

# ML 1

Topics:

- ① Bayes Decision Theory
- ② Parameter Estimation
- ③ Fisher Discriminant
- ④ PCA
- ⑤ SVM
- ⑥ Learning Theory & kernels
- ⑦ model selection
- ⑧ NN 1
- ⑨ NN 2
- ⑩ Decision Trees & Random Forests
- ⑪ Boosting
- ⑫ Clustering
- ⑬ KRR
- ⑭ XAI
- ⑮ Applications

## Lecture 1: Bayesian Decision Theory

What is Machine Learning?

Designing algorithms (machines) that efficiently convert finite sets of data (observations) into models with predictive capability.

Why Machine Learning?

Autonomous Decision Making: Substitute/support human decision making for achieving gains in efficiency in practical applications.

Knowledge Discovery: Finding laws that explain empirical phenomena.

### Bayes Decision Theory

A theoretical framework for building decision models that:

-Yields optimal classifiers (under some favourable conditions).

Helps us to understand how:

-properties of the data distribution (e.g. Gaussian-distributed),

-prior probabilities of classes,

-the criterion to optimize (e.g. maximum classification accuracy),

affect the classifier qualitatively and quantitatively.

**Posterior:**  
What is the prob. of  
being a certain class  
after observing  $x$ ?

features  
 $x_1$  length     $x_2$  lightness

salmon  $w_1$   
sea bass  $w_2$

Notation:

- $\omega_1, \omega_2, \dots$  set of classes (e.g. salmon, sea bass, ...),
- $x \in \mathbb{R}^d$  vector of observations (e.g.  $x_1$  is the length and  $x_2$  is the lightness).

We are also given the probability laws:

- $P(\omega_j)$ : probability of being of class  $j$ , **prior**
- $p(x | \omega_j)$ : density function of measurements for each class,
- $p(x)$ : density function of measurements (marginalized),

but what we are truly interested in is:

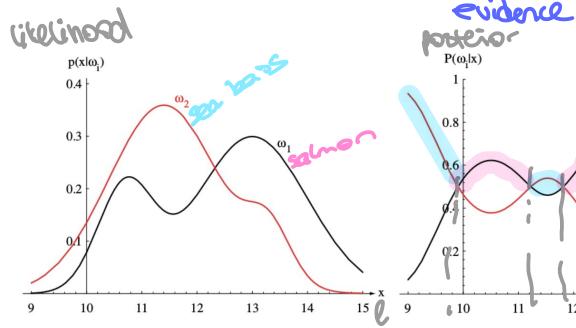
- $P(\omega_j | x)$ : probability of being of a certain class after observing  $x$ .

$p(\omega_j)$ : Salmon probability  
 $p(x | \omega_j)$ : size, length of salmon class  
 $p(x)$ : lengths of salmons  
 $p(\omega_j | x)$ : now I have some lengths, what is the prob. that these are in salmon class?  
**posterior**

Bayes Theorem:

$$P(\omega_j | x) = \frac{p(x | \omega_j) \cdot P(\omega_j)}{p(x)}$$

posterior  
likelihood prior  
evidence posterior



## Optimally Accurate Classifier

optimal decision function:  $\arg \max_j P(\omega_j | x)$

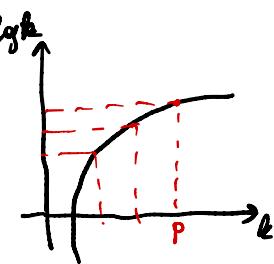
Alternative Formulas of the decision!

$$= \arg \max_j \frac{p(x | \omega_j) \cdot p(\omega_j)}{p(x)}$$

$$= \arg \max_j p(x | \omega_j) \cdot p(\omega_j)$$

$$= \arg \max_j \log [p(x | \omega_j) \cdot p(\omega_j)]$$

$$= \arg \max_j \log p(x | \omega_j) + \log p(\omega_j)$$



- Salmon: 11.5-ish (avg.) ...  $11.5 < l < 12.75 \rightarrow$  sea bass
- Sea bass: 13.5-ish (avg.)  $9.5 < l < 11 \rightarrow$  salmon  
 $11 < l < 12 \rightarrow$  sea bass  
 $l < 11 \rightarrow$  salmon

## Multivariate Normal Distributions

integrates to 1

$$p(x | \omega_j) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left( -\frac{1}{2} (x - \mu_j)^\top \Sigma_j^{-1} (x - \mu_j) \right)$$

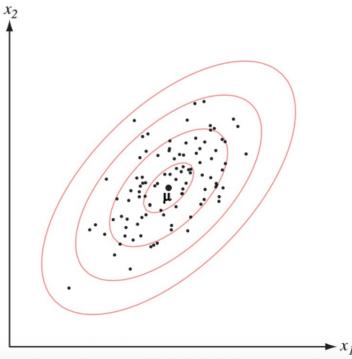


Image source: Duda et al. (2000)

- $\mu_j$  is the mean (center of the data distribution) location
- $\Sigma_j$  is the covariance (elongation of the data distribution and correlation between dimensions).

