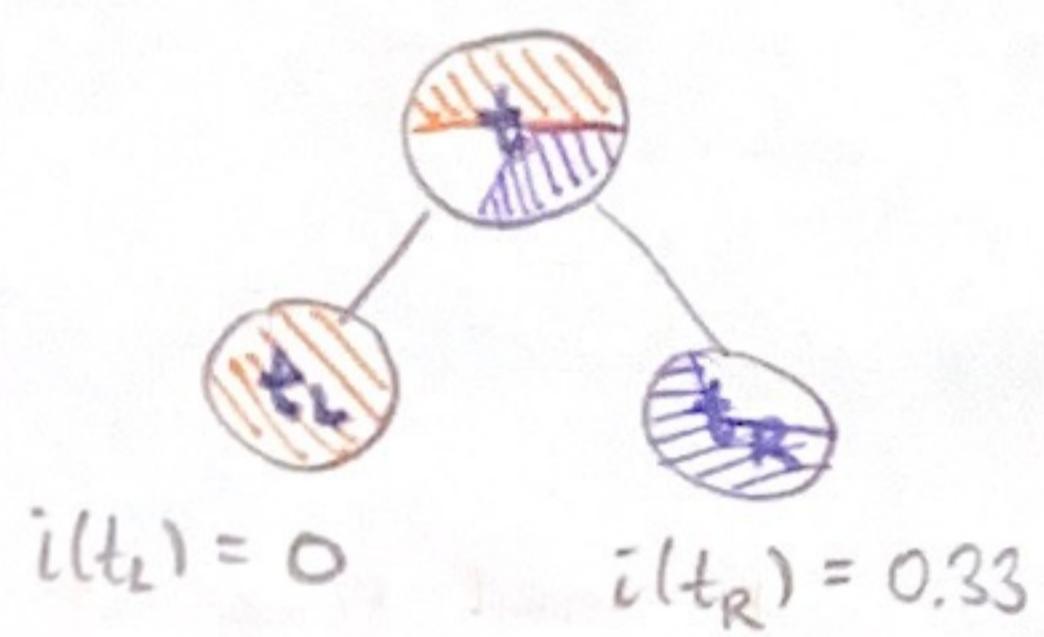
Inference

$$\Rightarrow \hat{y} = \operatorname{argmax}_c p(y=c|R) = \operatorname{argmax}_c n_{c,R}$$

Number of points in a cluster and region

- How are decision trees built?

$i(t) = 0.5$ → Greedy top-down approach (choosing the best split at each node)



\Rightarrow Split the node if it improves the misclassification rate

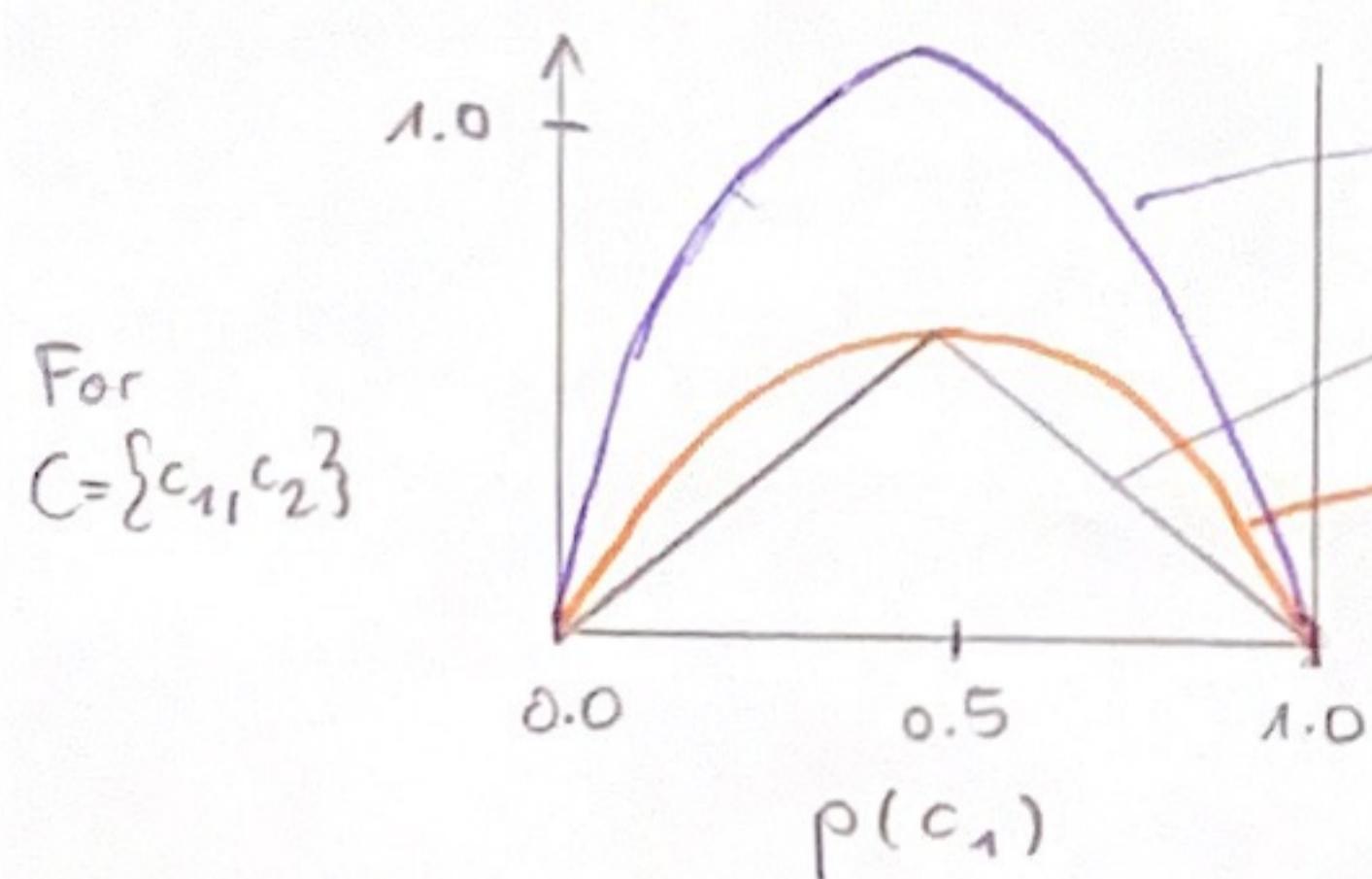
$$i_E(t) = 1 - \max_c p(y=c|t)$$

split node proportions

$$\Rightarrow \text{Improvement of split: } \Delta i(s,t) = i_E(t) - p_L \cdot i(t_L) - p_R \cdot i(t_R)$$

\Rightarrow Repeat the heuristic (split where $\Delta i(s,t)$ is largest)

- Impurity measures: ("criterion") (for categorical features)



Entropy: $i_H(t) = - \sum_{c_i \in C} \pi_{c_i} \log_2 \pi_{c_i}$

Gini Index: $i_G(t) = \sum_{c_i \in C} \pi_{c_i} (1 - \pi_{c_i}) = 1 - \sum_{c_i \in C} \pi_{c_i}^2$

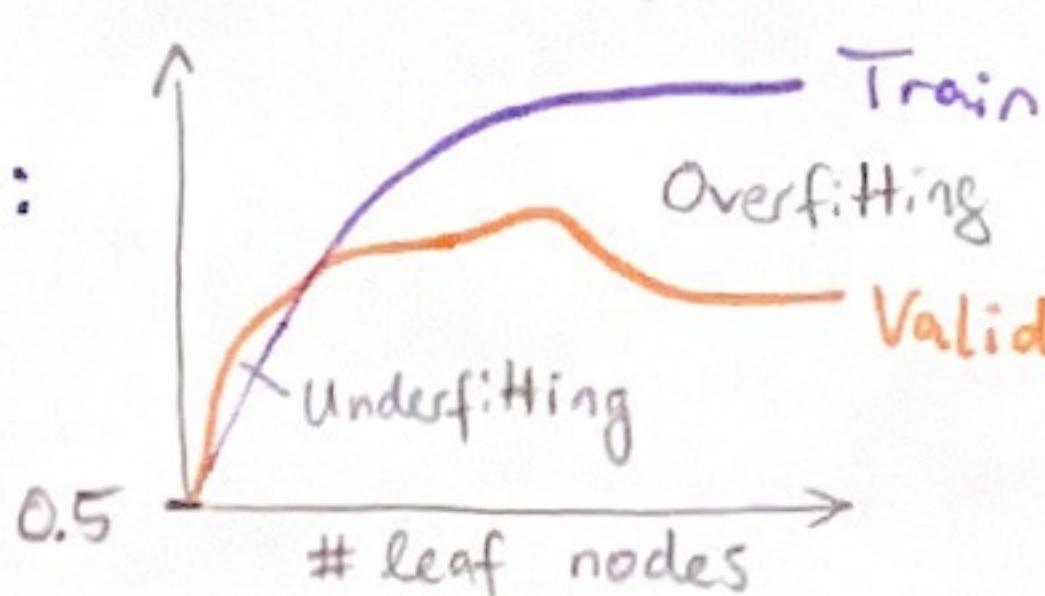
with $\pi_c = p(y=c|t)$

exponentially

\Rightarrow Advantage: The error $i(t)$ for unpure nodes is greater than for pure nodes.

(Gini Index is a bit faster than Entropy)

- Overfitting:



Underfitting Overfitting Validation \rightarrow tells how well the model generalizes

- K-fold Cross-Validation:

- * Split data into K folds (e.g., $10 \Rightarrow$ train 10 Models)
- ↪ Use $K-1$ for training and 1 for eval.
- ⇒ Average over all folds to estimate $\xrightarrow{\text{best hyperparams}}$
 $\xrightarrow{\text{model}}$

[Extreme Case: leave-one-out-cross-validation (LOOCV)
= N-fold CV if $N = \# \text{ of samples}$]

- When to stop growing a decision tree?

- pre-pruning
- Stopping criteria:
 - * distribution in branch is pure: $i(t) = 0 \ (\rightarrow \text{overfitting})$
 - * max. depth
 - * # of samples in node below threshold
 - * benefit of splitting below threshold $\Delta i(s, t) \leq t_A$
 - * accuracy on the validation set
 - Alternative: Grow fully + (Post-) pruning

- Reduced-Error Pruning:

- 1) Get the entire tree's error estimate on the valid. set: $\text{err}_{D_v}(T)$
- 2) For each node, calculate: $\text{err}_{D_v}(T/T_t)$ pruned tree at node t
- 3) Prune at node with highest error reduction
- 4) Repeat until $\forall t : \text{err}_{D_v}(t) < \text{err}_{D_v}(T/T_t)$

- Regression in DTs:

- * Compute mean at the leafs
- * MSE as splitting heuristic

- What's an ensemble? → Aggregate prediction from multiple models.

↪ Reduces Variance

e.g. Random Forests
↪ also use a sampled subset of features per tree

* Bagging (bootstrap aggregating):

- Train separate models on sampled sub-data-sets
- combine predictions

Rule of thumb:

	# of features
Classif.	$\log_2 f_d$
Regr.	$\log_2(d)$

* Boosting:

- Incrementally train weak classifiers that correct previous mistakes ("a chain")
- Give higher weight to misclassified examples

ML 03 : Probabilistic Inference

(3)

- Maximum Likelihood Estimation:

Example { Random experiment with 3 coin flips. With $p(\Theta) = \Theta$.
 The probability of a sequence of flips given by the parameters of some model:
 observed data, D } Θ the observations
 $p(HHTT | \Theta_1, \Theta_2, \Theta_3) =$ The likelihood of seeing H, T, T

\Rightarrow Find the Θ_i 's such that $p(HHTT | \Theta_1, \Theta_2, \Theta_3)$ is as large as possible.

Max. Likelihood: Find the parameters that maximise the likelihood of our observation.

Assumption 1: The flips are identical distribution.

$$\Rightarrow p_i(F_i = f_i | \Theta_i) = p(F_i = f_i | \Theta)$$

Assumption 2: The flips are independent.

$$\Rightarrow p(HHTT | \Theta) = p(H | \Theta) \cdot p(T | \Theta) \cdot p(T | \Theta) \\ = \prod_{i=1}^3 p(F_i = f_i | \Theta)$$

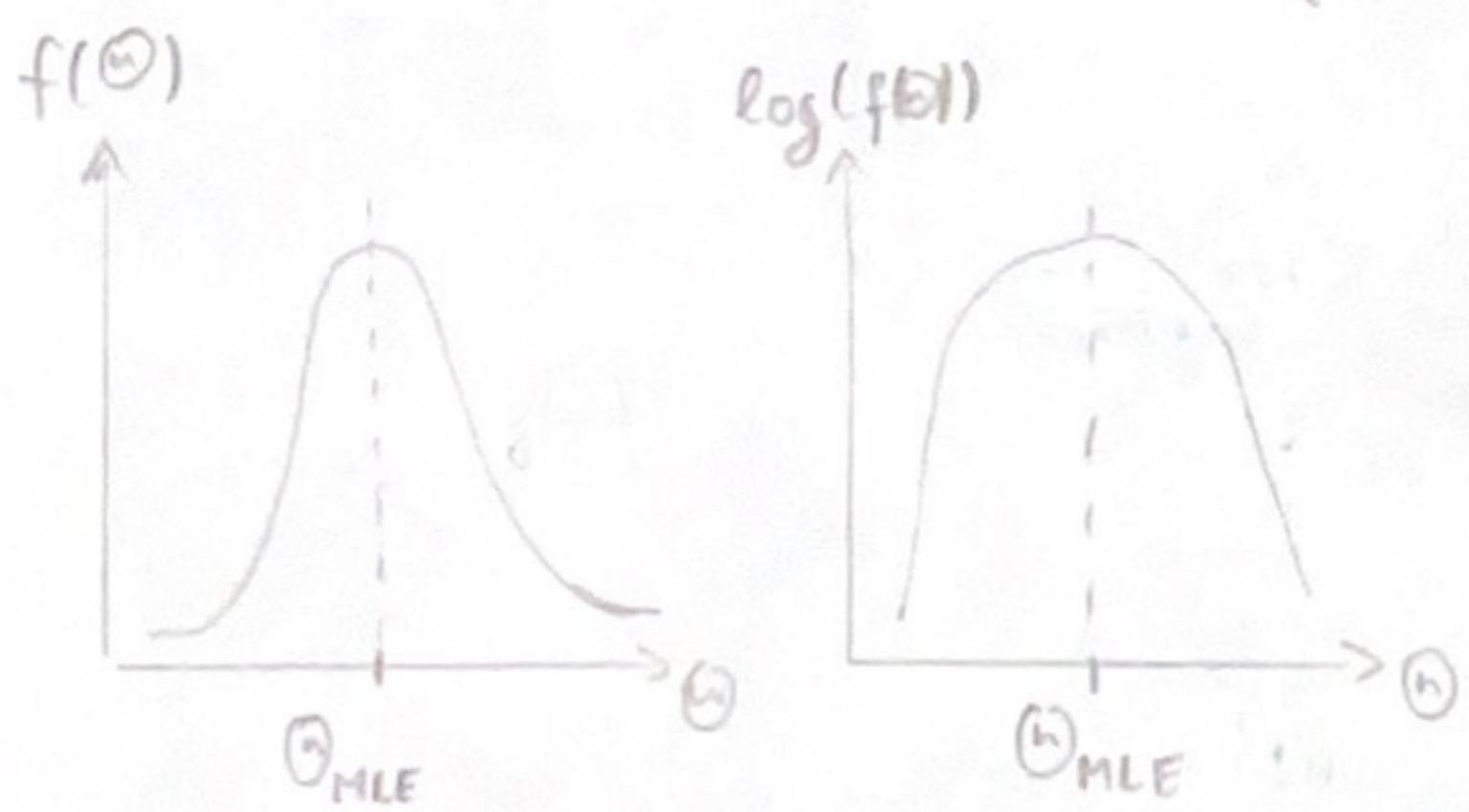
Total assumption:

The flips are independent and identically distributed (i.i.d.).

Note:
 This is not
 a probability
 distribution
 (does not integrate
 to 1)

\Rightarrow Interpret $p(D | \Theta)$ as a function $f(\Theta)$:

Find the maximum of this function:



Monotonic functions preserve critical points

$$\Theta_{MLE} = \underset{\Theta \in [0,1]}{\operatorname{argmax}} f(\Theta) = \underset{\Theta \in [0,1]}{\operatorname{argmax}} \log(f(\Theta)) \\ \dots \\ = \frac{|T|}{|H| + |T|} \text{ for the given Bernoulli coin example}$$

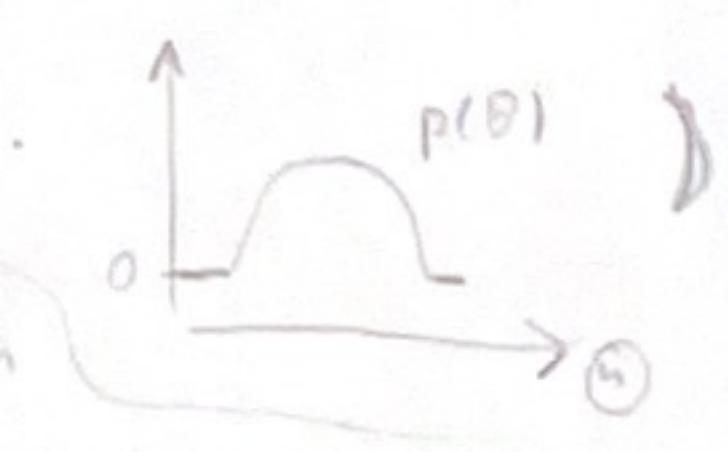
- Bayesian inference:

\rightarrow A prior distribution $p(\Theta)$ represents our beliefs about Θ before we observe any data.

(= basically a distribution/likelihood of values for Θ , e.g.

constraints

- 1) Does not depend on data
- 2) $p(\Theta) \geq 0 \quad \forall \Theta$
- 3) $\int p(\Theta) d\Theta = 1$



→ The posterior distribution encodes our beliefs in Θ after observing the data D . likelihood of D given Θ , prior, prior distribution

$$p(\Theta | D) = \frac{p(D | \Theta) \cdot p(\Theta)}{p(D)} \quad \text{Evidence distribution / marginal likelihood}$$

posterior \propto likelihood · prior

* The likelihood dominates the posterior with more data.

- Maximum a posteriori estimation (MAP): MLE ignores prior beliefs and performs poorly if little data is available.

$$\Rightarrow \Theta_{\text{MAP}} = \underset{\Theta}{\operatorname{argmax}} \ p(\Theta | D)$$

$$= \underset{\Theta}{\operatorname{argmax}} \frac{p(D | \Theta) p(\Theta)}{p(D)} = \underset{\Theta}{\operatorname{argmax}} p(D | \Theta) p(\Theta)$$

↑
the constant does not matter

The prior $p(\Theta)$ is chosen to be a Beta distribution (for simpler calculations)

$$\text{Beta}(\Theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \Theta^{a-1} (1-\Theta)^{b-1}, \Theta \in [0,1]$$

distribution parameters ↓
a, b > 0 Gamma function
 $\Gamma(n) = (n-1)!$ for $n \in \mathbb{N}$

$$\begin{aligned} \Rightarrow \Theta_{\text{MAP}} &= \underset{\Theta}{\operatorname{argmax}} \ p(D | \Theta) \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \Theta^{a-1} (1-\Theta)^{b-1} \\ &= \underset{\Theta}{\operatorname{argmax}} \underbrace{p(D | \Theta)}_{\substack{\text{by the Bernoulli-Experiment} \\ \# \text{ of tails}}} \cdot \Theta^{a-1} (1-\Theta)^{b-1} \\ &= \Theta^{|T|+a-1} (1-\Theta)^{|H|+b-1} \end{aligned}$$

$$= \underset{\Theta}{\operatorname{argmax}} \Theta^{|T|+a-1} (1-\Theta)^{|H|+b-1}$$

$$= \underset{\Theta}{\operatorname{argmax}} (|T|+a-1) \log \Theta + (|H|+b-1) \log (1-\Theta)$$

$$\downarrow = \frac{|T|+a-1}{|H|+|T|+a+b-2} \quad \begin{array}{l} \text{(solution for the given} \\ \text{Bernoulli coin example)} \end{array}$$

The data dominates this
for a lot of data

- How to estimate the posterior distribution:

→ One wants to know the posterior distribution $p(\Theta | D)$:

↳ How certain are we in our estimate $\hat{\Theta}_{MAP}$?

→ What is the probability that Θ is in some interval?

$$\left[p(\Theta | D) \propto \Theta^{|T|+a-1} (1-\Theta)^{|H|+b-1} \right]$$

⇒ We need the normalization constant, such that

$$\int p(\Theta | D) d\Theta = 1$$

Method 1: Compute Integral and divide by it

Method 2: Compare the unnormalized posterior to a known PDF (e.g., Beta distribution)

$$\Rightarrow \underline{p(\Theta | D) = \text{Beta}(\Theta | a+|T|, b+|H|)}$$

Solution for
the given
coin example

* With more data the posterior becomes more peaked
→ more certain of estimate for Θ

- Posterior predictive distribution:

With the outcome of the next flip $f \in \{0, 1\}$

$$P(F=f | D, a, b) = p(f | D, a, b)$$

$$\stackrel{\text{sum rule of probability}}{=} \int_0^1 p(f, \Theta | D, a, b) d\Theta$$

$$\stackrel{\text{conditional independence}}{=} \int_0^1 p(f | \Theta, D, a, b) p(\Theta | D, a, b) d\Theta$$

$$\stackrel{?}{=} \int_0^1 p(f | \Theta) p(\Theta | D, a, b) d\Theta$$

$$\stackrel{...}{=} \frac{(|T|+a)^f (|H|+b)^{1-f}}{|T|+a+|H|+b}$$

$$p(f|\Theta) = \text{Ber}(f|\Theta)$$

$$= \Theta^f (1-\Theta)^{1-f}$$

$$\stackrel{?}{=} \text{Ber}\left(f \mid \frac{|T|+a}{|T|+a+|H|+b}\right)$$

The approach that
computes the full Bayesian
posterior distribution

„Fully Bayesian“

⇒ With lots of data, the differences in predictions between MLE, MAP and fully Bayesian become less noticeable.