## Classifier for Gaussians ( $\Sigma_1 = \Sigma_2$ )

Recall: The optimal classifier is  $\arg \max_j [\log p(x | \omega_j) + \log P(\omega_j)]$ , and we have the data distributions:

$$p(x | \omega_j) = \frac{\log}{\sqrt{(2\pi)^d \det(\Sigma_j)}} \exp \left( -\frac{1}{2} (x - \mu_j)^\top \Sigma_j^{-1} (x - \mu_j) \right)$$

↑ does not depend on  $\omega_j$  ↗ because  $\Sigma_1 = \Sigma_2$ , then we have a single matrix

$$= \arg \max_j \left[ -\frac{1}{2} (x - \mu_j)^\top \Sigma_j^{-1} (x - \mu_j) + \log P(\omega_j) \right]$$

$$= \arg \max_j \left[ -\frac{1}{2} x^\top \Sigma_j^{-1} x + \underbrace{x^\top \Sigma_j^{-1} \mu_j}_{\text{const does not depend on } j} - \frac{1}{2} \mu_j^\top \Sigma_j^{-1} \mu_j + \log P(\omega_j) \right]$$

log of a product = sum of the log

log of exponential = identity

adding something constant to argmax (exp function), don't change anything

## Classifier for Gaussians ( $\Sigma_1 = \Sigma_2$ ) $\hookrightarrow$ db is linear

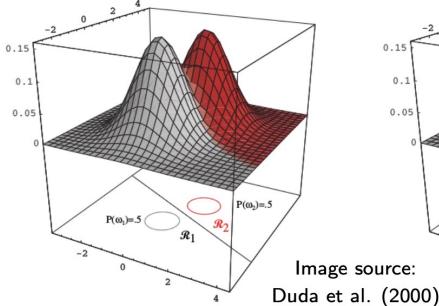
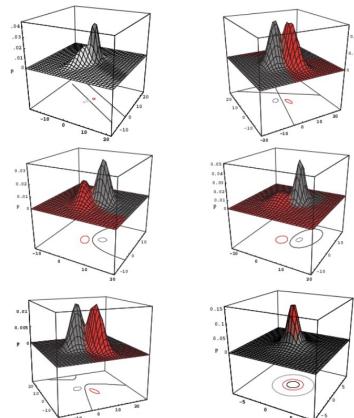


Image source:  
Duda et al. (2000)

- Decision boundary is linear and oriented by mean and covariance.
- Offset is controlled by class prior probabilities.

## Classifier for Gaussians ( $\Sigma_1 \neq \Sigma_2$ ) $\hookrightarrow$ then the db is quadratic



- When covariances  $\Sigma_1$  and  $\Sigma_2$  are not the same, the decision boundary is quadric instead of linear. Quadrics include circle, ellipse, parabola, hyperbola, and degenerate forms.

Image source: Duda et al. (2000)

# Classifying Binary Data

- Assume that our data is binary, i.e.  $x \in \{0, 1\}^d$ , with each dimension generated *independently* according to some Bernoulli distribution:

$$P(x_i = 0 | \omega_j) = 1 - q_{ij}$$

$$P(x_i = 1 | \omega_j) = q_{ij}$$

where  $q_{ij}$  are the parameters.

- The probability of the whole multivariate observation can be written as:

$$P(x | \omega_j) = \prod_{i=1}^d [q_{ij}x_i + (1 - q_{ij}) \cdot (1 - x_i)]$$

- Question: How to express the optimal decision boundary

$$\arg \max_j P(\omega_j | x)$$

## Minimum Cost Decisions

### Example: Buying a Car

- Suppose you would like to purchase a second-hand car. After observing the car (collecting a vector of measurements  $x$ ), you assess that it has a defect with probability

$$P(\text{defect} | x) = 0.1 \quad P(\text{no defect} | x) = 0.9$$

- Concretely, the decision you need to take is *not to classify* the whether the car has a defect or not, but whether to *buy* the car or not.
- For this, we need to evaluate the cost of each scenario, e.g.

$$\begin{aligned} \text{cost (buy | defect)} &= 100.0 \\ \text{cost (buy | no defect)} &= -20.0 \\ \text{cost (not buy | defect)} &= 0.0 \\ \text{cost (not buy | no defect)} &= 0.0 \end{aligned}$$

and take the action with the *lowest expected cost*.