ML 04: Linear Regression

- What's the linear regression problem?

Given: \rightarrow observations $X = \{x_1, \dots, x_N\}$
targets $y = \{y_1, \dots, y_N\}$

Find: Mapping from Inputs to targets: $y_i \approx f(x_i)$
 \rightarrow which generalizes to new x

- What's the lin. model?

bias weights

$$f_{\vec{w}}(\vec{x}_i) = w_0 + \vec{w}^T \vec{x}_i$$

$$= (w_0, w_1, \dots, w_D) \begin{pmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_D \end{pmatrix} = \vec{w}^T \vec{x}$$

Notation:

Assume the bias term is always absorbed and write \vec{w} and \vec{x} .

- How is the ~~so far~~ optimal lin. model found?

1) Loss function, e.g. least squares: $E_{LS}(\vec{w}) = \frac{1}{2} \sum_{i=1}^N (\vec{w}^T \vec{x}_i - \vec{y}_i)^2$

\hookrightarrow Measures the misfit between the model and observations

2) Objective: $w^* = \underset{w}{\operatorname{argmin}} E_{LS}(w)$

\hookrightarrow multiple observations \vec{x}_i^T stacked into a matrix

$$= \underset{w}{\operatorname{argmin}} \frac{1}{2} (\vec{x} \vec{w} - \vec{y})^T (\vec{x} \vec{w} - \vec{y})$$

3) Optimal solution: $\nabla_w E_{LS}(\vec{w}) = \frac{\text{slide}}{\text{...}} = \vec{x}^T \vec{x} \vec{w} - \vec{x}^T \vec{y} = 0$

$$\Rightarrow w^* = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{y}$$

- What if the dependency between y and \vec{x} is not linear?

\Rightarrow Choose a universal function approximator.

f is still lin. in \vec{w} but not lin. in x

$$f_{\vec{w}}(\vec{x}) = w_0 + \sum_{j=1}^M w_j \phi_j(\vec{x})$$

Absorb bias $\phi_0 = 1$

$$= \vec{w}^T \vec{\phi}(\vec{x})$$

with $\phi_j(x)$ as some non-linear function,
e.g. $\phi_1(x) = \sin(x)$
"transform the basis"

\rightarrow Typical basis functions:

* Polynomials $\phi_j(x) = x^j$

* Gaussian $\phi_j(x) = e^{-\frac{(x-\mu_j)^2}{2s^2}}$

* Logistic Sigmoid $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$

with $\sigma(a) = \frac{1}{1+e^{-a}}$

Design matrix:

$$\Phi = \begin{pmatrix} \vec{\phi}_0(x_1) & \phi_1(x_1) & \dots & \phi_M(x_1) \\ \vec{\phi}_0(x_2) & \phi_1(x_2) & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \vec{\phi}_0(x_N) & \phi_1(x_N) & \dots & \phi_M(x_N) \end{pmatrix} \in \mathbb{R}^{N \times (M+1)}$$

$$= D \quad 2) \text{ Objective: } E_{LS}(\vec{\omega}) = \frac{1}{2} (\Phi \vec{\omega} - \vec{y})^T (\Phi \vec{\omega} - \vec{y})$$

$$3) \text{ Solution: } \omega^* = (\Phi^\top \Phi)^{-1} \Phi^\top \tilde{y}$$

Design matrix

How is the degree of the polynomial chosen?

M → too small \Rightarrow underfitting } \Rightarrow Use training & validation
→ too large \Rightarrow overfitting } sets to optimize performance

Regularization with

→ Controlling overfitting with regularization:
L2-regularization / Ridge Regression regularization strength

Ridge Regression

Ridge Regression

$$E_{\text{ridge}}(\vec{w}) = \frac{1}{2} \sum_{i=1}^N (\vec{w}^T \vec{\phi}(\vec{x}_i) - \vec{y}_i)^2 + \frac{\lambda}{2} \|\vec{w}\|_2^2$$

squared
 L^2 norm

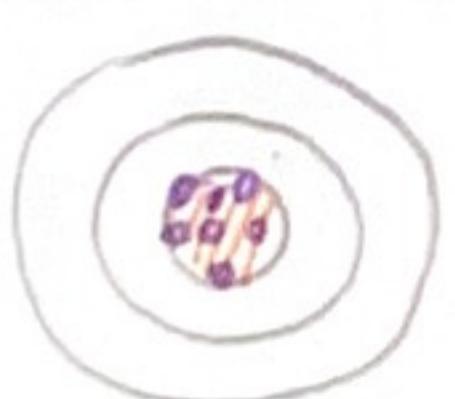
→ Larger regularization strength λ
leads to smaller weights \vec{w}

-What's the tradeoff between bias and variance?

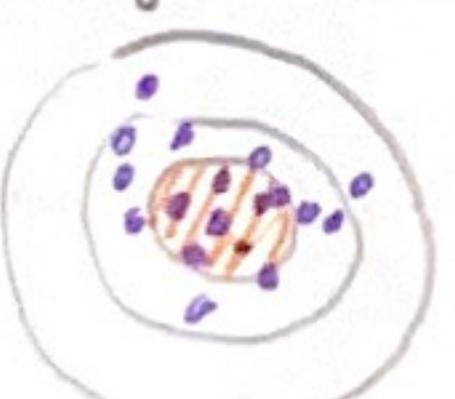
total error } **Bias**: expected error due to model mismatch

Variance: variation due to high randomness in the data

Low variance



High Variance



("Sensitivity of model towards noise in the data")

High Bias Low Bias



Each dot is the solution of a trained model (on a different subset of the training data)

- model is too flexible and thus captures noise in the data
 - ↳ this is overfitting
 - ↳ If the model has a high capacity or λ is too low

→ Trade-off: 1) Select a model with high capacity
2) keep the variance low by choosing a good λ

ML 04: Probabilistic linear regression

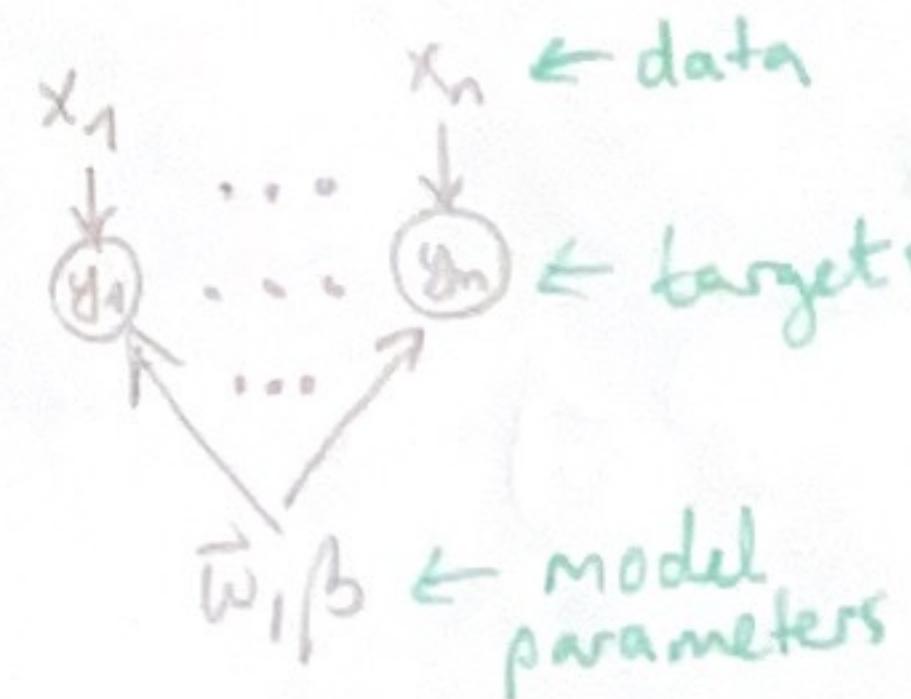
* Given data is sampled from a function plus noise

$$\Rightarrow y_i = f_{\vec{w}}(\vec{x}_i) + \underbrace{\epsilon_i}_{\text{noise}, \epsilon_i \sim N(0, \beta^{-1})} \quad \text{with } \beta = \frac{1}{\sigma^2}$$

$$\Rightarrow y_i \sim N(f_{\vec{w}}(\vec{x}_i), \beta^{-1}) \quad \text{"fixed precision"}$$

- With MLE:

Likelihood



Likelihood of a single sample: $p(y_i | f_{\vec{w}}(\vec{x}_i), \beta) = N(y_i | f_{\vec{w}}(\vec{x}_i), \beta^{-1})$

=> Likelihood of entire dataset: $p(\vec{y} | X, \vec{w}, \beta) = \prod_{i=1}^N p(y_i | f_{\vec{w}}(\vec{x}_i), \beta)$

independent

Maximum Likelihood

$$\Rightarrow \vec{w}_{ML}, \beta_{ML} = \underset{\vec{w}, \beta}{\operatorname{argmax}} \quad p(\vec{y} | X, \vec{w}, \beta)$$

\ln is monotonic

$$= \underset{\vec{w}, \beta}{\operatorname{argmin}} -\ln(p(\vec{y} | X, \vec{w}, \beta))$$

= $E_{ML}(\vec{w}, \beta)$ maximum likelihood error function

$$\Rightarrow E_{ML}(\vec{w}, \beta) = -\ln \left[\prod_{i=1}^N \underbrace{\frac{\beta^{-1}}{\sqrt{2\pi}} \exp \left(-\frac{\beta}{2} (\vec{w}^T \vec{\phi}(\vec{x}_i) - y_i)^2 \right)}_{\text{Gaussian}} \right]$$

expand + log rules

likelihood of all data

$$= \frac{\beta}{2} \sum_{i=1}^N [(\vec{w}^T \vec{\phi}(\vec{x}_i) - y_i)^2] - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi$$

Optimize the likelihood w.r.t. \vec{w}

with respect to

$$\Rightarrow \vec{w}_{ML} = \underset{\vec{w}}{\operatorname{argmin}} E_{ML}(\vec{w}, \beta)$$

$$= \dots = \underset{\vec{w}}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^N (\vec{w}^T \vec{\phi}(\vec{x}_i) - y_i)^2 = \underset{\vec{w}}{\operatorname{argmin}} E_{LS}(\vec{w})$$

This is the least squares

error fn!

=> Maximizing the likelihood is equivalent to minimizing the least squares error function.

Optimize the likelihood w.r.t. β

$$\beta_{ML} = \underset{\beta}{\operatorname{argmin}} E_{ML}(\vec{w}_{ML}, \beta)$$

already minimizes for \vec{w}

$$\Rightarrow \frac{\partial}{\partial \beta} E_{ML}(\vec{w}_{ML}, \beta) = \dots \stackrel{!}{=} 0$$

Solve for β

$$\Rightarrow \frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{i=1}^N (\vec{w}_{ML}^T \vec{\phi}(\vec{x}_i) - y_i)^2$$

This is the mean residual/error

=> Advantage compared to "regular" lin. regression:

Besides the prediction we get a confidence value for it.

- With MAP:

Connection to the coin flip example:

	train data	likelihood	prior	posterior
Coin	$D = X$	$p(D \theta)$	$p(\theta a, b)$	$p(\theta D)$
Regr.	$D = \{X, \vec{y}\}$	$p(\vec{y} X, \vec{w}, \beta)$	$p(\vec{w} \cdot)$	$p(\vec{w} X, \vec{y}, \beta, \cdot)$

RECAP $\Rightarrow D$

$$p(\vec{w}|X, \vec{y}, \beta) = \frac{\underbrace{p(\vec{y}|X, \vec{w}, \beta)}_{\text{likelihood}} \cdot \underbrace{p(\vec{w}|\cdot)}_{\text{prior}}}{\underbrace{p(\vec{y}|X, \beta, \cdot)}_{\text{Normalization constant}}} \propto \text{likelihood} \cdot \text{prior}$$

Pick prior
 $p(\vec{w}|\cdot)$

$$\left\{ p(\vec{w}|\alpha) = N(\vec{w}|0, \alpha^{-1}I) = \left(\frac{\alpha}{2\pi}\right)^{\frac{M}{2}} \exp\left(-\frac{\alpha}{2}\vec{w}^T\vec{w}\right) \right.$$

of elements in \vec{w}
precision

Zero mean
→ small values of w
→ reduces overfitting

optimize the posterior w.r.t. \vec{w}

(This ex. assumes that β is given)

$$\left[\begin{array}{l} \vec{w}_{MAP} = \underset{\vec{w}}{\operatorname{argmax}} \ p(\vec{w}|X, \vec{y}, \beta) \\ \text{plug in + log rules + remove constant} \Rightarrow \dots = \underset{\vec{w}}{\operatorname{argmin}} - \underbrace{\ln[p(\vec{y}|X, \vec{w}, \beta)] - \ln[p(\vec{w}|\alpha)]}_{E_{MAP}(\vec{w})} \\ E_{MAP} = \dots \propto E_{\text{ridge}}(\vec{w}) + \text{const.} \\ \text{If the prior is Gaussian} \end{array} \right]$$

- With fully Bayesian approach:

→ Computes the full posterior distribution

Conjugate prior

Since the prior and likelihood are Gaussians, the posterior $\Rightarrow p(\vec{w}|D) = N(\vec{w}|\vec{\mu}, \Sigma)$ is also a Gaussian! (in the given example)

$$p(\vec{w}|D) \propto p(\vec{y}|X, \vec{w}, \beta) \cdot p(\vec{w}|\alpha)$$

$$p(\vec{w}|D) = N(\vec{w}|\vec{\mu}, \Sigma) \quad \text{with} \quad \vec{\mu} = \beta \Sigma \Phi^T \vec{y}$$

$$\Sigma^{-1} = \alpha I + \beta \Phi \Phi^T$$

- How to predict new data?

With MLE:

$$p(\hat{y}_{\text{new}} | \vec{x}_{\text{new}}, \vec{w}_{ML}, \beta_{ML}) = N(\hat{y}_{\text{new}} | \vec{w}_{ML}^T \vec{\Phi}(\vec{x}_{\text{new}}), \beta_{ML}^{-1})$$

With MAP:

$$p(\hat{y}_{\text{new}} | \vec{x}_{\text{new}}, \vec{w}_{MAP}, \beta) = N(\hat{y}_{\text{new}} | \vec{w}_{MAP}^T \vec{\Phi}(\vec{x}_{\text{new}}), \beta^{-1})$$

predictive distributions

$$\begin{aligned} p(\hat{y}_{\text{new}} | \vec{x}_{\text{new}}, D) &= \int p(\hat{y}_{\text{new}}, \vec{w} | \vec{x}_{\text{new}}, D) d\vec{w} \\ &= \dots \end{aligned}$$

Use the full posterior distribution to comp. the posterior predictive distribution:

also shows variance

ML 05: Linear Classification

- What's the classification problem?

Given: Observations $X = \{x_1, x_2, \dots, x_N\}, x_i \in \mathbb{R}^D$

Possible classes $C = \{1, \dots, Q\}$

Labels $y = \{y_1, y_2, \dots, y_N\}, y_i \in C$

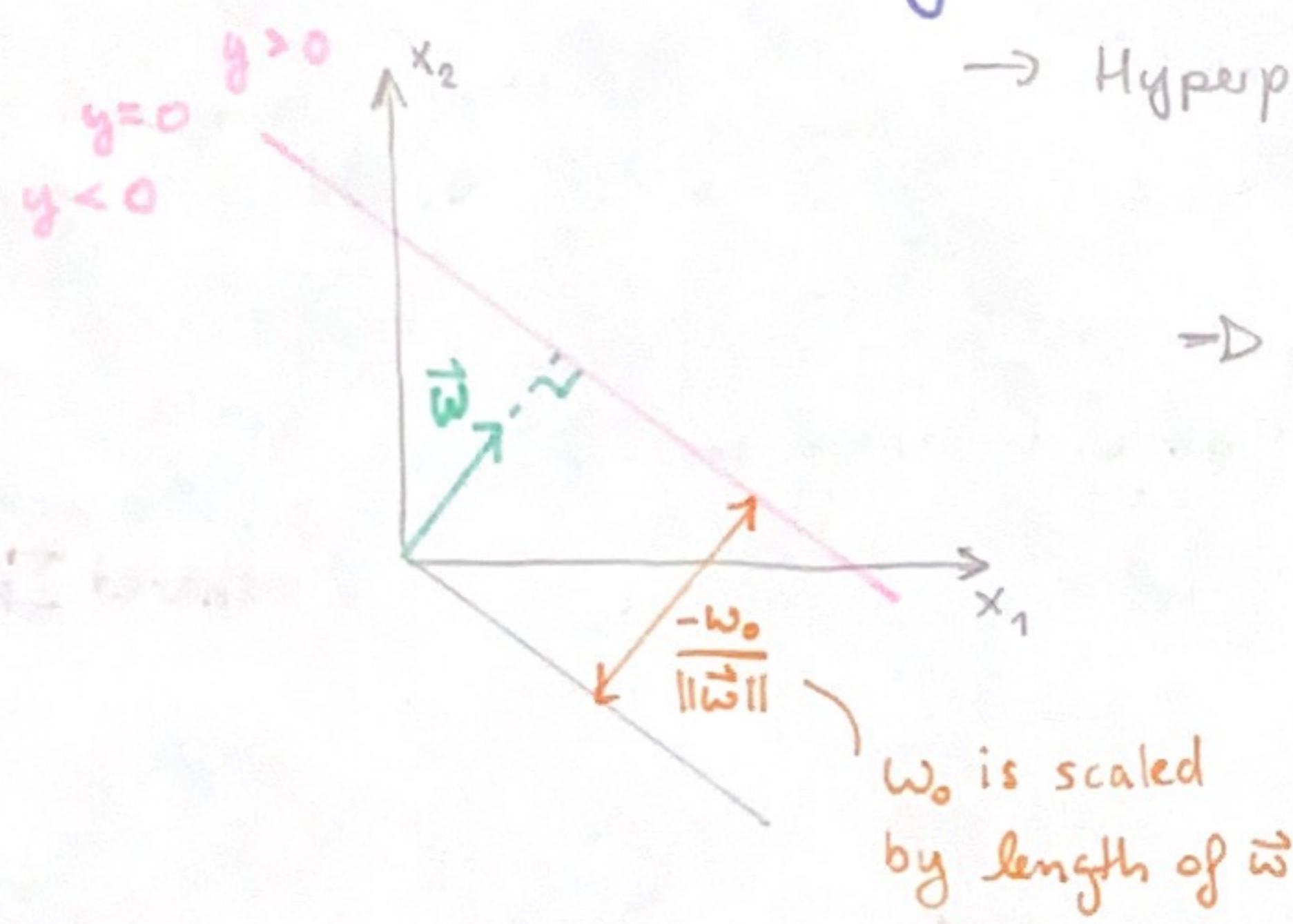
Find: function $f: \mathbb{R}^D \rightarrow C$ that maps x_i to y_i

- Zero-One-Loss: \rightarrow denotes the number of misclassified examples

$$l_{01}(\hat{y}_i, \tilde{y}_i) = \sum_{i=1}^N I(\hat{y}_i \neq y_i)$$

Indicator = $\begin{cases} 1, & \text{if condition is true} \\ 0, & \text{else} \end{cases}$

- Hyperplane as a decision boundary:



\rightarrow Hyperplane is defined by normal \vec{w} and offset w_0

$$\Rightarrow \vec{w}^T \vec{x} + w_0 \begin{cases} = 0, & \text{if } \vec{x} \text{ on plane} \\ > 0, & \text{if } \vec{x} \text{ on normal's side} \\ < 0, & \text{else} \end{cases}$$

If a hyperplane can perfectly separate all points into two classes, the data set is lin. separable

- Perceptron for binary classification:

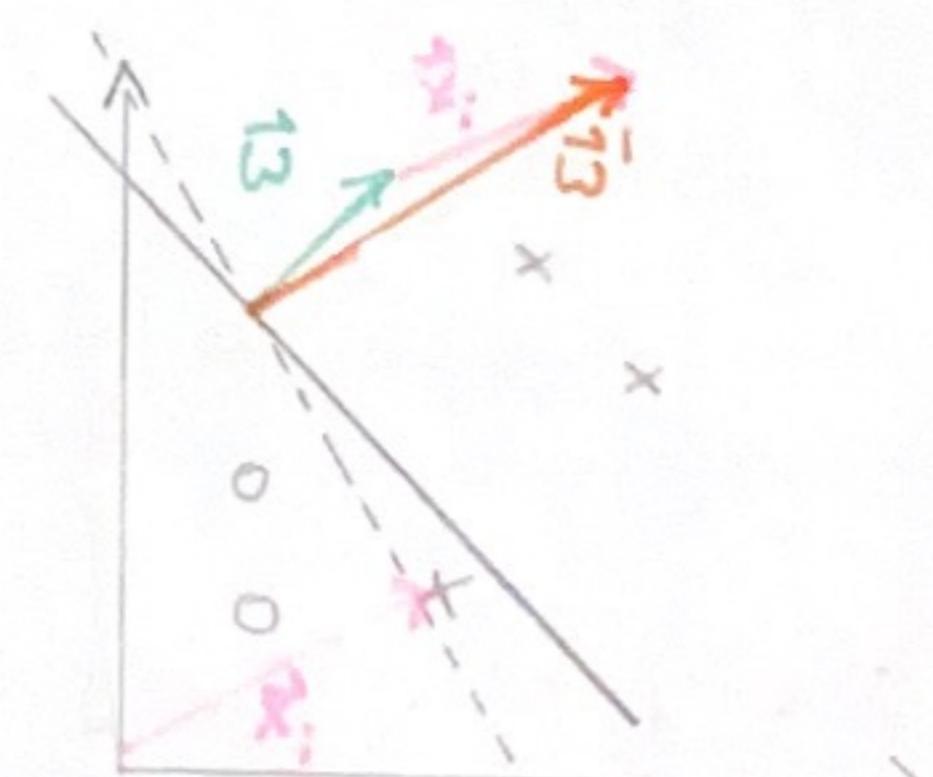
$$\hat{y} = f(\vec{w}^T \vec{x} + w_0) \text{ with } f(t) = \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{else} \end{cases}$$

Training: Init. $\vec{w}, w_0 \leftarrow \vec{0}$

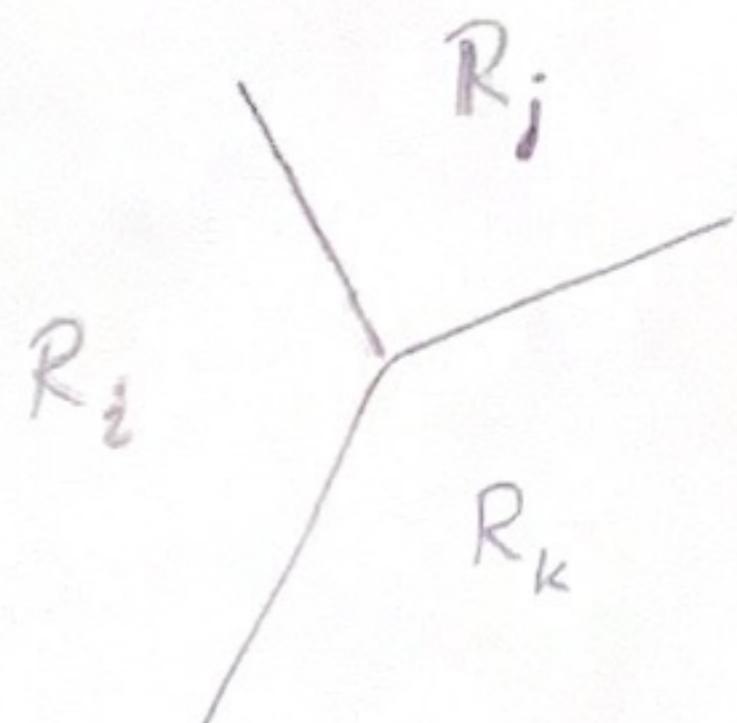
iterate until convergence

For each misclassified sample update

$$\vec{w} \leftarrow \begin{cases} \vec{w} + \vec{x}; & \text{if } y_i = 1 \\ \vec{w} - \vec{x}; & \text{if } y_i = 0 \end{cases}$$

$$w_0 \leftarrow \begin{cases} w_0 + 1 & \text{if } y_i = 1 \\ w_0 - 1 & \text{if } y_i = 0 \end{cases}$$


- Multi-class discriminant:



$$\hat{y} = \arg \max_c f_c(\vec{x}) \quad \text{with} \quad f_c(\vec{x}) = \vec{w}_c^T \vec{x} + w_{oc}$$

Assign \vec{x} to the class which produces the largest $f_c(\vec{x})$

Each class has a normal \vec{w}_c

\rightarrow need basis func. ϕ that maps data to a space where it is lin. separable
 \Rightarrow doesn't work if the data is not lin. Separable

- Probabilistic models for classification:

- 1) Generative models $p(y=c|\vec{x}) = \frac{p(\vec{x}|y=c) \cdot p(y=c)}{p(\vec{x})}$
- 2) Discriminative models

- Probabilistic generative models for lin. classification:

Idea: Obtain posterior: $p(y=c|\vec{x}) \propto \underbrace{p(\vec{x}|y=c)}_{\text{class conditional}} \cdot \underbrace{p(y=c)}_{\text{class prior}}$

\rightarrow distribution of samples in class c \rightarrow probability of a point belonging to class c

Choose prior: $y \sim \text{Categorical}(\vec{\theta})$

1.) Choose parametric models
 $p(y=c|\vec{\theta})$
 and
 $p(\vec{x}|y=c, \vec{\psi})$

$$\Rightarrow p(y) = \prod_{c=1}^C \theta_c^{\mathbb{I}(y=c)} \quad \text{so that } p(y=c) = \theta_c$$

$$\Rightarrow \hat{\theta}_c^{\text{MLE}} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}(y_i=c) \quad \text{"# of samples in class divided by # of all data"}$$

Choose class conditionals: \vec{x} is continuous

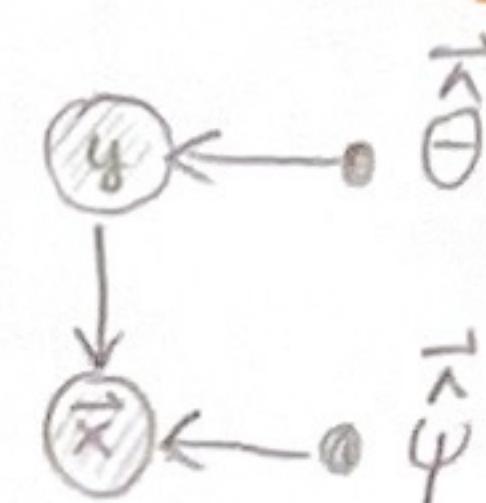
Multivariate normal for each class

$$\Rightarrow p(\vec{x}|y=c) = \mathcal{N}(\vec{x}|\vec{\mu}_c, \Sigma) \quad \begin{matrix} \checkmark & \text{separate means but shared } \Sigma \end{matrix}$$

$$= \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left\{-\frac{1}{2} (\vec{x} - \vec{\mu}_c)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_c)\right\}$$

2.) Estimate parameters $\{\vec{\psi}, \vec{\theta}\}$ from data (e.g. with MLE)

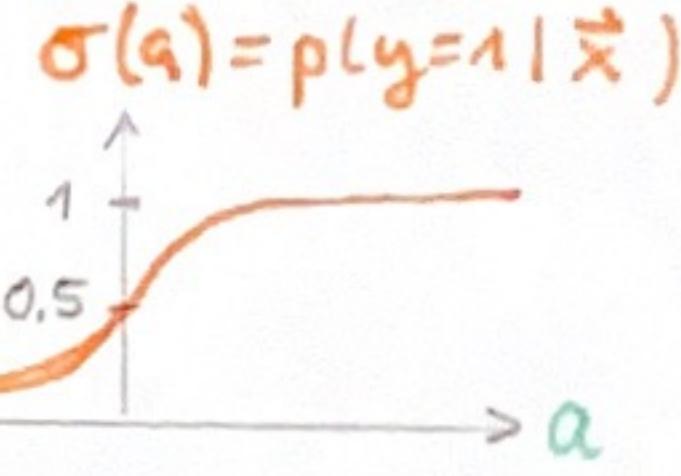
(3.) Possibility to generate new data:



1. Sample $y_{\text{new}} \sim p(y|\vec{\theta})$

2. Sample $\vec{x}_{\text{new}} \sim p(\vec{x}|y=y_{\text{new}}, \vec{\psi})$

$$\Rightarrow y|\vec{x} \sim \text{Bernoulli}(\sigma(\vec{w}^T \vec{x} + w_0))$$



Calculate $p(y=c|\vec{x})$:

Assume two classes $C = \{0, 1\}$

$$p(y=1|\vec{x}) = \frac{p(\vec{x}|y=1)p(y=1)}{p(\vec{x}|y=0)p(y=0)}$$

$$a = \log \frac{p(\vec{x}|y=1)p(y=1)}{p(\vec{x}|y=0)p(y=0)} =$$

$$\frac{p(\vec{x}|y=1)p(y=1)}{p(\vec{x}|y=0)p(y=0)} = \frac{1}{1 + \exp(-a)} := \sigma(a)$$

4.) Inference with LDA

Linear Discriminant Analysis (LDA) with $a = -\frac{1}{2}(\vec{x} - \vec{\mu}_1)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_1) + \log p(y=1) + \frac{1}{2}(\vec{x} - \vec{\mu}_0)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_0)$

$$= \vec{w}^T \vec{x} + w_0 \stackrel{!}{=} 0 \Rightarrow \text{Decision Boundary}$$

$$= \vec{w}^T \vec{x} + w_0 = -\frac{1}{2} \vec{\mu}_1^T \Sigma^{-1} \vec{\mu}_1 + \frac{1}{2} \vec{\mu}_0^T \Sigma^{-1} \vec{\mu}_0 + \log \frac{p(y=1)}{p(y=0)}$$

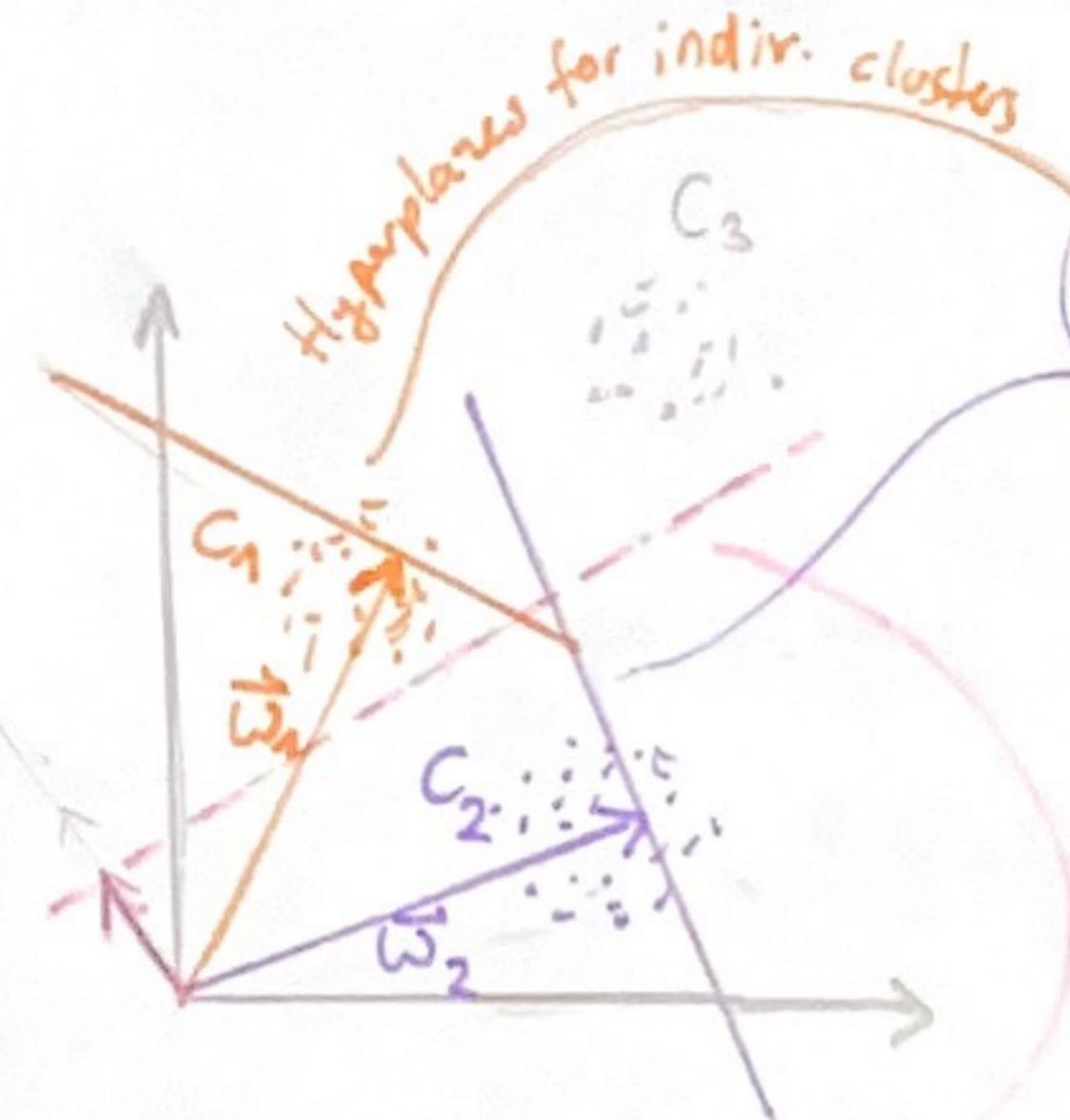
LDA for $C=2$: Posterior dist., is sigmoid of lin. func. of \vec{x} (8)

LDA for $C > 2$: $p(y=c|\vec{x}) = \frac{p(\vec{x}|y=c) p(y=c)}{\sum_{c'=1}^C p(\vec{x}|y=c') p(y=c')}$

$$\omega_{co} = -\frac{1}{2} \vec{\mu}_c^T \Sigma^{-1} \vec{\mu}_c + \log p(y=c)$$

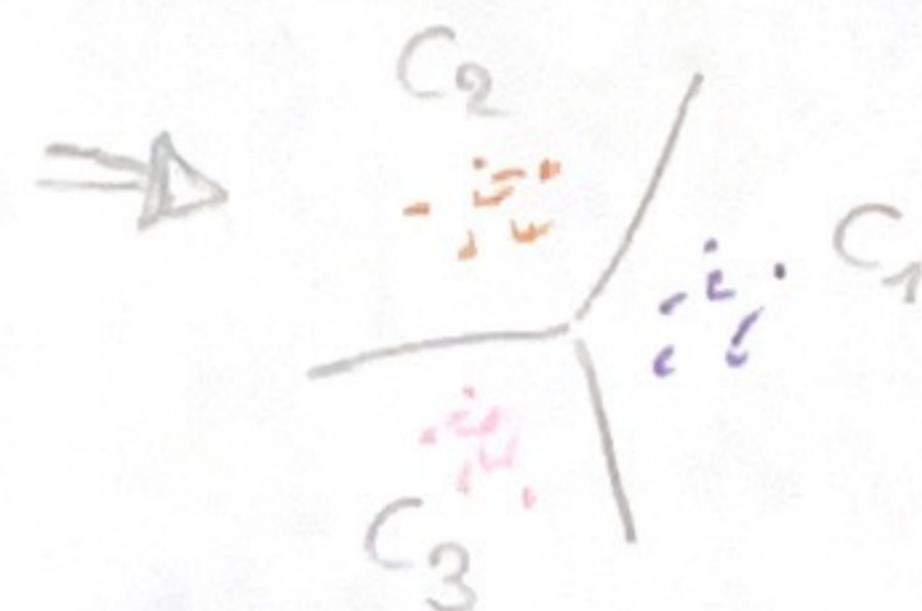
$$= \frac{\exp(\vec{\omega}_c^T \vec{x} + \omega_{c0})}{\sum_{c'=1}^C \exp(\vec{\omega}_{c'}^T \vec{x} + \omega_{c'0})}$$

each class has its own function



Decision boundary for 1 & 2
 $\vec{\omega}_1^T \vec{x} + \omega_{10} = \vec{\omega}_2^T \vec{x} + \omega_{20}$
 $\Rightarrow (\vec{\omega}_1 - \vec{\omega}_2)^T \vec{x} + \omega_{10} - \omega_{20} = 0$

new Hyperplane



IS Softmax

$$\sigma(\vec{x})_i = \frac{\exp(x_i)}{\sum_{k=1}^K \exp(x_k)}$$

$$\sigma: \mathbb{R}^K \rightarrow \Delta^{K-1}$$

Probability Simplex
 (Space of non-negative vectors that sum up to one)

In practice

$$p(y|\vec{x}) = \sigma(W\vec{x} + \vec{\omega}_0)$$

"Logits"

probabilities

- Naive Bayes (Variant II):

→ chooses a different class conditional

$$\approx D \quad p(\vec{x}|y=c) = \mathcal{N}(\vec{x}|\vec{\mu}_c, \Sigma_c)$$

Covariance matrices are individual but diagonal

Assumption: d features of a sample \vec{x} are conditionally independent

$$p(x_1, \dots, x_d|y=c) = \prod_{i=1}^d p(x_i|y=c)$$

Example for two classes $C = \{0, 1\}$

$$\approx D \quad a = \log \frac{p(\vec{x}|y=1) p(y=1)}{p(\vec{x}|y=0) p(y=0)}$$

$$= \dots = \vec{x}^T W_2 \vec{x} + \vec{\omega}_2^T \vec{x} + \omega_2$$

Slide 28

ADVANTAGE OF Naive Bayes:

→ due to the cond. independence of features

choose a (univariate) distribution for each feature

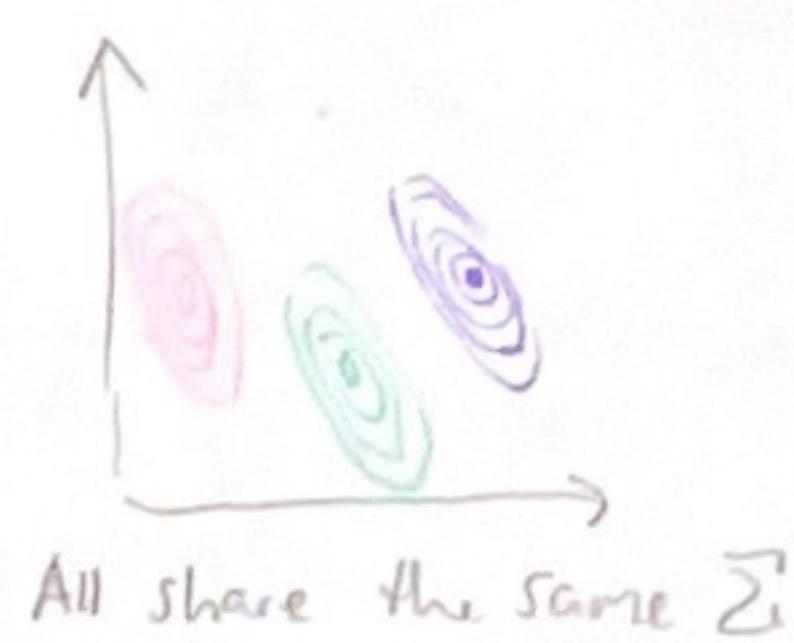
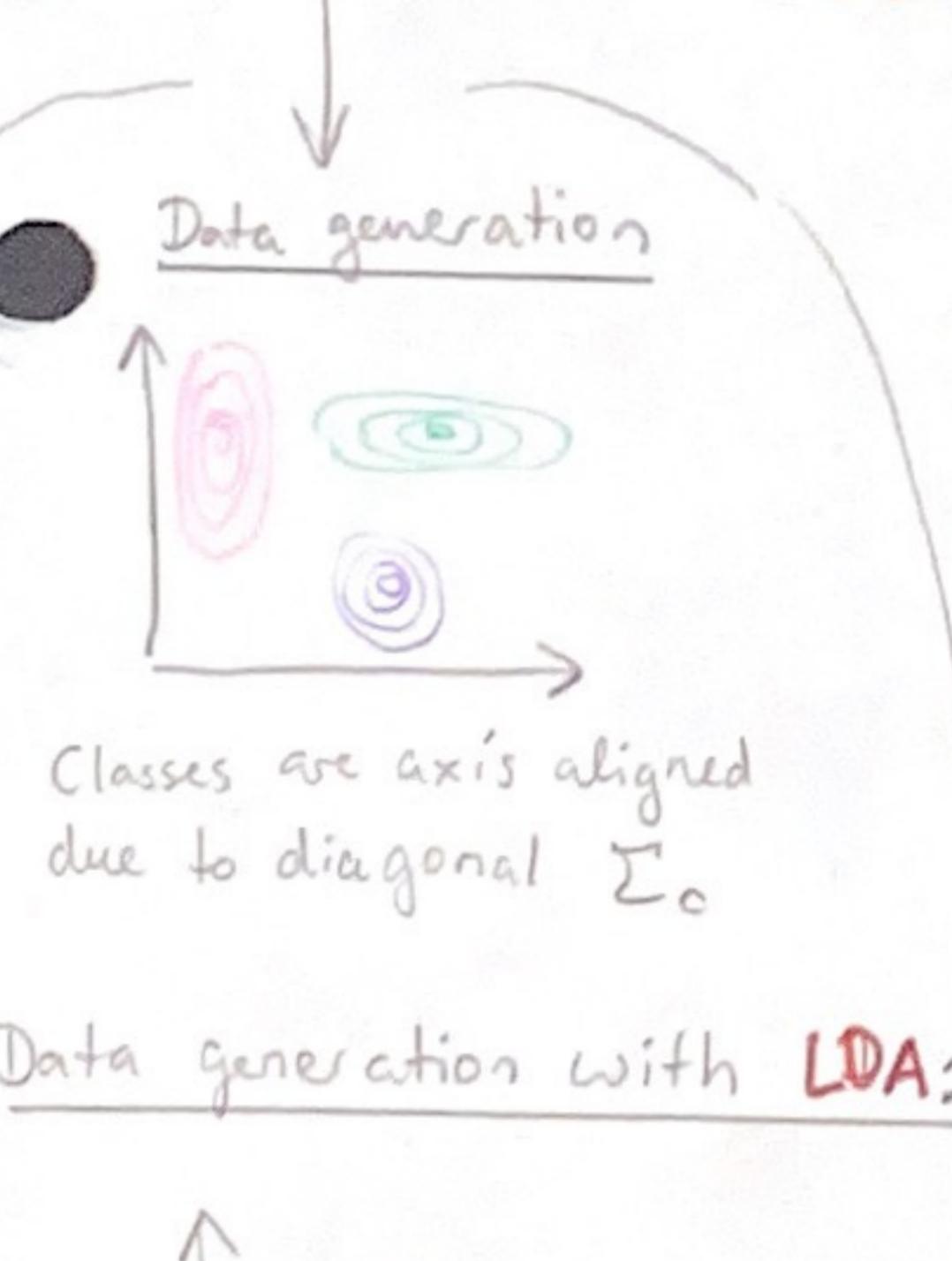
x_1 - Gaussian

x_2 - Categorical

= ... / = ... / = ...
 ↓ quadratic term

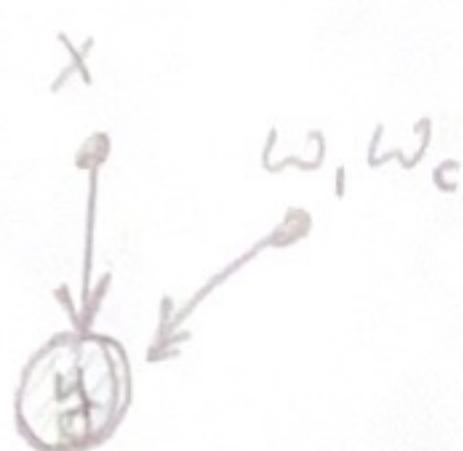
→ quadratic decision boundary

(LDA leads to lin. decision boundaries)



- Probabilistic discriminative model:

→ Instead of \vec{w}, w_0 depending on μ_0, μ_1, Σ , they are chosen directly



Example: logistic regression (binary)

$$y | \vec{x} \sim \text{Bernoulli}(\sigma(\vec{w}^T \vec{x} + w_0)) \text{ with } \sigma(a) = \frac{1}{1+e^{-a}}$$

$$\begin{aligned} \Rightarrow p(\vec{y} | \vec{w}, X) &= \prod_{i=1}^N p(y_i | \vec{x}_i, \vec{w}) \quad \text{absorbs } w_0 \\ &= \prod_{i=1}^N \sigma(\vec{w}^T \vec{x}_i)^{y_i} (1 - \sigma(\vec{w}^T \vec{x}_i))^{(1-y_i)} \end{aligned}$$

Error func

$$\Rightarrow E = -\log p(\vec{y} | X, \vec{w})$$

$$= -\sum_{i=1}^N \underbrace{y_i \log \sigma(\vec{w}^T \vec{x}_i)}_{:= \text{Binary cross entropy loss}} + (1-y_i) \log (1 - \sigma(\vec{w}^T \vec{x}_i))$$

equivalent
to MLE

$$\Rightarrow \vec{w}^* = \underset{\vec{w}}{\operatorname{argmin}} E(\vec{w})$$

→ No closed form (use optimization)

With weight regularization: $E(\vec{w}) = -\log p(\vec{y} | \vec{w}, X) + \lambda \|\vec{w}\|_F^2$

Example: Multiclass Logistic regression

$$p(Y | \vec{w}, X) = \prod_{i=1}^N p(\vec{y}_i | \vec{x}_i, \vec{w}) = \prod_{i=1}^N \prod_{c=1}^C p(y_i=c | \vec{x}_i, \vec{w}_c)$$

Binary matrix / $Y \in \{0, 1\}^{N \times C}$

→ One-hot-encoding

$$\Rightarrow E(\vec{w}) = -\log p(Y | \vec{w}, X) = -\sum_{i=1}^N \sum_{c=1}^C y_{ic} \log p(y_i=c | \vec{x}_i, \vec{w}_c)$$

$$= -\sum_{i=1}^N \sum_{c=1}^C y_{ic} \underbrace{\frac{\exp(\vec{w}_c^T \vec{x}_i)}{\sum_{c=1}^C \exp(\vec{w}_c^T \vec{x}_i)}}_{\text{Softmax}}$$

:= Cross Entropy loss
w/ one-hot-encoding

$$\hookrightarrow y_{ic} = \begin{cases} 1, & \text{if sample } i \text{ in class } c \\ 0, & \text{else} \end{cases}$$

ML 06 : Optimization

(9)

↪ General task: Find solution $\hat{\theta}^*$ minimizing function

$$f: \hat{\theta}^* = \underset{\hat{\theta} \in X}{\operatorname{argmin}} f(\hat{\theta})$$

domain of $\hat{\theta}$ (constraints on parameters) objective func.

- Convexity:

→ Set X is convex, if $\forall x, y \in X$: $\lambda x + (1-\lambda)y \in X$ for $\lambda \in [0,1]$

\Rightarrow



Not convex:



Region above a convex function is a convex set.

→ $f(x)$ is a convex function on convex set X iff

$$\left[\forall x, y \in X : \underbrace{\lambda f(\vec{x}) + (1-\lambda)f(\vec{y})}_{\text{for } \lambda \in [0,1]} \geq f(\lambda \vec{x} + (1-\lambda)\vec{y}) \right]$$

not convex

→ For a convex func: Each local minimum is a global minimum

$$f(\vec{y}) - f(\vec{x}) \geq \frac{f((1-t)\vec{x} + t\vec{y}) - f(\vec{x})}{t}$$

$$\text{for } t \rightarrow 0: \frac{[f(\vec{x} + \epsilon) - f(\vec{x})]}{\epsilon \cdot (y-x)}$$

If this condition holds

for all $x, y \in X$, then

f is convex on X .

$$f(\vec{y}) - f(\vec{x}) \geq (\vec{y} - \vec{x})^\top \nabla f(\vec{x})$$

- Verifying convexity:

- of func. {
- ① A twice diff. func. is convex iff its Hessian is positive semi-definite,
 - ② Prove condition above.
 - ③ Show that the func. is obtained from simple convex func. that preserve convexity

Convexity-preserving Operations:

$$h(\vec{x}) = f_1(\vec{x}) + f_2(\vec{x})$$

$$h(\vec{x}) = \max(f_1(\vec{x}), f_2(\vec{x}))$$

$$h(\vec{x}) = c \cdot f_1(\vec{x}) \quad \text{if } c > 0$$

$$h(\vec{x}) = c \cdot f_1(\vec{x}) \quad \text{if } c < 0$$

$$h(\vec{x}) = f_1(A\vec{x} + \vec{b})$$

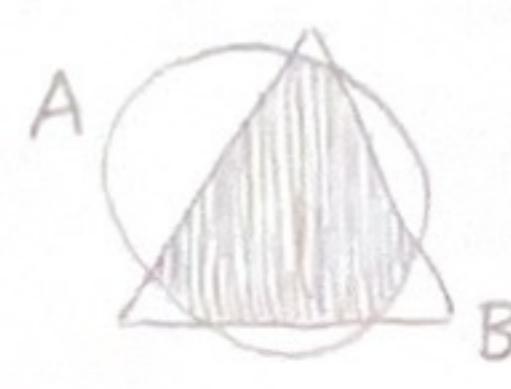
$$h(\vec{x}) = m(f_1(\vec{x}))$$

convex & non-decreasing

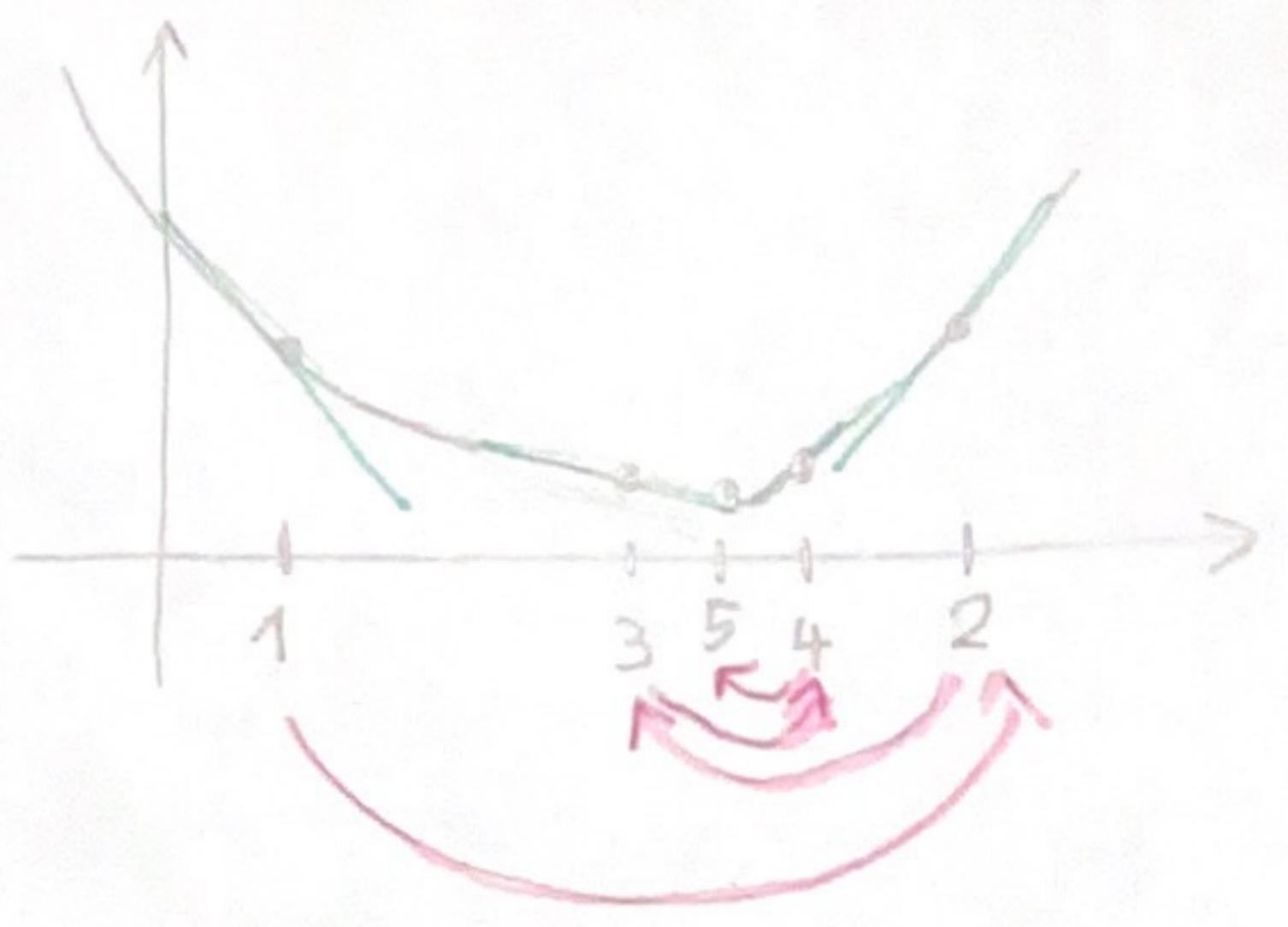
of sets

- {
- ① Prove definition
 - ② Apply intersection rule

→ If A and B are convex sets,
 $A \cap B$ is convex



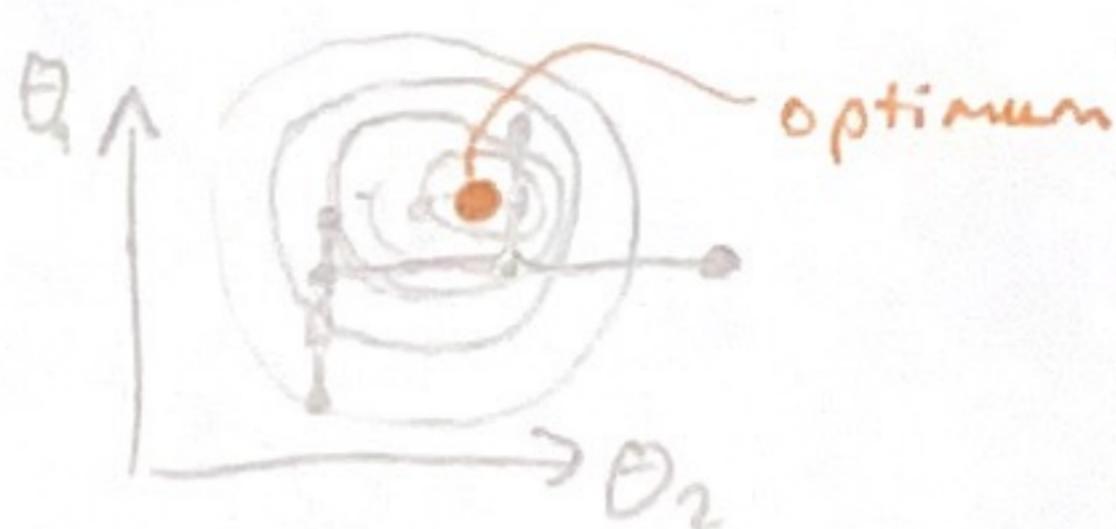
- Optimization in a One-dim. problem:



Goal: For diff. f search for θ with $\nabla f(\theta) = 0$

→ Do interval bisection

(in a multivariate problem do interval bisection in one dim., then in the next one, ...)

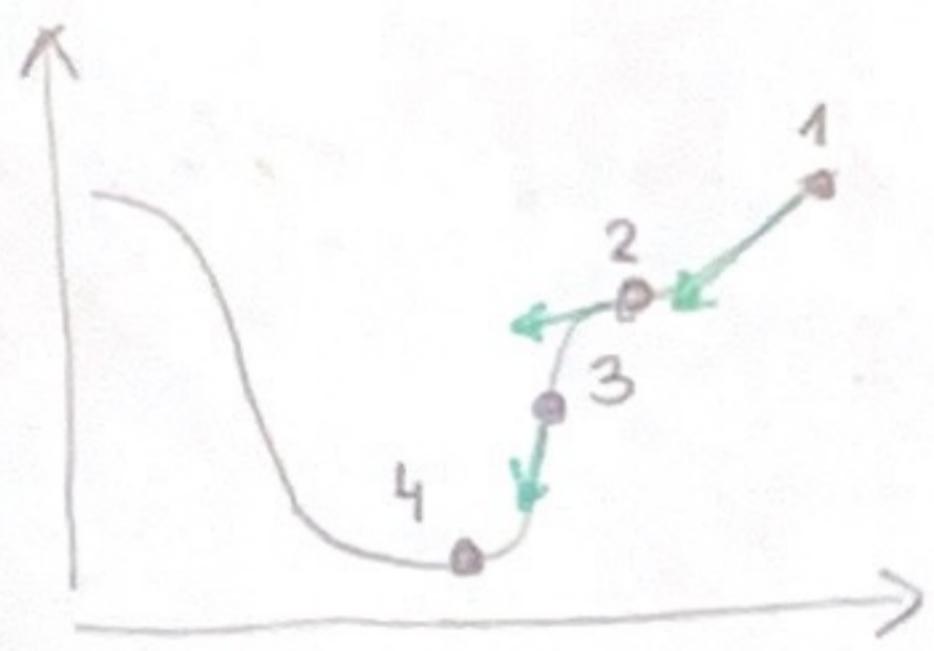


- Optimization in Multi-dim. problems: Gradient descent (GD)

→ Idea: Gradient points in the direction of steepest ascent

Approximates the objective function locally

GD with Line Search



→ works on non-convex functions

given a starting point $\vec{\theta} \in \text{Dom}(f)$

repeat

1. $\Delta \vec{\theta} = -\nabla f(\vec{\theta})$ Get direction of steepest descent.

2. Line Search. $t^* = \underset{t>0}{\operatorname{argmin}} f(\vec{\theta} + t\Delta \vec{\theta})$ Solve one-dim. Line-Search problem

3. Update. $\vec{\theta} = \vec{\theta} + t^* \cdot \Delta \vec{\theta}$

until stopping criterion is met

Residual error $\rho = f(\vec{\theta}^{(k)}) - p^* \leq \varepsilon$ after at most $\frac{\log((f(\vec{\theta}^{(0)}) - p^*)\varepsilon)}{\log(1/c)}$ steps

→ ~~Fast~~ Parallel / distributed implementation

Often problems are in the form $f(\vec{\theta}) = \sum_i L_i(\vec{\theta}) + g(\vec{\theta})$

iterates over data

decompose

→ ~~distribute~~ distribute gradient by using the sum rule

→ Faster algorithm (actually the standard one)

→ too small: slow convergence

→ Set a learning rate τ

→ too high: alg. oscillates w/o converg.

→ use learning rate schedule

→ Momentum GD (incorporates history of previous gradients)

$$\vec{m}_t \leftarrow \tau \nabla f(\vec{\theta}_t) + \gamma \vec{m}_{t-1}$$

$$\vec{\theta}_{t+1} \leftarrow \vec{\theta}_t + \vec{m}_t$$

accelerates the search in directions where many previous gradient points in the same direction

- AdaGrad:

- Different learning rate per parameter

- Adam: $\vec{m}_t = \beta_1 \vec{m}_{t-1} + (1-\beta_1) \nabla f(\vec{\theta}_t)$ (1. Moment)

$\vec{v}_t = \beta_2 \vec{v}_{t-1} + (1-\beta_2) (\nabla f(\vec{\theta}_t))^2$ (2. Moment)

$$\Rightarrow \vec{\theta}_{t+1} = \vec{\theta}_t - \frac{\tau}{\sqrt{\vec{v}_t} + \epsilon} \hat{m}_t \quad \text{with } \hat{m}_t = \frac{\vec{m}_t}{1-\beta_1^t} \quad , \quad \hat{v}_t = \frac{\vec{v}_t}{1-\beta_2^t}$$

GD = First Order Method

- Newton - Method (Second-Order-Method)

Givens - Convex function f

- Nonneg. second derivative: $\underbrace{\nabla^2 f(\vec{\theta}) \succeq 0}_{:= \text{Hessian is positive semi-definite}}$

\vdash Hessian is positive semi-definite

\rightarrow Taylor Expansion at point $\vec{\theta}_t$

$$f(\vec{\theta}_t + \vec{\delta}) = \underbrace{f(\vec{\theta}_t) + \vec{\delta}^T \nabla f(\vec{\theta}_t)}_{:= g(\vec{\delta})} + \frac{1}{2} \vec{\delta}^T \nabla^2 f(\vec{\theta}_t) \vec{\delta} + O(\vec{\delta}^3)$$

Minimize $\Rightarrow \vec{\theta}_{t+1} \leftarrow \vec{\theta}_t - \underbrace{[\nabla^2 f(\vec{\theta}_t)]^{-1}}_{\tau := \text{learning rate}} \nabla f(\vec{\theta}_t) \quad \Rightarrow \nabla g(\vec{\delta}) = 0$

- \rightarrow for low-dimensional methods
- ⊕ good convergence rate
 - ⊕ Gradient requires $O(d)$ data (memory)
 - ⊖ Hessian requires $O(d^2)$ data
 - ⊖ Update step is $O(d^3)$ & tricky to parallelize
↳ inverse

- Stochastic Optimization: (~~minimizes~~ with computational costs)

\rightarrow Stochastic Gradient Descent

$\text{SGD computes noisy estimate of gradients on subset (minibatch)}$

$$\frac{1}{n} \left(\sum_{i=1}^n L_i(\vec{\theta}) \right) = \underset{i \in \{1, \dots, n\}}{\mathbb{E}} [L_i(\vec{\theta})] \approx \left(\frac{1}{|S|} \sum_{j \in S} (L_j(\vec{\theta})) \right) \quad \text{with } S \subseteq \{1, \dots, n\}$$

write Loss as empirical risk minimization

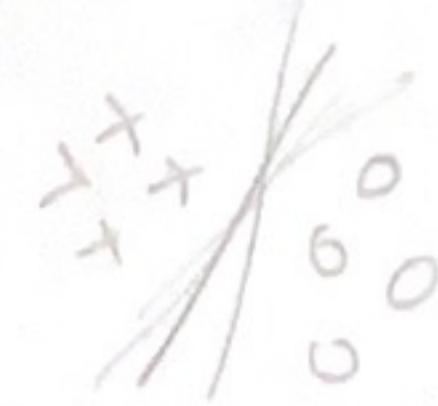
Scale Loss on subset such that it matches the loss on all data

Approximate estimation on subset S of the data

$$\Rightarrow \sum_{i=1}^n L_i(\vec{\theta}) \approx \frac{n}{|S|} \sum_{j \in S} L_j(\vec{\theta})$$

Example: Perceptron (binary lin. classifier)

$$\delta(\vec{x}) = \begin{cases} 1 & \text{if } \vec{w}^T \vec{x} + b > 0 \\ -1 & \text{else} \end{cases}$$



$$\Rightarrow \vec{w}^*, b^* = \underset{\vec{w}, b}{\operatorname{argmin}} \sum_i L(y_i, \vec{w}^T \vec{x}_i + b)$$

↓ prediction

$$\text{Hinge-loss } L(u, v) = \max(0, 1 - uv)$$

$$\text{ground truth} = \begin{cases} 1 - uv, & \text{if } uv < 1 \text{ incorrect prediction} \\ 0, & \text{else correct prediction} \end{cases}$$

$$\Rightarrow \nabla_w L(y_i, \vec{w}^T \vec{x}_i + b) = \begin{cases} -y_i \vec{x}_i, & \text{if } uv < 1 \\ 0, & \text{else} \end{cases}$$

$$\nabla_b L(y_i, \vec{w}^T \vec{x}_i + b) = \begin{cases} -y_i, & \text{if } uv < 1 \\ 0, & \text{else} \end{cases}$$

SGD
 with minibatch of size 1
 initialize $\vec{w} = 0$ and $b = 0$
 repeat
 if $y_i \cdot (\vec{w}^T \vec{x}_i + b) < 1$ then
 $\vec{w} \leftarrow \vec{w} + \tau \cdot n \cdot y_i \cdot \vec{x}_i$ and $b \leftarrow b + \tau \cdot n \cdot y_i$
 end if
 until all classified correctly

SGD: Expectation of residual error decreases with speed
 # of iterations $t \sim E[\rho]^{-1}$ residual

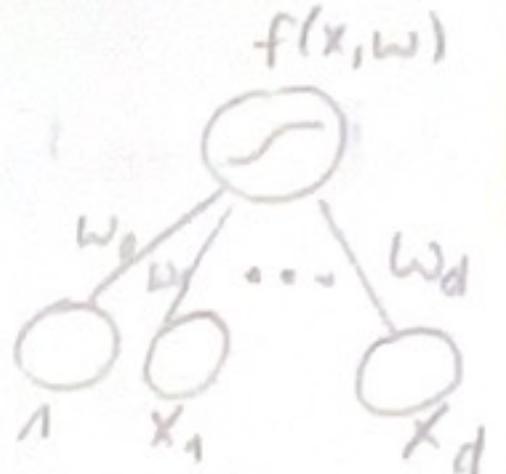
GD: speed has $t \sim \log \rho^{-1}$ } GD has faster convergence speed, but each iteration takes longer

- Distributed Learning:

- Data Parallelism (multiple workers operate on the same model and update parameters)
- Model parallelism (Use inherent model parallelisms and share on multiple workers (e.g. matrix ops) and share data)

ML 07 : Deep Learning I

11



- How to handle basis functions? *non-linearity*

→ Regular logistic regression only for lin. separable data $y|x \sim \text{Bernoulli}(o(w^T x))$

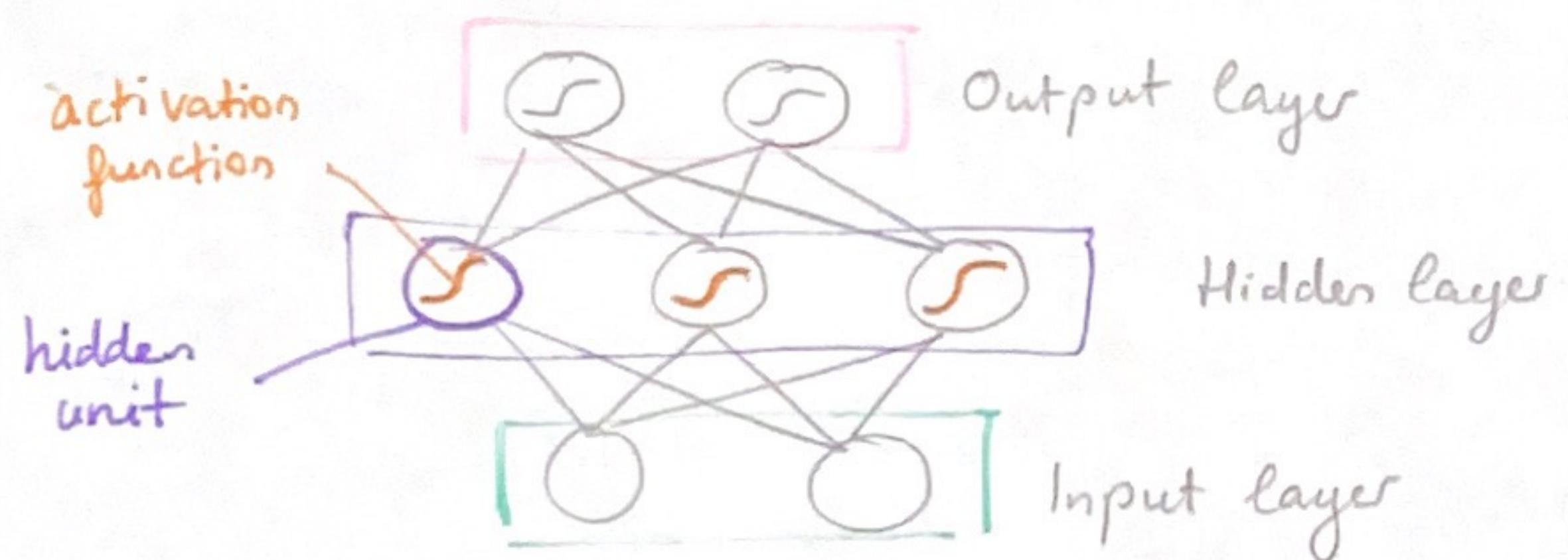
→ Use basis functions: $f(x, y) = \sigma\left(w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\vec{x})\right) = \sigma(\vec{w}^\top \vec{\phi}(\vec{x}))$
 ↳ maps samples into a space where they're lin. separable *basis function*

⇒ Then find the best weights with the binary cross entropy loss

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{n=1}^N - (y_n \log f(x_n, w) + (1-y_n) \log(1-f(x_n, w)))$$

⇒ Is feed-forward NN with 1 hidden layer

- Multi-layered perceptron: (MLP) (feed-forward; fully-connected)

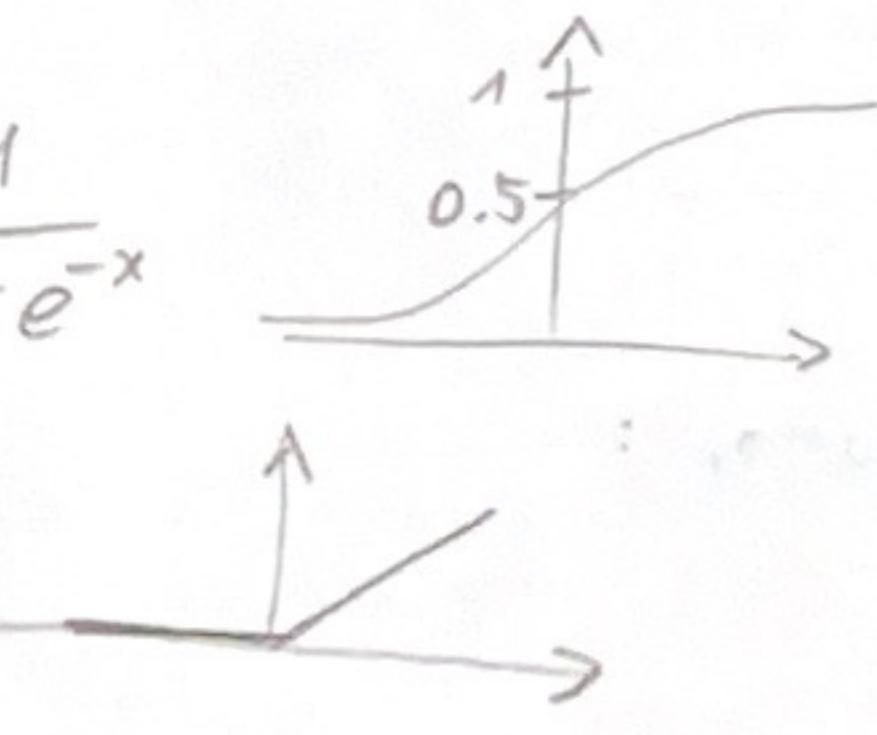


$$f(x, w) = \sigma_2(w_2 \sigma_1(w_1 \sigma_0(w_0 x)))$$

w/ 2 hidden layers

Activation func.:

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



$$\text{ReLU: } \max(0, x)$$



- Why use nonlinear activation funcs?

→ Lin. layers result in lin. transf.: $f(x, w) = w_L (w_{L-1} (\dots (w_0 \vec{x}) \dots)) = w^\top \vec{x}$

- Universal approximation theorem:

→ With enough hidden units an MLP with a lin. output layer and one hidden layer can approximate any continuous function (over a closed subset of \mathbb{R}^D)

- Why add more layers?

→ Use more layers with fewer hidden units instead of one wide layer
 ⇒ fewer parameters

→ In practice: Deep networks are easier to train

- What needs to be changed in a NN for different tasks?
 - the activation func in the final layer and the loss func

Prediction target	$p(y x)$	Final layer	Loss function
Binary	Bernoulli	Sigmoid	Binary cross entropy
Discrete	Categorical	Softmax	Cross entropy
Continuous	Gaussian	Identity	Squared error

The loss func
should allow
to do gradient
based training

- How to compute the gradient for GD?

1) Work out manually (by hand)

$$2) \text{ Numerically: } \frac{\partial E_n}{\partial w_{ij}} = \frac{E_n(w_{ij} + \epsilon) - E_n(w_{ij})}{\epsilon} + O(\epsilon)$$

→ expensive and
maybe not stable

Computing this means
evaluating the NN

→ $O(|W|)$ for every
parameter → $O(|W|^2)$

3) Symbolic differentiation: Automates deriving if by hand.
→ "Writing down" expression is expensive

4) Automatic differentiation / Backpropagation

→ Eval. $\nabla_w E(W)$ at current W

→ $O(|W|)$

- Backprop:

Example: $f(x) = \frac{2}{\sin(\exp(-x))} = d(c(b(a(x))))$

1. Write func as
composition of
modules

Comp. graph:



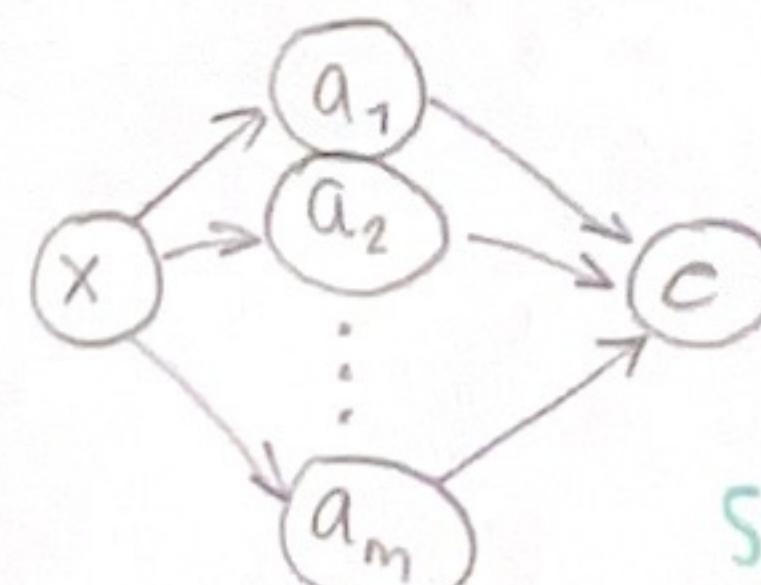
2. Do a forward pass and cache
intermediate values for given input \vec{x}

3. Compute local derivatives for \vec{x}

4. Compute global derivative

Derivative: $\frac{\partial f}{\partial x} = \frac{\partial d}{\partial c} \cdot \frac{\partial c}{\partial b} \cdot \frac{\partial b}{\partial a} \cdot \frac{\partial a}{\partial x}$

- Multivariate chain rule:



$$\frac{\partial c}{\partial x} = \sum_{l=1}^m \frac{\partial c}{\partial a_l} \cdot \frac{\partial a_l}{\partial x}$$

Sum of the derivatives along all paths.

- Jacobian & Gradient:

$\vec{a}: \mathbb{R}^n \rightarrow \mathbb{R}^m$

Jacobian: $\frac{\partial \vec{a}}{\partial \vec{x}} = \begin{bmatrix} \frac{\partial a_1}{\partial x_1} & \dots & \frac{\partial a_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial a_m}{\partial x_1} & \dots & \frac{\partial a_m}{\partial x_n} \end{bmatrix}$

$= (\nabla_{\vec{x}} a_i)^T \rightarrow$ "the gradient is the transpose of the Jacobian"

$\nabla_{\vec{x}} a$ has the same shape as \vec{x}

Gradient: $\nabla_{\vec{x}} a = \left(\frac{\partial a}{\partial \vec{x}} \right)^T = \left[\frac{\partial a}{\partial x_1} \dots \frac{\partial a}{\partial x_m} \right]^T$

- Chain rule in matrix form:

$$\frac{\partial c}{\partial x_j} = \sum_{i=1}^m \frac{\partial c}{\partial a_i} \frac{\partial a_i}{\partial x_j} \quad \text{In matrix form} \quad \frac{\partial c}{\partial \vec{x}} = \underbrace{\frac{\partial c}{\partial \vec{a}}}_{= \nabla_{\vec{x}} c} \underbrace{\frac{\partial \vec{a}}{\partial \vec{x}}}_{= \left(\frac{\partial \vec{a}}{\partial \vec{x}} \right)^T} \quad \text{with } \left[\frac{\partial \vec{a}}{\partial \vec{x}} \right]_{ij} = \frac{\partial a_i}{\partial x_j}$$

→ Matrix calculus: $f: \mathbb{R} \rightarrow$

	scalar	vector	matrix
scalar	sc	vec	matr.
vector	vec	matr	3-way-tensor
matrix	matr	3-way-tensor	4-way-tensor

$\left(\frac{\partial \vec{a}}{\partial \vec{w}} \right)_{ijk} = \frac{\partial a_i}{\partial w_{jk}}$ e.g.

shape of the derivative of the func

- Accumulating the gradient (Example):

Given feed-forward NN

$$\begin{cases} \vec{a} = W\vec{x} + \vec{b} & \leftarrow \text{affine layer} \\ \vec{h} = \sigma(\vec{a}) \\ \hat{y} = V\vec{h} + c \\ E = (\hat{y} - y)^2 \end{cases}$$

→ In the affine layer calculate:

forward (W, \vec{x}, \vec{b}): compute $W\vec{x} + \vec{b}$

backward ($\frac{\partial E}{\partial \vec{a}}$): Compute
grads w.r.t. the
output of the
affine layer

would be 3-way-tensor

Derived on sl. 43

$$\frac{\partial E}{\partial W} = \frac{\partial E}{\partial \vec{a}} \frac{\partial \vec{a}}{\partial W} = \vec{x} \frac{\partial E}{\partial \vec{a}}$$

$$\frac{\partial E}{\partial \vec{x}} = \frac{\partial E}{\partial \vec{a}} \frac{\partial \vec{a}}{\partial \vec{x}} = \frac{\partial E}{\partial \vec{a}} W$$

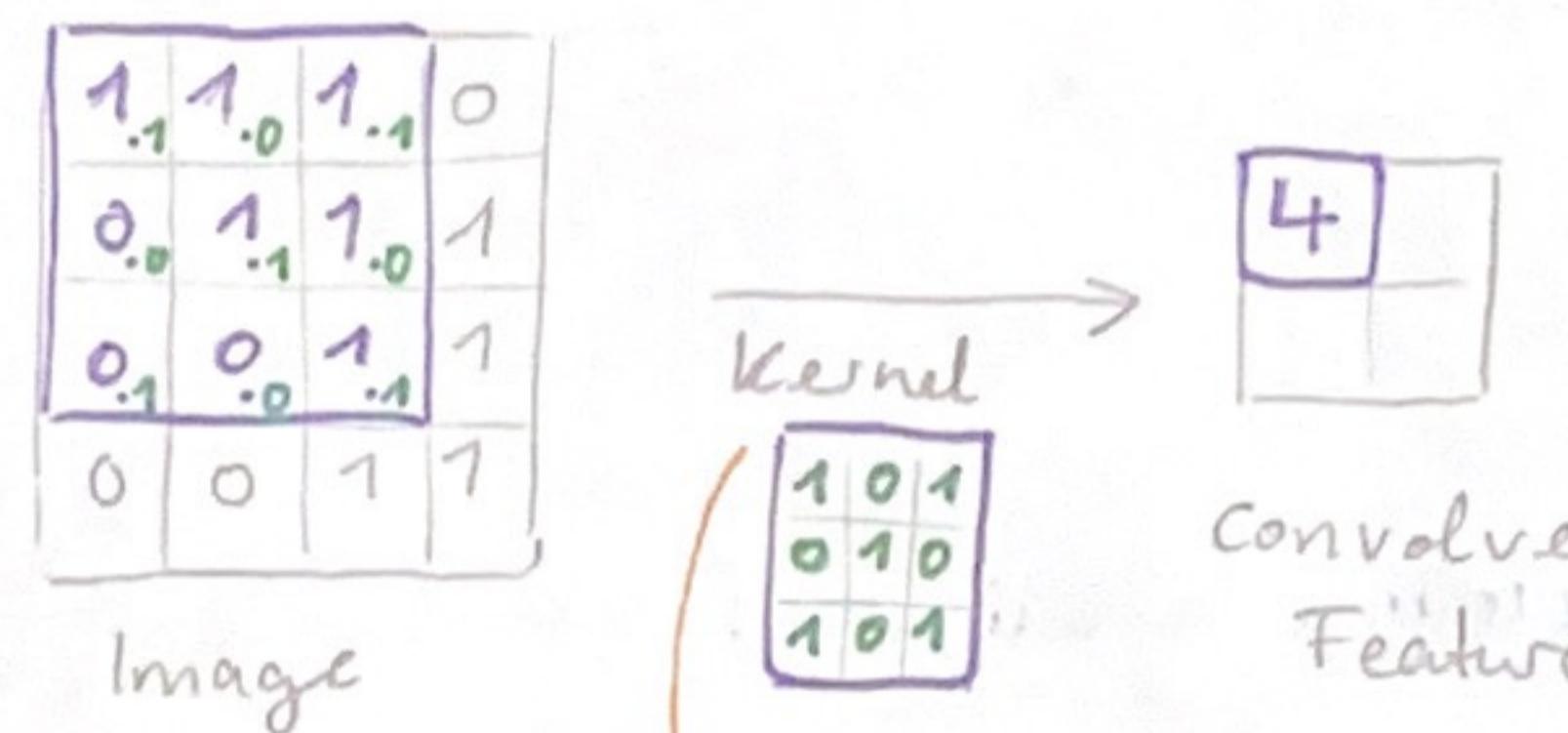
$$\frac{\partial E}{\partial \vec{b}} = \frac{\partial E}{\partial \vec{a}} \frac{\partial \vec{a}}{\partial \vec{b}} = \frac{\partial E}{\partial \vec{a}}$$

We can comp.
the gradient
w/o comp.
the Jacobian!

- Convolution:

- continuous definition: $(x * k)(t) = \int_{-\infty}^{\infty} x(\tau) k(t-\tau) d\tau$
"weighted avg of input x with weights/kernel/filter k at each point in time t "
- discrete definition: $(x * k)(t) = \sum_{\tau=-\infty}^{\infty} x(\tau) k(t-\tau)$
- 2-D "convolution" in CNNs: $\hat{x}(i,j) = \sum_{l=1}^L \sum_{m=1}^M x(i+l, j+m) k(l,m)$
 output of convolution \uparrow eval kernel over corresponding input
 is actually cross-correlation

Convolutions often act on multiple channels;
 $\Rightarrow L \times M \times C_{in} \times C_{out}$ parameters



Kernel weights are shared (& learned)

* Padding Schemes: \rightarrow VALID: No padding. Reduces output size to

$$D_{l+1} = \frac{(D_l - K)}{\text{input size along a dim}} + 1$$

Kernel width

\rightarrow SAME: Add padding to preserve the input size
 $P = \lfloor K/2 \rfloor$

\rightarrow FULL: Add $K-1$ values on each side
 \rightarrow increase output size

* Strides: (the distance between positions the kernel is applied)

$$\Rightarrow D_{l+1} = \left\lfloor \frac{D_l + 2P - K}{S} \right\rfloor + 1$$

\rightarrow Strides $S > 1$ downsample the signal

* Pooling: (calculate summary statistics in sliding window \rightarrow downsampling)
 \rightarrow Max / mean / L_p -norm pooling

- Weight initialization:

- * Weight symmetry: Hidden units with the same weights & biases will always get the exact same gradients
 \rightarrow never learn different features
 \rightarrow break symmetry w/ small random values
- * Weight scale: Activation func may saturate for hidden unit w/ large fan-in
 \rightarrow vanishing or exploding gradients w/ wrong weight scales
 \rightarrow init. weights w/ good mean & variance

- Xavier-Glorot initialization: \rightarrow idea: preserve mean & variance of an incoming i.i.d. signal in the forward & backward pass

\Rightarrow By controlling the mean and variance of the weight distribution



In regression tasks:

For initial weights of last layer, consider the target scale

Weight matrices:

$$\rightarrow \text{Normal with } \begin{cases} \mu = 0 \\ \text{Var}(W) = \frac{2}{\text{fan-in} + \text{fan-out}} \end{cases}$$

- Vanishing / Exploding gradients:

Example: Say, a NN multiplies the input t times by W

$$\Rightarrow W^t = (V \text{diag}(D) V^{-1})^t = V \text{diag}(W)^t V^{-1}$$

$\Rightarrow \begin{cases} \text{if } D_{ii} < 1, \text{ then } D_{ii}^t \text{ is small} \rightarrow \text{vanishing gradient} \\ \text{if } D_{ii} > 1, \text{ then } D_{ii}^t \text{ explodes} \rightarrow \text{exploding gradient} \end{cases}$

* Saturated activation funcs also cause exploding/vanishing gradients.

\Rightarrow Solutions: - Use batch-norm, other activations, etc.
- gradient clipping

- Regularization: \rightarrow prevents overfitting

* L₂-norm: enforces small weights

* L₁-norm: enforces sparsity

* Data augmentation, injecting noise, parameter sharing, dropout

randomly disable neurons during

& forward/backward pass

\rightarrow sampling from different architectures with shared weights

Machine Learning 09: SVM and Kernels

14

- Linear classifier:

$$h(\vec{x}) = \text{sign}(\vec{w}^T \vec{x} + b) = \begin{cases} -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ +1 & \text{if } z > 0 \end{cases}$$



a wide margin around the dividing hyperplane makes it more likely to correctly classify new samples
 \Rightarrow Find hyperplane with maximum margin

- Linear classifier w/ margin:

* Require $\vec{w}^T \vec{x} + (b-s) > 0$ for all points from class +1

and $\vec{w}^T \vec{x} + (b+s) < 0$ for all points from class -1

$$\Rightarrow \text{Size of Margin } m = -\underbrace{\frac{b-s}{\|\vec{w}\|}}_{\substack{\text{signed distance from} \\ \text{the origin to the} \\ \text{purple hyperplane}}} - \underbrace{\left(-\frac{b+s}{\|\vec{w}\|}\right)}_{\substack{\text{---} \\ \text{green} \\ \text{---}}} = \underbrace{\frac{2s}{\|\vec{w}\|}}_{\substack{\text{---} \\ \text{---}}}$$

signed distance from
the origin to the
purple hyperplane

m only depends on this ratio.
w.l.o.g., we set s=1

\Rightarrow scale $\|\vec{w}\|$ and b

$s=1$

\Rightarrow Require the constraints:

$$\vec{w}^T \vec{x}_i + b \begin{cases} \geq +1 & \text{for } y_i = +1 \\ \leq -1 & \text{for } y_i = -1 \end{cases}$$

$$\Leftrightarrow y_i (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \text{for all } i$$

with the optimization problem $\underset{y, b, \vec{w}}{\arg \max} (m)$

- SVM's optimization problem:

$$\left\{ \begin{array}{l} \vec{w}^*, b^* = \underset{\vec{w}, b}{\arg \max} \frac{2}{\sqrt{\vec{w}^T \vec{w}}} \\ \text{constrained convex optimization problem} \\ \downarrow \text{squaring is a monotonic func} \\ = \underset{\vec{w}, b}{\arg \min} \frac{\vec{w}^T \vec{w}}{2} \\ = \frac{1}{2} \|\vec{w}\|^2 := f_0(\vec{w}, b) \\ \text{constraint: } f_i(\vec{w}, b) = y_i (\vec{w}^T \vec{x}_i + b) - 1 \geq 0 \quad \text{for } i=1, \dots, N \end{array} \right.$$

- Optimization w/ inequality constraints:

Given $f_0: \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_i: \mathbb{R}^d \rightarrow \mathbb{R}$

minimize $f_0(\vec{\theta})$ subject to $f_i(\vec{\theta}) \leq 0 \quad \forall i$

* A point $\vec{\theta}$ is **feasible** if it satisfies the constraints

* Optimal value is the **minimum** p^* and the point where it's obtained is the **minimizer** $\vec{\theta}^*$

$$\Rightarrow p^* = f_0(\vec{\theta}^*)$$

Lagrange multiplier

Lagrangian: $L(\vec{\theta}, \vec{\alpha}) = f_0(\vec{\theta}) + \sum_{i=1}^M \alpha_i f_i(\vec{\theta})$

Lagrange dual function: $g(\vec{\alpha}) = \min_{\vec{\theta} \in \mathbb{R}^d} L(\vec{\theta}, \vec{\alpha}) = \min_{\vec{\theta} \in \mathbb{R}^d} \left(f_0(\vec{\theta}) + \sum_{i=1}^M \alpha_i f_i(\vec{\theta}) \right)$

$\vec{\theta}$ is unconstrained in this eq., since f_i is contained
is concave in $\vec{\alpha}$

$\begin{matrix} \uparrow p^* \\ -g(\vec{\alpha}) \end{matrix} \Rightarrow$ For every $\vec{\alpha}$, the unconstrained $g(\vec{\alpha})$ is a lower bound on the optimal value of the constrained problem:

$$\min_{\substack{\vec{\theta} \in \mathbb{R}^d \\ f_i(\vec{\theta}) \leq 0}} f_0(\vec{\theta}) = f_0(\vec{\theta}^*) \geq f_0(\vec{\theta}^*) + \underbrace{\sum_{i=1}^M \alpha_i f_i(\vec{\theta}^*)}_{\leq 0, \text{ since } \alpha_i \geq 0, \text{ by def. }} = L(\vec{\theta}^*, \vec{\alpha}) \geq \min_{\vec{\theta} \in \mathbb{R}^d} L(\vec{\theta}, \vec{\alpha}) = g(\vec{\alpha})$$

$$\Rightarrow \forall \vec{\alpha} \quad f_0(\vec{\theta}^*) \geq g(\vec{\alpha})$$

\Rightarrow the maximum d^* of the Lagrangian dual problem is the best possible lower bound on p^*

* Weak duality: $g(\vec{\alpha}) \leq p^* \quad \forall \vec{\alpha} \geq \vec{0}$

\Rightarrow duality gap: $p^* - d^* \geq 0$

Strong duality: $\underbrace{g(\vec{\alpha}^*)}_{\substack{\text{max of Lagrange} \\ \text{dual problem}}} = p^* = \underbrace{f_0(\vec{\theta}^*)}_{\substack{\text{minimum of} \\ \text{original problem}}}$

Lagrange dual problem:

$$\underset{\vec{\alpha}}{\text{maximize}} \quad g(\vec{\alpha}) \quad \text{subject to } \alpha_i \geq 0 \quad \forall i$$

- Constraint optimization problem of SVM:

1. Lagrangian: $L(\vec{w}, b, \vec{\alpha}) = \frac{1}{2} \vec{w}^T \vec{w} - \sum_{i=1}^N \alpha_i [y_i (\vec{w}^T \vec{x}_i + b) - 1]$

2. Minimize w.r.t. \vec{w} and b

$$\nabla_{\vec{w}} L(\vec{w}, b, \vec{\alpha}) = \vec{w} - \sum_{i=1}^N \alpha_i y_i \vec{x}_i = 0 \Rightarrow \vec{w} = \sum_{i=1}^N \alpha_i y_i \vec{x}_i \quad (1)$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i \neq 0$$

$$\Rightarrow g(\vec{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \vec{x}_i^T \vec{x}_j$$

subject to $\sum_{i=1}^N \alpha_i y_i = 0, \alpha_i > 0 \quad \forall i$

3. Solve $\underset{\vec{\alpha}}{\operatorname{argmax}} g(\vec{\alpha})$.

\Rightarrow Find $\vec{\alpha}^*$ via a QP-solver

\Rightarrow Get \vec{w} from (1)

\Rightarrow Get b from:

$$f_0(\vec{\theta}^*) + \underbrace{\sum_i \alpha_i f_i(\vec{\theta}^*)}_{\text{for an } \vec{x}_i \text{ for which } \alpha_i \neq 0} = f_0(\vec{\theta}^*) \quad \left. \begin{array}{l} \text{strong} \\ \text{duality} \\ \downarrow \\ \alpha_i^* f_i(\vec{\theta}^*) = 0 \end{array} \right\} \text{Complementary slackness condition: holds for strong dualities:}$$

$\vec{x}_i \neq 0$, the corresponding constraint $f_i(\vec{w}, b)$ must be zero

4. From $\alpha_i^* f_i(\vec{\theta}^*) = 0 \Rightarrow \alpha_i [y_i (\vec{w}^T \vec{x}_i + b) - 1] = 0 \quad \forall i$

(a training sample only contributes to \vec{w} if $\alpha_i \neq 0$, such that $y_i (\vec{w}^T \vec{x}_i + b) = 1$ (it lies on the margin!))

\Rightarrow Then, sample x_i with $\alpha_i \neq 0$ is called support vector

5. \Rightarrow Classification:

$$h(\vec{x}) = \operatorname{sign} \left(\sum_{i=1}^N \alpha_i y_i \vec{x}_i^T \vec{x} + b \right)$$

Since most α_i are zero, we only store the \vec{x}_i where $\alpha_i \neq 0$. ("solution is sparse")

- Soft margin support vector machines: "can handle noisy data/outliers"

Idea: Relax constraints, but punish relaxation

• \rightarrow slack variable $\xi_i \geq 0$ for every sample \vec{x}_i
 (the distance by which the constraint is violated in units of $\|\vec{w}\|$.)

$$\Rightarrow \vec{w}^T \vec{x}_i + b \begin{cases} \geq 1 - \xi_i & \text{for } y_i = +1 \\ \leq -1 + \xi_i & \text{for } y_i = -1 \end{cases} \quad = y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i \geq 0$$

$\xi = 0.3$, since $\vec{w}^T \vec{x} + b = y$

\Rightarrow Cost function becomes:

$$f_0(\vec{w}, b, \xi) = \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \xi_i$$

(Constrained)

Optimization problem: minimize $f_0(\vec{w}, b, \xi) = \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \xi_i$

$$\left. \begin{array}{l} \text{subject to} \\ \bullet y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i \geq 0 \\ \bullet \xi_i \geq 0 \end{array} \right\} \forall i$$

$$\Rightarrow \text{Lagrangian: } \textcircled{1} L(\vec{w}, b, \xi, \alpha, \mu) = \underbrace{\frac{1}{2} \vec{w}^T \vec{w}}_{f_0(\vec{w})} + C \sum_{i=1}^N \xi_i - \underbrace{\sum_{i=1}^N \alpha_i (y_i(\vec{w}^T \vec{x}_i + b) - 1 + \xi_i)}_{\text{fi}} - \underbrace{\sum_{i=1}^N \mu_i \xi_i}_{\text{fi}}$$

two sets of constraints

$$\textcircled{2} \text{ Min. } L(\vec{w}, b, \xi, \alpha, \mu) \text{ w.r.t. } \vec{w}, b, \xi$$

$$\nabla_{\vec{w}} L = \vec{w} - \sum_{i=1}^N \alpha_i y_i \vec{x}_i \stackrel{!}{=} 0, \quad \frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i \stackrel{!}{=} 0$$

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \mu_i \stackrel{!}{=} 0 \Rightarrow \alpha_i = C - \mu_i$$

$$\text{Dual feasibility: } \mu_i, \alpha_i \geq 0 \quad \xrightarrow{\quad} \quad 0 \leq \alpha_i \leq C$$

$$\textcircled{3} \text{ Dual problem: maximize } g(\vec{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \vec{x}_i^T \vec{x}_j$$

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

\rightarrow Hinge loss formulation:

From the constrained optimization problem: $\xi_i = \begin{cases} 1 - y_i(\vec{w}^T \vec{x}_i + b), & \text{if } y_i(\vec{w}^T \vec{x}_i + b) \leq 1 \\ 0 & \text{else} \end{cases}$

"SVMs are L2-regularized perceptron with a hinge-loss function"

$$= \max(0, 1 - y_i(\vec{w}^T \vec{x}_i + b))$$

$$\Rightarrow \text{Rewrite as unconstrained problem: } \min_{\vec{w}, b} \frac{1}{2} \vec{w}^T \vec{w} + C \sum_{i=1}^N \max(0, 1 - y_i(\vec{w}^T \vec{x}_i + b))$$

\rightarrow optimize with GD

- Feature space:

- * We use basis functions $\Phi(\cdot)$ to make models nonlinear and classify non-lin. separable data

$$\phi: \mathbb{R}^D \rightarrow \mathbb{R}^M, \vec{x}_i \mapsto \phi(\vec{x}_i)$$

- Kernels:

Rewrite the Lagrange dual function as

$$g(\vec{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j k(\vec{x}_i, \vec{x}_j) \quad \text{with } k(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i)^T \phi(\vec{x}_j).$$

*K
referred to as the kernel trick*

- Advantage of kernels:
- 1) Encode similarity between arbitrary data (numerical, objects, graphs, ...)
 - 2) Easier to evaluate...

Mercer's theorem

A kernel is **valid**, if it gives rise to a symmetric, PSD kernel matrix K for some input data X .

$$\text{Kernel matrix} = \text{Gram matrix} = K = \begin{pmatrix} k(x_1, x_1) & \dots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_N) \end{pmatrix}$$

or, a kernel is valid, if it can be represented as a dot product in some feature space.

⇒ if the kernel is not valid, our optimization problem may become not convex

* Kernel preserving functions:

- $k(\vec{x}_1, \vec{x}_2) = k_1(\vec{x}_1, \vec{x}_2) + k_2(\vec{x}_1, \vec{x}_2)$
- $k(\vec{x}_1, \vec{x}_2) = c \cdot k_1(\vec{x}_1, \vec{x}_2)$ with $c > 0$
- $k(\vec{x}_1, \vec{x}_2) = k_1(\vec{x}_1, \vec{x}_2) \cdot k_2(\vec{x}_1, \vec{x}_2)$
- $k(\vec{x}_1, \vec{x}_2) = k_3(\phi(\vec{x}_1), \phi(\vec{x}_2))$ with k_3 on $\mathcal{X}' \subseteq \mathbb{R}^M$ and $\phi: \mathcal{X} \rightarrow \mathcal{X}'$
- $k(\vec{x}_1, \vec{x}_2) = \vec{x}_1^T A \vec{x}_2$ with $A \in \mathbb{R}^{N \times N}$ symmetric and PSD

- Examples of kernels:

- Polynomial: $k(\vec{a}, \vec{b}) = (\vec{a}^T \vec{b})^P$ or $(\vec{a}^T \vec{b} + 1)^P$
- Gaussian: $k(\vec{a}, \vec{b}) = \exp\left(-\frac{\|\vec{a} - \vec{b}\|^2}{2\sigma^2}\right)$ hyperparameters
- Sigmoid: $k(\vec{a}, \vec{b}) = \tanh(\kappa \cdot \vec{a}^T \vec{b} - \delta)$ for $\kappa, \delta > 0$
is not PSD, but works well in practice

- Classifying new points w/ a kernelized SVM:

$$\underbrace{\alpha^* f_i(\vec{\theta}^*)}_{\text{complementary slackness condition}} = 0 \Leftrightarrow y_i \left(\sum_{j \mid \vec{x}_j \in S^c} \alpha_j y_j k(\vec{x}_i, \vec{x}_j) + b \right) = 1$$

set of support vectors with $\xi_i = 0$
(lie on the margin)

$$\Rightarrow b = y_i - \left(\sum_{j \mid \vec{x}_j \in S^c} \alpha_j y_j k(\vec{x}_i, \vec{x}_j) \right)$$

$$\Rightarrow h(\vec{x}) = \text{sign} \left(\sum_{j \mid \vec{x}_j \in S^c} \alpha_j y_j k(\vec{x}_j, \vec{x}) + b \right)$$

- Approaches for multiclass data:

- 1) One-vs-rest: $\begin{matrix} \text{SVM} \\ \downarrow \\ C \end{matrix}$ classifiers for C classes
→ Winner is the class where the dist from the hyperplane is max
- 2) One-vs-one: $\binom{C}{2}$ classifiers and evaluate all. The winner
is the class with the majority vote (weighted by the distance from the margin)
all possible pairings

ML 10 : Dimensionality reduction and matrix factorization

- Unsupervised learning: → Find hidden/latent structure in the data

↳ Can be viewed as compression (compress data point from a higher dim. to a lower dim. latent space, e.g. its cluster label)

→ Used for: * Visualization (understanding data)

* preprocessing (Dim. reduction)

* Unsupervised pretraining (Autoencoders)

* High-dim. data is problematic: 1) curse of dim.

2) costly computation

3) Highly correlated features cause problems w/ some algorithms

- Dim. Reduction via lin. transformations:

→ want to capture intrinsic value of data in any coordinate system

→ Use orthogonal basis transformation (and then potentially discard dims)

⇒ Transform data vector into new coordinate system

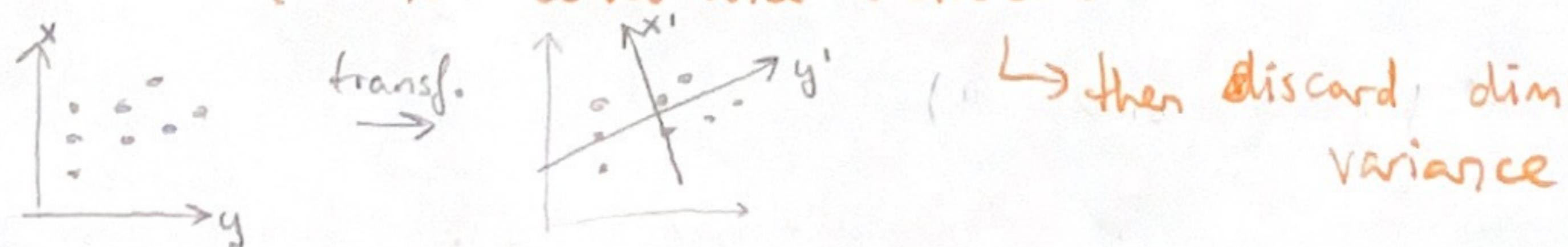
$$(x')^T = x^T \cdot F \quad \text{orthonormal transformation matrix } \in \mathbb{R}^{d \times k}$$

$$\Rightarrow X' = X \cdot F \quad \text{for all points } x_i$$

→ After the transf. discard dimensions w/ low variance.

- Principal Component Analysis: (PCA)

Goal { → Find a transf., such that features in the new coordinate system are lin. uncorrelated (→ the covariance between the new dim. is 0)

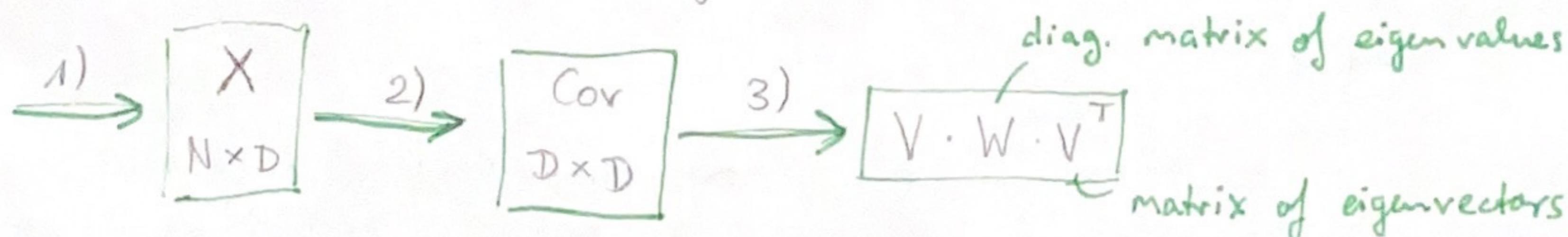


Approach:

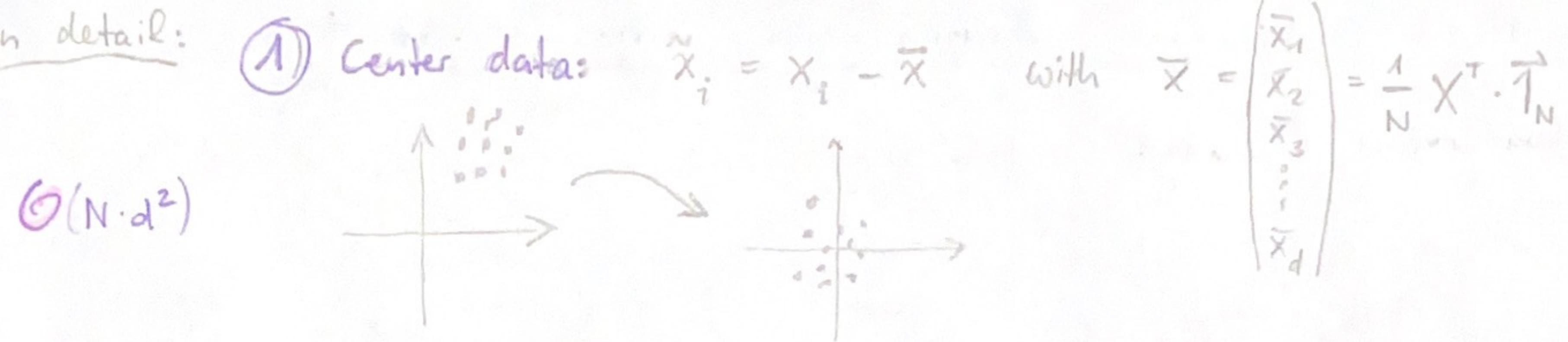
1) Center data (zero mean)

2) Comp. covariance matrix

3) Use the eigenvector decomposition to transf. the coord. sys.



In detail:



$\mathcal{O}(N \cdot d^2)$

$\mathcal{O}(d^3)$

② Compute covariance matrix:

Variance in the j th dim. of X : $\text{Var}(X_j) = \frac{1}{N} \sum_{i=1}^N (x_{ij} - \bar{x}_j)^2 = \frac{1}{N} \cdot X_j^T X_j - \bar{x}_j^2$

Covariance between

dim. j_1 and j_2 : $\text{Cov}(X_{j_1}, X_{j_2}) = \frac{1}{N} \sum_{i=1}^N (x_{ij_1} - \bar{x}_{j_1})(x_{ij_2} - \bar{x}_{j_2}) = \frac{1}{N} X_{j_1}^T X_{j_2} - \bar{x}_{j_1} \bar{x}_{j_2}$

"sigma" $\Rightarrow \Sigma_X = \begin{bmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_d) \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & & \vdots \\ \vdots & & \ddots & \vdots \\ \text{Cov}(X_d, X_1) & \dots & \dots & \text{Var}(X_d) \end{bmatrix}$

Since the mean of \tilde{X} is zero: $\Sigma_{\tilde{X}} = \frac{1}{N} \underbrace{X^T X}_{\text{sometimes omitted}} - \underbrace{\bar{x} \bar{x}^T}_{=0}$

$\mathcal{O}(N \cdot d \cdot k)$

③ Transf. the coord. sys. such that the covariances between the new axes are zero:

$$\text{Cov}' = \begin{bmatrix} \text{Var}(1)' & 0 & \dots & 0 \\ 0 & \text{Var}(2)' & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \text{Var}(d)' \end{bmatrix}$$

\Rightarrow Eigen decomposition:

$$\Sigma_{\tilde{X}} = \Gamma \cdot \Lambda \cdot \Gamma^T$$

\Rightarrow New coordinate sys. is defined by eigenvectors

"gamma"; orthonormal matrix w/ eigenvectors (normalized) of $\Sigma_{\tilde{X}}$ as columns

"lambda"; diagonal matrix of eigenvalues

$$\text{Covariance after transf: } \Sigma_{\tilde{X}'} = \Gamma^T \Sigma_{\tilde{X}} \Gamma = \Gamma^T \Gamma \Lambda \Gamma^T \Gamma = \Lambda$$

$$\Gamma = \begin{bmatrix} 1 & 1 & 1 & 1 \end{bmatrix}, \quad \Lambda = \begin{bmatrix} 10 & & & \\ & 7 & & \\ & & 3 & \\ & & & 0.1 \end{bmatrix}$$

Γ truncated

\Rightarrow Reduce dimensionality by keeping only the columns in Γ corresponding to the k largest eigenvalues (because the eigenvalue says how much the vector is scaled!) has the desired form

$$\Rightarrow Y_{\text{reduced}} = \tilde{X} \cdot \Gamma_{\text{reduced}}$$

Eigenvectors = Principal components

k is often 90% of the energy

$$\rightarrow \sum_{i=1}^k \lambda_i \geq 0.9 \sum_{i=1}^d \lambda_i$$

\rightarrow There are alternative views of PCA that treat it e.g. as an optimization problem

- Singular Value Decomposition: (SVD) → Finds best low-rank approximation to given matrix (18)

* Data often lies in a low-dim manifold embedded in a higher-dim space
 ↪ measure this dim. through the rank of the data matrix

* The data is noisy

→ Find best low-rank approximation

$$\text{Frobenius-Norm} \quad \|A - B\|_F^2 = \sum_{i=1}^N \sum_{j=1}^D (a_{ij} - b_{ij})^2$$

Euclidean dist. of flattened matrices

= Constraint optimization problem

$$\min_B \|A - B\|_F^2$$

constraint: $\text{rank}(B) = k$

→ Definition:

$$\min(\mathcal{O}(nd^2), \mathcal{O}(n^2d))$$

$n \times d$

$$A = U \cdot \Sigma \cdot V^T \quad \forall A \in \mathbb{R}^{n \times d}$$

U and V are "concepts"
 a.k.a. latent dimensions
 a.k.a. latent factors

$$\begin{aligned} & U \in \mathbb{R}^{n \times r} \quad \text{left singular vectors} \\ & V \in \mathbb{R}^{d \times r} \quad \text{right singular vectors} \\ & \Sigma \in \mathbb{R}^{r \times r} \quad \text{singular values} \quad \left. \begin{array}{l} \text{diagonal w/ positive values sorted} \\ \text{in decreasing order} \end{array} \right\} \text{column orthonormal} \end{aligned}$$

$$= \sum_{i=1}^r \underbrace{\sigma_i \cdot \underbrace{u_i}_{\substack{n \times 1 \\ \text{rank } 1}} \circ \underbrace{v_i^T}_{1 \times d}}_{\substack{n \times d \\ \text{rank } r}}$$

⇒ Represent the data as a linear combination of these latent concepts!

* The "strength" values σ in Σ weight the respective concept

* Interpretation: V^T : rows are the new coordinate axes (basis vectors)

σ : the "spread"/variance in the direction of a coordinate axis

$U \cdot \Sigma$: The coordinates of the points in the new coordinate system

→ Set smallest singular values σ to zero and ignore the corresponding dimensions

$$\begin{array}{c} \text{rank } k \\ \boxed{B} \\ \text{approx. of } A \end{array} = \boxed{U} \cdot \boxed{\Sigma} \cdot \boxed{V^T}$$

The reduced data is
 $P = U \Sigma = A \cdot V$

- Comparison:

→ Given data X (centered) → $\Sigma = U \Sigma V^T$ and new data is $XV = U \Sigma V^T$ (or truncated V)

SVD: $X = U \Sigma V^T$ and new data is $XV = U \Sigma V^T$ (or truncated V)

PCA: computes eigen decomposition of cov. matrix

$$\rightarrow \Sigma = X^T X$$

$$\rightarrow \text{Eigen decomposition: } X^T X = \Gamma \Lambda \Gamma^T$$

Projected data is $X \cdot \Gamma$ (or truncated Γ)

→ PCA and SVD are equivalent:

$$X^T X = (U \Sigma V^T)^T U \Sigma V^T = V \Sigma^T U^T U \Sigma V^T = V \Sigma^2 V^T$$

new
 $\Sigma^2 \stackrel{\Delta}{=} \Gamma \Lambda$

$$\Rightarrow V = \Gamma$$

$$\Sigma^2 = \Lambda$$

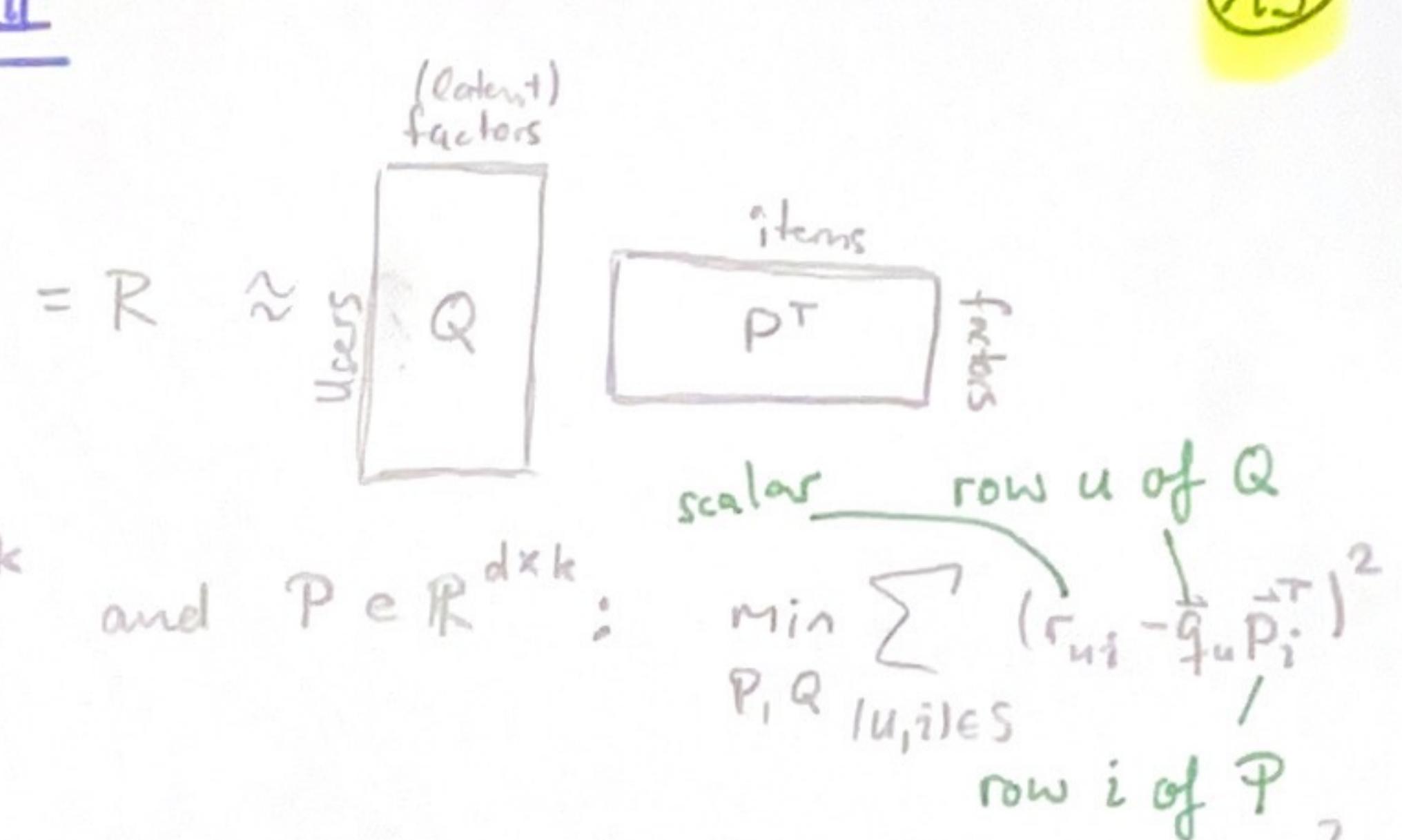
ML 11: Dim. Reduction & Matrix Factorization II

(19)

- Latent Factor Models : Given :

Netflix
recommendation
challenge

items								
1	3	1	2	3	4	5	6	7
6	3	1	2	3	4	5	6	7
2	7	6	3	4	5	1	2	3
7	8	2	1	3	4	5	6	7
8	9	3	2	1	4	5	6	7



Differences to PCA & SVD $\left\{ \begin{array}{l} \rightarrow \text{Only sum over existing entries : } S = \{(u, i) | r_{ui} \neq \text{missing}\} \\ \rightarrow \text{Don't require columns of } P \text{ and } Q \text{ to be orthogonal \& unit length} \end{array} \right.$

- Alternating optimization : (a.k.a. block coordinate optimization)

Pseudo Code

1. Init $P^{(0)}, Q^{(0)}$, $t=0$
2. $P^{(t+1)} = \underset{P}{\operatorname{argmin}} f(P, Q^{(t)})$
3. $Q^{(t+1)} = \underset{Q}{\operatorname{argmin}} f(P^{(t)}, Q)$
4. $t = t+1$
5. go to 2 until convergence

alternately keep one variable fixed and solve for the other



- Initialization:
 - With mean / random / ... values
 - Use result from SVD where empty entries are replaced with 0

$$Q = U\Sigma \quad \text{and} \quad P^T = V^T$$

• Solve: $P^{(t+1)} = \underset{P}{\operatorname{argmin}} f(P, Q^{(t)}) = \underset{P}{\operatorname{argmin}} \sum_{(u, i) \in S} (r_{ui} - \vec{q}_u \vec{p}_i^T)^2$

Since Q is fixed,
we can solve for
each \vec{p}_i independently

$$= \underset{i=1 \dots d}{\operatorname{argmin}} \sum_{u \in S_{*,i}} (r_{ui} - \vec{q}_u \vec{p}_i^T)^2$$

equivalently for $\vec{q}_u^{(t+1)}$

ordinary least squares regression problem $S_{*,i} = \{u | (u, i) \in S\}$

$$\min_w \sum_{j=1}^n (y_j - \vec{w}^T \vec{x}_j)^2 \rightarrow \vec{w}^* = \left(\frac{1}{n} \sum_{j=1}^n \vec{x}_j \vec{x}_j^T \right)^{-1} \cdot \frac{1}{n} \sum_{j=1}^n \vec{x}_j y_j$$

$$\Rightarrow \vec{p}_i^T = \left(\frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \vec{q}_u^T \vec{q}_u \right)^{-1} \cdot \frac{1}{|S_{*,i}|} \sum_{u \in S_{*,i}} \vec{q}_u^T r_{ui}$$

equivalently for \vec{q}_u

* Advantages & Drawbacks:

- + Solution to difficult optimization problem
- + simple to implement & efficient on sparse data
- solution is only approximation
- Solution is highly dependant on initial solution

- With SGD : objective : $\mathcal{L} = \sum_{(u,i) \in S} (r_{ui} - \vec{q}_u \vec{p}_i^\top)$
- 1) Pick random user u and item i with rating r_{ui}
 - 2) Compute gradients w.r.t. parameters = $\frac{\partial \mathcal{L}}{\partial \vec{q}_u}$ and $\frac{\partial \mathcal{L}}{\partial \vec{p}_i}$
 - 3) Update parameters : $\vec{q}_u \leftarrow \vec{q}_u - \eta \frac{\partial \mathcal{L}}{\partial \vec{q}_u}$, $\vec{p}_i \leftarrow \vec{p}_i - \eta \frac{\partial \mathcal{L}}{\partial \vec{p}_i}$
- learning rate
- Rating prediction :
- Estimate missing rating for user u and item i : $\hat{r}_{ui} = \vec{q}_u \cdot \vec{p}_i^\top$
- row u of Q row i of P
- User & item biases:
- $$\min_{P, Q} \sum_{\substack{(u, i) \in S \\ b_u, b_i, b}} (r_{ui} - (\vec{q}_u \vec{p}_i^\top + b_u + b_i + b))^2$$
- user bias item bias overall bias
(e.g. this user gives overly neg. reviews)
- Challenge: Overfitting:
- # of regression parameters = number of latent factors = k
 - # of data points = cardinality of $S_{*,*} / S_{u,*}$
 - \Rightarrow Use Regularization
- overfit if large k and less data
- $$\min_{P, Q} \sum_{(u, i) \in S} (r_{ui} - \vec{q}_u \vec{p}_i^\top)^2 + \left[\lambda_1 \sum_u \|\vec{q}_u\|^2 + \lambda_2 \sum_i \|\vec{p}_i\|^2 \right]$$
- optim. { → Closed form solution from ridge regression
→ SGD
- L2 vs L1-Regularization:
- L2: tries to "shrink" all values equally.
 - large values are highly penalized
 - it is unlikely that values are exactly 0
 - L1: enforces sparsity
- * We want sparsity
- Less data to store
 - More intuitive, that sparse input data comes from sparse signal
- Further Matrix Factorization stuff:
- Used for dim. reduction + Visualization (compute and plot 2-D or 3D (data analysis) latent factors)
 - * Often the data contains only non-negative values
 - SVD may lead to factor that have neg. values
 - difficult to interpret + this may predict neg. values
 - \Rightarrow Use Non-neg. matrix factorization

- Neighbor graph methods:

- Disadvantage of Matrix factorization methods (PCA, SVD): preserve global structure, but may loose local structure

- Neighbor graph methods try to preserve the neighborhood of each point (local structure)

- Procedure**
1. Construct neighborhood graph of high-dimensional data
 2. (Randomly) initialize points in a low-dim. space
 3. Optimize the coordinates of points in low-dim. space s.t. similarities align.
(between the random points in low-dim data and the orig. points in high dim. data)
-

"adjacency matrix of the graph"
"pairwise similarity matrix"
"pairwise distance matrix"

- t-SNE : t-distributed stochastic neighbor embedding

* High-dim similarities for input \vec{x}_i : $P_{j|i} =$

the probability of choosing point j when choosing according to the similarity to point i
only (very close)

→ by setting σ_i we get effective neighbors, e.g.

↳ set σ_i for each data point s.t. it has a fixed # of neighbors
for a large σ_i

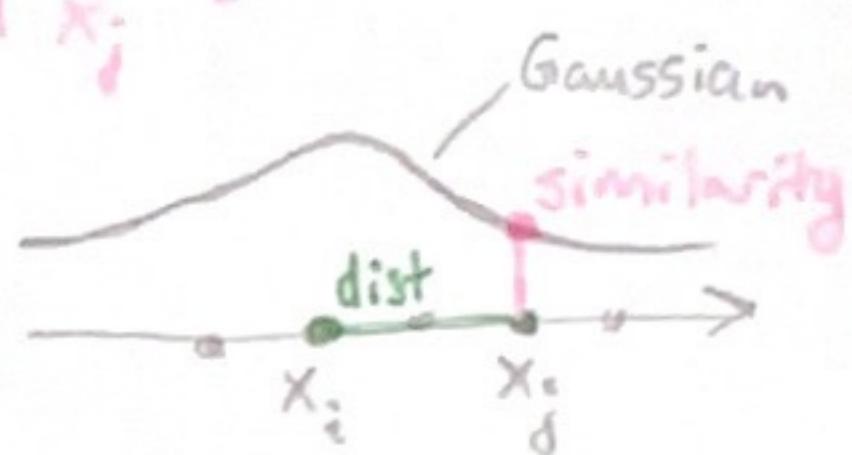
$$P_{ii} = 0, P_{ij} = \frac{P_{ilj} + P_{jli}}{2n} \quad (\text{fixed perplexity})$$

$$\frac{\exp\left(-\frac{\|\vec{x}_i - \vec{x}_j\|}{2\sigma_i^2}\right)}{\sum_{k \neq i} \exp\left(-\frac{\|\vec{x}_i - \vec{x}_k\|}{2\sigma_i^2}\right)}$$

dist. between \vec{x}_i and \vec{x}_j

convert to distribution

similarity between \vec{x}_i and \vec{x}_j



example for 1-D data,
(if the dist is smaller, the similarity will be larger)

→ we could use other ways to measure the similarity

* Low-dim similarities for parameters \vec{y}_i : $q_{ij} = \frac{(1 + \|\vec{y}_i - \vec{y}_j\|^2)^{-1}}{\sum_k \sum_{l \neq k} (1 + \|\vec{y}_i - \vec{y}_l\|^2)^{-1}}, q_{ii} = 0$
t-distribution

⇒ Make probabilities P_{ij} Match q_{ij} by changing \vec{y}_i :

$$\min_{\vec{y}_i} KL(P||Q) = \sum_i \sum_{i \neq j} P_{ij} \log\left(\frac{P_{ij}}{q_{ij}}\right)$$

→ KL-Divergence:

- $KL(P||Q) \geq 0 \quad \forall P, Q$ and $KL(P||Q) = 0 \iff P = Q$
- $KL(P||Q) \neq KL(Q||P)$

→ t-SNE is SotA for visualizing high-dim. data

↳ be careful when interpreting cluster sizes and distances

→ sensitive to hyperparams

→ Not good for more than 3 dimensions (→ not used for dim. reduction)

- Autoencoders (Non-lin. dim. reduction):

* PCA / SVD can only capture lin. structure : $\mathbf{Y} = \tilde{\mathbf{X}} \cdot \mathbf{P}$

→ But data may lie in a non-lin. low-dim. manifold

→ Autoencoder finds low-dim. representation by learning to reconstruct the input
 $f(\vec{x}, w) := \hat{x} \stackrel{!}{=} \vec{x}$

min. reconstruction error between \hat{x} and \vec{x} $\left\{ \min_w \frac{1}{N} \sum_{i=1}^N \|f(\vec{x}_i, w) - \hat{x}_i\|^2 \right.$

→ Find latent representation $\vec{z} \in \mathbb{R}^L$ which is a compact representation of $\vec{x} \in \mathbb{R}^D$

encoder: $f_{\text{enc}}(\vec{x}) = \vec{z}$

decoder: $f_{\text{dec}}(\vec{z}) \approx \vec{x}$

* Linear Autoencoders are PCA

$$\begin{aligned} f_{\text{enc}}(\vec{x}, w_1) &= \vec{x} w_1 \\ f_{\text{dec}}(\vec{x}, w_2) &= \vec{x} w_2 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} f_{\text{dec}}(f_{\text{enc}}(\vec{x})) = \vec{x} w_1 w_2 = \vec{x} w$$

$$\Rightarrow \min_w \frac{1}{N} \sum_{i=1}^N \|\vec{x}_i w - \vec{x}_i\|^2 = w^* \text{ s.t. the rank of } w^* \text{ is } L$$

↓
optimal solution through PCA : $w^* = P$

Summary:

Dim. Red. ~~Matrix factorization~~

- require less storage
- more efficient data processing
- PCA & SVD give optimal low-rank approximation

Matrix fact. methods:

- allows to handle missing values

Autoencoders:

- Works on a lot of data (images, ...)

ML 12: Clustering

→ group objects \vec{x}_i into K clusters based on their similarity

* is unsupervised learning

* Used for gene expression analysis, data compression, visualization, ...

→ Given data $X = \{\vec{x}_1, \dots, \vec{x}_n\}$ find

for each object \vec{x}_i find the assignment to a cluster $z_i \in \{1, \dots, K\}$ s.t.:

- objects within clusters are similar to each other

- objects between clusters are dissimilar to each other

- K-means algorithm:

→ distance-based: → maximize similarity or minimize dissimilarity (distance)

* Manhattan-distance: $d(\vec{x}_i, \vec{x}_j) = \sum_d |\vec{x}_{id} - \vec{x}_{jd}| = \|\vec{x}_i - \vec{x}_j\|_1$

* Euclidean distance: $d(\vec{x}_i, \vec{x}_j) = \sqrt{\sum_d (\vec{x}_{id} - \vec{x}_{jd})^2} = \|\vec{x}_i - \vec{x}_j\|_2$

* Mahalanobis distance: $d(\vec{x}_i, \vec{x}_j) = \sqrt{(\vec{x}_i - \vec{x}_j)^T \Sigma^{-1} (\vec{x}_i - \vec{x}_j)}$ if $\Sigma = I$ this is euclidean dist.

Model:

- Each cluster is defined by its centroid $\vec{\mu}_k \in \mathbb{R}^D$

- Cluster indicator $\vec{z}_i \in \{0, 1\}^K$ with $z_{ik} = 1 \Leftrightarrow \vec{x}_i$ belongs to cluster i
one-hot encoding

- K-means-objective (distortion measure)

$$J(X, z, \mu) = \sum_{i=1}^N \sum_{k=1}^K z_{ik} \|\vec{x}_i - \vec{\mu}_k\|_2^2 \stackrel{\text{def}}{=} \text{dist. of every point to its cluster's centroid}$$

→ Goal: find best centroids and cluster assignments

$$z^*, \mu^* = \underset{z, \mu}{\operatorname{argmin}} J(X, z, \mu)$$

Solution: - Solving the joint optim. problem is difficult → solve two sub problems

→ Lloyd's algorithm:

1) Initialize the centroids $\mu = \{\mu_1, \dots, \mu_K\}$

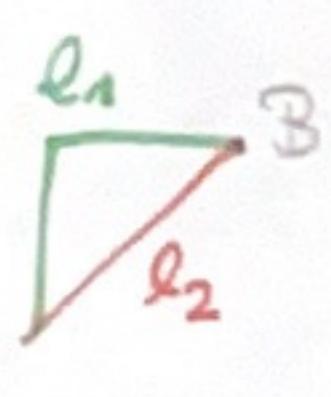
2) Update cluster indicators (solve $\min_z J(X, z, \mu)$)

$$z_{ik} = \begin{cases} 1 & \text{if } k = \arg \min_j \|\vec{x}_i - \vec{\mu}_j\|^2 \\ 0 & \text{else} \end{cases} \quad \left. \begin{array}{l} \text{assign every point} \\ \text{to its nearest cluster} \\ \text{centroid} \end{array} \right\}$$

3) Update centroids (solve $\min_\mu J(X, z, \mu)$)

$$\vec{\mu}_k = \frac{1}{N_k} \sum_{i=1}^N z_{ik} \vec{x}_i \quad \text{where} \quad N_k = \sum_{i=1}^N z_{ik} \quad \left. \begin{array}{l} \text{set centroids to mean} \\ \text{of all points in the} \\ \text{cluster} \end{array} \right\}$$

4) Return to 2) until convergence

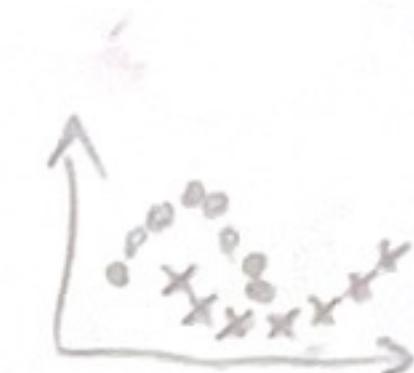


→ K-means++ algorithm: (Used to initialise good centroids)

- 1) Choose the first centroid $\vec{\mu}_1$ uniformly at random among the data points
- 2) For each point \vec{x}_i compute $D_i^2 = \|\vec{x}_i - \vec{\mu}_1\|_2^2$
- 3) Sample the next centroid $\vec{\mu}_k$ from $\{\vec{x}_i\}$ w/ probability proportional to D_i^2
- 4) Recompute $D_i^2 = \min \underbrace{\{\|\vec{x}_i - \vec{\mu}_1\|_2^2, \dots, \|\vec{x}_i - \vec{\mu}_k\|_2^2\}}_{\text{dist for point } \vec{x}_i \text{ to every centroid}}$
- 5) Return to 3) until all K initial centroids are chosen

→ Advantages & Disadvantages:

- | | |
|------------------|--|
| model issues | + cheap to compute w/ Lloyd's alg. : $O(nK)$ |
| } | - Quality depends on distance func a lot |
| | - Can't detect clusters w/ overlapping hulls |
| algorithm issues | - Sensitivity to outliers |
| } | - Extreme Sensitivity to initialization |



- Gaussian Mixture Model: (GMM)

general prob. view of clustering

Probabilistic Model: $p(x|\theta) = \sum_z p(x,z|\theta) = \sum_z p(x|z,\theta) p(z|\theta)$ since z are never observed → "latent variables"

and optimize w.r.t. θ

$$\theta^* = \arg \max_{\theta} p(x|\theta)$$

Gaussian mixture model: $p(\vec{x}, \vec{z}|\theta) = p(\vec{x}|\vec{z}, \theta) \cdot p(\vec{z}|\theta)$

with $p(\vec{z}|\theta) = \text{Cat}(\pi)$ Cluster prior is categorical

$p(\vec{x}|\vec{z}_k=1, \theta) = \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$ Each cluster has its own multivariate normal distribution

$$\Rightarrow \theta = \{\pi, \mu, \Sigma\}$$

→ Generative process of GMMs: 1) Draw one-hot cluster indicator $\vec{z} \sim \text{Cat}(\pi)$

2) Draw sample $\vec{x} \sim \mathcal{N}(\mu_k, \Sigma_k)$ if $z_k = 1$

marginalization over z

$$\# \text{ GMM Likelihood: } p(\vec{x}|\vec{\pi}, \mu, \Sigma) = \sum_{k=1}^K p(z_k=1|\vec{\pi}) p(\vec{x}|z_k=1, \vec{\mu}_k, \Sigma_k)$$

$$= \sum_{k=1}^K \pi_k \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$$

i.i.d. data

$$\Rightarrow p(X|\vec{\pi}, \mu, \Sigma) = \sum_{i=1}^N p(\vec{x}_i|\vec{\pi}, \mu, \Sigma)$$

(b)

"The likelihood is a mixture of k different Gaussians"

→ Inference in GMM:

Given parameters $\{\vec{\pi}, \vec{\mu}, \Sigma\}$, we want to assign points to clusters

⇒ compute the posterior distr. of cluster indicators $p(z|x, \pi, \mu, \Sigma)$

Bayes' rule

$$\Rightarrow p(z_{ik}=1 | \vec{x}_i, \pi, \mu, \Sigma) = \frac{p(z_{ik}=1 | \vec{\pi}) p(\vec{x}_i | z_{ik}=1, \vec{\mu}_k, \Sigma_k)}{p(\vec{x}_i | \vec{\pi}, \mu, \Sigma)}$$

$$p(z|\vec{x}, \theta) = \frac{p(x|z, \theta) p(z|\theta)}{p(x|\theta)}$$

$$= \frac{\pi_k N(\vec{x}_i | \vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(\vec{x}_i | \vec{\mu}_j, \Sigma_j)}$$

$\hat{\gamma}(z_{ik}) = \text{the responsibility of component } k \text{ for observation } i$

→ Learning in GMM:

"How likely does the datapoint belong to the cluster?"

We want to determine optimal parameters (find values that max. the log-likelihood)

$$\pi^*, \mu^*, \Sigma^* = \underset{\pi, \mu, \Sigma}{\operatorname{argmax}} \sum_{i=1}^N \log \left(\sum_{k=1}^K \pi_k N(\vec{x}_i | \vec{\mu}_k, \Sigma_k) \right)$$

The sum over the log makes it difficult to solve this!

Instead of optimizing $\log p(X|\theta)$ (which we can't) we optimize a proxy objective:

$$\mathcal{L}(\gamma_t|\theta) := \underbrace{E_{z \sim \gamma_t(z)} [\log p(x, z | \pi, \mu, \Sigma)]}_{\text{is a lower bound of the log-likelihood (proof in lecture)}} = \underbrace{2 \sum_{i=1}^N \sum_{k=1}^K \gamma_t(z_{ik}) \log p(x_i, z_{ik} | \theta)}_{\text{sum over log which we can solve}}$$

→ Expectation Maximization (EM) algorithm: (is alternating optim.)

Appl. of the EM-Alg. in general:

Given a problem:

$$\max_{\theta} \log p(x|\theta) = \log \sum_z p(x, z|\theta)$$

E-step:

Eval the posterior

$$\gamma_t = p(z|x, \theta^{(t)})$$

1) Initialize model parameters: $\{\pi^{(0)}, \mu_1^{(0)}, \dots, \mu_K^{(0)}, \Sigma_1^{(0)}, \dots, \Sigma_K^{(0)}\}$

2) E Step

Evaluate the responsibilities

$$\gamma_t(z_{ik}) = p(z_{ik}=1 | \vec{x}_i, \pi^{(t)}, \mu^{(t)}, \Sigma^{(t)})$$

same as inference

3) M Step (max. the proxy objective w.r.t. the parameters)

$$\mu_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_{ik}) \vec{x}_i$$

$$\Sigma_k^{(t+1)} = \frac{1}{N_k} \sum_{i=1}^N \gamma_t(z_{ik}) (\vec{x}_i - \vec{\mu}_k^{(t+1)}) (\vec{x}_i - \vec{\mu}_k^{(t+1)})^\top$$

$$\pi_k^{(t+1)} = \frac{N_k}{N}$$

with $N_k = \sum_{i=1}^N \gamma_t(z_{ik})$

M-step:

max. the expected joint log-likelihood $\log p(x, z|\theta^{(t+1)})$ w.r.t. θ under the current beliefs $\gamma_t(z)$

4) Return to 2) until the proxy objective converges

- Summary of EM:

* It is not possible to optimize the log-likelihood due to latent variables Z

$$\log(p(x|\theta)) = \log\left(\sum_z p(x|z,\theta) p(z|\theta)\right)$$

* It is "easy" to optimize the joint probability

$$\log(p(x,z|\theta)) = \log(p(z|\theta)) + \log(p(x|z,\theta))$$

⇒ Main idea of EM:

1) Use our current beliefs $p(z|x,\theta^{(t)})$ to "pretend" that we know Z

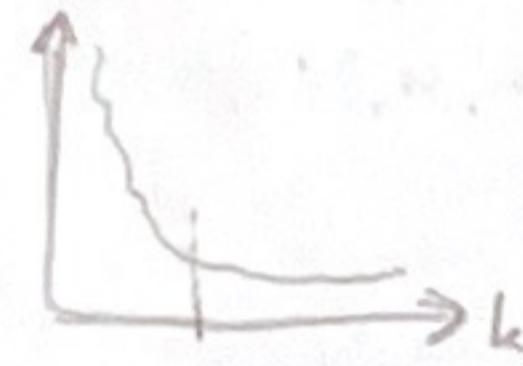
2) E step: update the responsibilities $\gamma_t(z) = p(z|x,\theta^{(t)})$

3) M step: update parameters $\theta^{(t+1)} = \arg \max_{\theta} E \left[\log p(x,z|\theta) \right]$

: take all previous $\gamma_t(z)$

- Choosing the number of clusters:

→ Heuristic methods: Elbow / knee heuristic: Plot within-cluster sum of squared distances for varying K



→ optimal K at knee

Gap statistic: Compare within-cluster variation to uniform data.

Silhouette: Per point difference of ^{the} within-cluster mean distance ~~to all~~ ^{for points in this cluster} to points in closest cluster. Maximize this for all points.

→ Probabilistic methods:

* Need a generative model that defines the data likelihood $\hat{L} = p(x|\hat{z}, \hat{\theta})$

Bayesian information criterion: (BIC)

optimal parameters

$$BIC = \underbrace{M \log \underbrace{n}_{\text{num free parameters}}}_{\text{num samples}} - 2 \log \hat{L}$$

Akaike information criterion: (AIC)

$$AIC = 2M - 2 \log \hat{L}$$

- Hierarchical clustering:

→ Agglomerative: Bottom up

→ Divisive: Top down

ML 13 - Privacy

- What do we have to consider beyond accuracy:

- * **Privacy** (models with large capacity can memorize data)
- * **Security** (fool classifier with adversarial input)
- * **Fairness**
- * **Explainability**
- * **Accountability** (can an outside editor ensure the system works as intended?)

- Differential privacy: (DP)

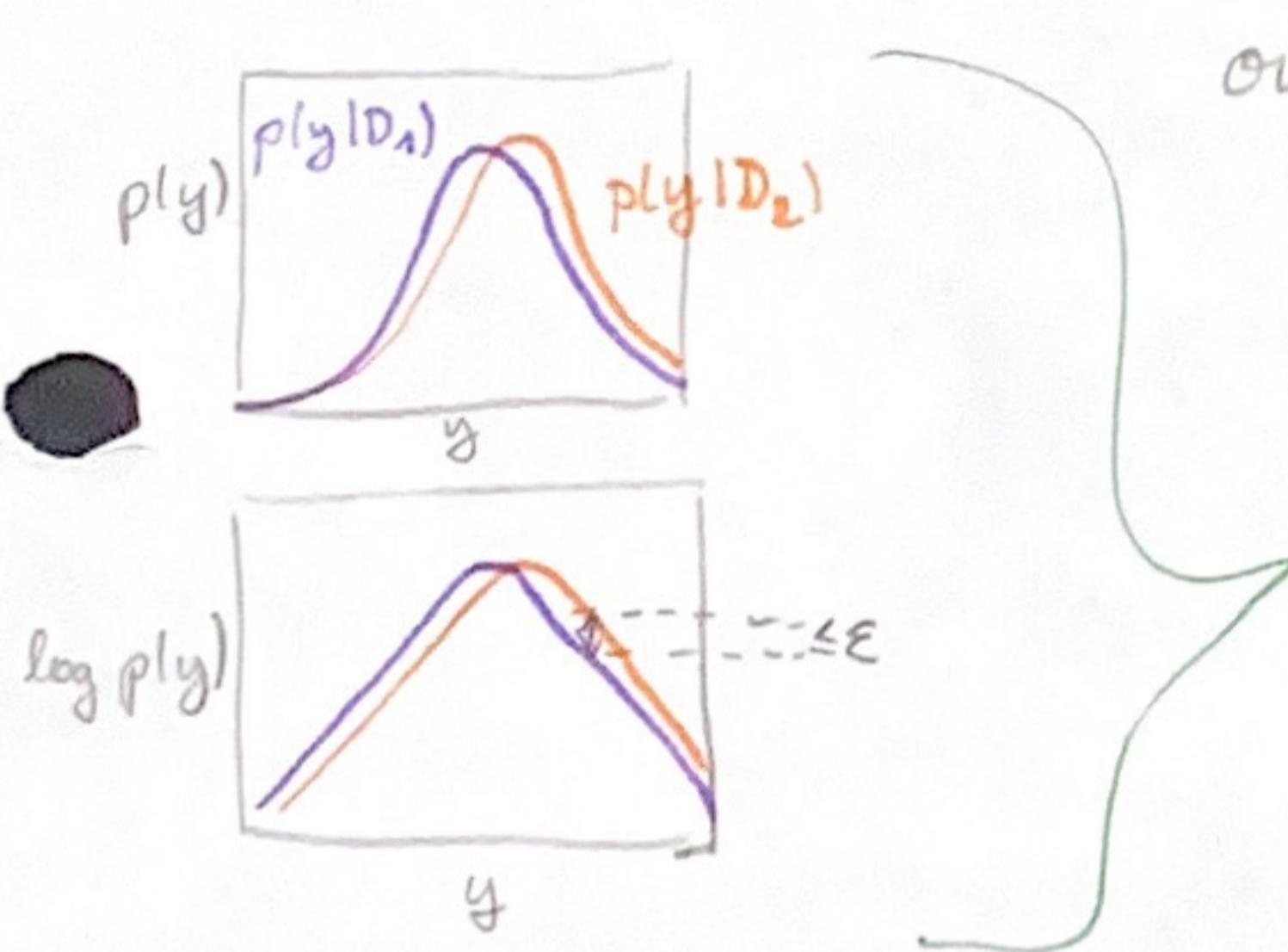
→ **Setup:** Input space X with symmetric neighbouring relation \simeq
Output space Y

Function of interest f and privacy parameter $\epsilon \geq 0$

→ **Randomized mechanism:**

→ computes the function f and adds noise

* $M_f: X \rightarrow Y$ is ϵ -differentially private if for all neighboring inputs $X \simeq X'$ and for all set of "a single instance is changed"



outputs $Y \subseteq Y$ we have:

$$e^{-\epsilon} \leq \frac{\Pr[M_f(X) \in Y]}{\Pr[M_f(X') \in Y]} \leq e^{\epsilon} \Rightarrow |\log(\cdot) - \log(\cdot)| \leq \epsilon$$

prob. of output is for changed input is very similar

the two distributions w/
changed data must be similar

distribution over noise
z for every single dim

- Laplace - Mechanism:

→ global sensitivity ~~function~~: of $f: X \rightarrow \mathbb{R}^d$:

$$\Delta_p = \sup_{X \simeq X'} \|f(X) - f(X')\|_p$$

captures what
is protected

measures by, the magnitude by which a single instance can change the output of the func in the worst case.

Output perturbation w/ Laplace noise:

- 1) A curator holds data $X = (x_1, \dots, x_n) \in \mathcal{X}$
- 2) The curator computes $f(X)$
- 3) They sample i.i.d. Laplace noise $z \sim \text{Lap}(0, \frac{\Delta_1}{\epsilon})^d$
- 4) They reveal the noisy data $f(X) + z$

DP for Machine Learning:

Goal: Train a model on dataset \mathcal{D} and release the optimal weights (from the private data) in public.

Setup:

- dataset \mathcal{D}

- parameters θ

- $f(\mathcal{D}) = \underset{\theta}{\operatorname{argmin}} L(\mathcal{D}, \theta)$ are the optimal weights

- $\mathcal{D} \approx \mathcal{D}'$ when they differ in one instance

→ Techniques:
"D neighbors D"

- * Perturb input (add noise to dataset)

- * Perturb weights (add noise to weights)

- * Perturb objective (Optimize $L(\mathcal{D}, \theta) + \theta^T z$)

- * Perturb gradients

→ Fundamental properties of DP.

- * Robustness to post-processing (operations after applying the DP-mechanism preserve DP)

- * Composition (Apply multiple DP-mechanisms)

- * Group privacy (If M is ϵ -DP w.r.t. $X \approx X'$, then applying M to data with t changed instances is $(t\epsilon)$ -DP)

Federated Learning:

→ learning a model w/o any centralized entity having access to all the data

- 1) Send weights from server to user
- 2) User computes loss + gradient locally
- 3) Send updated weights back (add randomization at user level)

- What causes bias?
 - * Tainted training data ("If there is already a bias in the data")
 - * Skewed sample (selection bias; initial predictions may influence future observations)
 - ↳ "What is in the data is as important as what is not in the data"
 - * Proxies: Some features are legally protected (e.g. race, etc.). But other features are highly correlated w/ these.
 - * Sample size disparity: Models tend to fit larger groups first and possibly trade off accuracy for smaller groups
 - * Limited features: Features for minority groups may be less informative

- Definitions of fairness:

- * **Group Fairness:** Treat all groups equally
- * **Individual Fairness:** Treat similar individuals similarly
- * **Counterfactual Fairness:** "What would the decision be if the individual (group) was changed s.t. it belongs to a different group?"

- How to achieve fairness:

Setup: $X \in \mathbb{R}^d \rightarrow$ features of indiv.
 $A \in \{a, b, \dots\} \rightarrow$ sensitive features
 $R = r(X, A) \in \{0, 1\} \rightarrow$ binary predictor
 $Y \in \{0, 1\} \rightarrow$ target

\Rightarrow assume X, Y, A are generated from underlying distribution

\Rightarrow Notation $P_a\{R\} = P\{R | A = a\}$

→ Fairness through unawareness:

- * **Use $R = r(X)$** instead of $R = r(X, A)$
 - + easy to use and implement
 - does not work due to correlations w/ other features

→ First Fairness Criterion: Independence:

- * R must be independent of $A \Rightarrow R \perp\!\!\!\perp A$
- * In binary classification: For all groups a, b we require $P_a\{R=1\} = P_b\{R=1\}$
 - Approx. version $\Rightarrow \frac{P_a\{R=1\}}{P_b\{R=1\}} \geq 1 - \epsilon$
 - same percentage of applicants receive loan

* Achieve Independence: 1) Post-processing \rightarrow decorrelate

2) Training-Time Constraint

(-) cannot use information A

\rightarrow Optimal predictor $R=Y$ is ruled out, since we need to cover each group equally

3) Pre-processing / representation learning

\hookrightarrow Map (X, A) to a low-dim. representation Z which is not correlated to A. Train on Z .

\rightarrow Second Fairness criterion: Separation:

* Require R and A to be conditionally independent given Y

$$\Rightarrow R \perp\!\!\!\perp A | Y$$

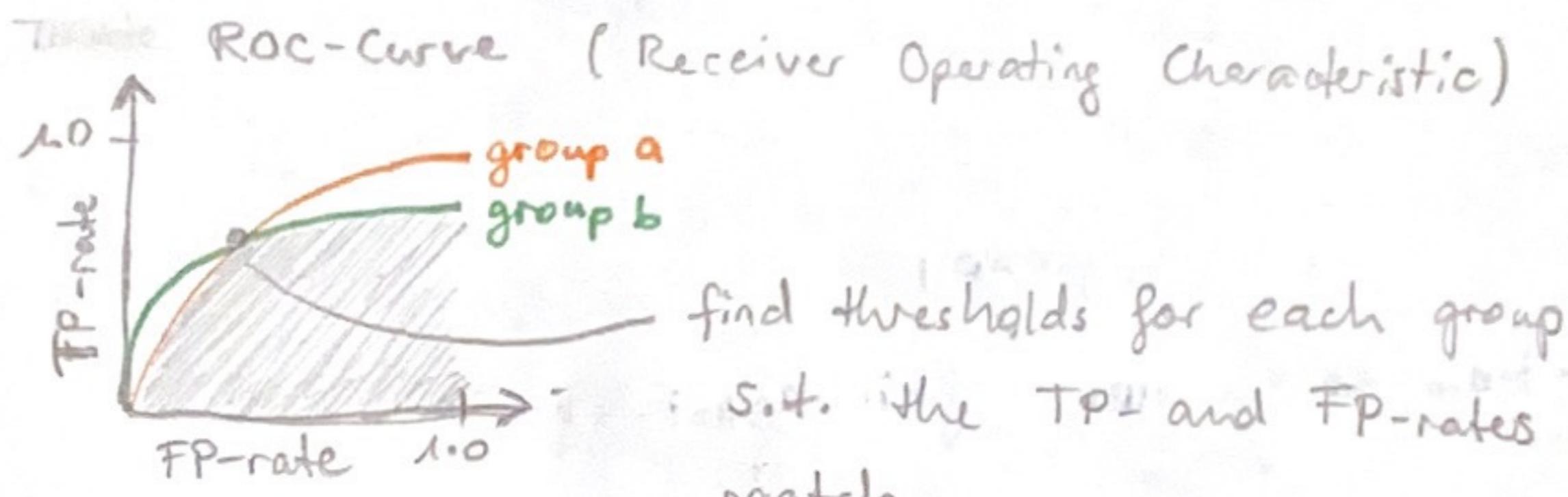
* In binary classification: (Enforce that the true positive and FP rates are equal for all groups)

$$P_a(R=1 | Y=1) = P_b(R=1 | Y=1) \quad \text{TP}$$

$$P_a(R=1 | Y=0) = P_b(R=1 | Y=0) \quad \text{FP}$$

* Commonly used relaxation: Only match the TP rate \rightarrow Equal opportunity

* Achieve Separation:



(+) Optimal predictor not ruled out

(-) does not help closing gap between two groups

\rightarrow Third Criterion: Sufficiency:

* Require X and A to be independent conditional on the target

$$\Rightarrow Y \perp\!\!\!\perp A | R$$

* In binary classification: (Enforce that the rate of the target is the same given same scores for different groups)

$$P_a(Y=1 | R=r) = P_b(Y=1 | R=r)$$

\Rightarrow Each group should be calibrated

$$\underbrace{P_a(Y=1 | R=r)}_{\text{the confidence of the classifier}} = r$$

the confidence of the classifier should match the probability of the outcome

(+) Equal chance of success ($Y=1$) given acceptance ($R=1$)

(+) To predict Y we do not need A when we have R

(-) may not help to close gaps

- Adversarial Examples / Robustness:

\hookrightarrow deliberate perturbations of the data to achieve some malicious goal
 ^{10.r.t.}
 \hookrightarrow worst-case noise