## Classification Accuracy Special Case

$$\text{Recall: } \lambda(\alpha_k | x) = \sum_{j=1}^C \lambda(\alpha_k | \omega_j) P(\omega_j | x)$$

Show that the problem of maximum accuracy classification is a special instance of expected cost minimization with a particular set of actions  $(\alpha_k)_k$  and a particular cost function  $\lambda(\alpha_k | \omega_j)$ .

$$\text{decide classify as } \omega_k$$

$$\lambda(\text{decide } \omega_k | \omega_j) = -f_{jk}$$

$$\arg \min_k \sum_j -f_{jk} \cdot P(\omega_j | x)$$

$$= \arg \min_k \{-P(\omega_k | x)\} = \arg \max_k P(\omega_k | x)$$

$$\arg \min(-) = \arg \max$$

# Classifying Binary Data

**Recall:** The optimal classifier is  $\arg \max_j [\log p(x | \omega_j) + \log P(\omega_j)]$ , and we have the data distributions:

$$\begin{aligned} \text{left: } P(x | \omega_j) &= \prod_{i=1}^d [q_{ij}x_i + (1 - q_{ij}) \cdot (1 - x_i)] \\ &= \arg \max_j \left( \sum_i \log (q_{ij}x_i + (1 - q_{ij})(1 - x_i)) + \log P(\omega_j) \right) \\ &= \arg \max_j \left( \sum_i x_i \underbrace{\log q_{ij}}_{a_{ij}} + (1 - x_i) \underbrace{\log (1 - q_{ij})}_{b_{ij}} + \log P(\omega_j) \right) \\ &= \arg \max_j (x^T a_j + (1 - x)^T b_j + \log P(\omega_j)) \\ &= \arg \max_j \{x^T (a_j - b_j) + 1^T b_j + \log P(\omega_j)\} \end{aligned}$$

## Minimum Cost Decisions

### General problem formulation:

- Let  $(\alpha_k)_k$  be the set of actions. The expected cost " $\lambda$ " of taking a certain action is given by

$$\lambda(\alpha_k | x) = \sum_{j=1}^C \lambda(\alpha_k | \omega_j) P(\omega_j | x)$$

where  $\lambda(\alpha_k | \omega_j)$  is the cost of taking action  $\alpha_k$  given the class  $\omega_j$ .

- The optimal action to take is therefore:  $\arg \min_k \lambda(\alpha_k | x)$

$$\lambda(\text{buy} | x) = 100 \cdot 0.1 + (-20) \cdot 0.9 = -8$$

$$\lambda(\text{not buy} | x) = 0 \cdot 0.1 + 0 \cdot 0.9 = 0$$

$$\arg \min \{\lambda(\text{buy}): -8, \lambda(\text{not buy}): 0\} = \text{buy}$$

## Measuring Classification Error

- So far, we have studied what the decision boundary should be in order to predict optimally.

- However, in certain cases, it is also important to determine what is the *expected error* of the classifier (e.g. to determine whether the classifier is good enough for practical use).

- The expected error is the probability that the data is of a different class than the one predicted, e.g. for a binary classifier:

$$P(\text{Err} | x) = \begin{cases} P(\omega_1 | x) & \text{if "decide } \omega_2" \\ P(\omega_2 | x) & \text{if "decide } \omega_1" \end{cases}$$

- For the Bayes optimal classifier, this reduces to

$$P(\text{Err} | x) = \min\{P(\omega_1 | x), P(\omega_2 | x)\}$$

- The expected error of this maximally accurate classifier is computed as the integral of its error probability over the distribution  $p(x)$ .

$$P(\text{Err}) = \int_x P(\text{Err} | x) p(x) dx \\ = \int_x \min\{P(\omega_1 | x), P(\omega_2 | x)\} p(x) dx$$

This is also known as the *Bayes error rate*.

- Generally, this integral cannot be solved analytically, because of the min function. Error must instead be evaluated numerically/empirically, or it can also be **bounded** analytically.

→ If we have model knowledge (posterior, prior and likelihood) then Bayes will have the lowest error rate.

## Bounding the Error of the Classifier

**Another simple bound**

$$P(\text{Err}) = \int_x \min\{P(\omega_1 | x), P(\omega_2 | x)\} p(x) dx \\ = \int_x \min\{P(\omega_1 | x)p(x), P(\omega_2 | x)p(x)\} dx$$

Recall:  $P(\omega_j | x)p(x) = p(x | \omega_j)P(\omega_j)$

$$\begin{aligned} &= \int_x \min\{p(x | \omega_1)P(\omega_1), p(x | \omega_2)P(\omega_2)\} dx \\ &\leq \int_x \min\{\sum_j p(x | \omega_j)\} P(\omega_1), \sum_j p(x | \omega_j) P(\omega_2) dx \\ &= (\sum_j \int_x p(x | \omega_j) dx) \cdot \min\{P(\omega_1), P(\omega_2)\} \\ &= 2 \cdot \min\{P(\omega_1), P(\omega_2)\} \end{aligned}$$

↑ sum of 2 classes  
↑ integrates to 1  
↑ depends on the class priors  
↑ very likely  
↑ very unlikely

**Additional insight:** The optimal classifier improves its accuracy when one class prior probability is strongly dominant over the other class.

## Bounding the Error of the Classifier

**Very basic bound**

- Observe for binary classification that  $P(\omega_2 | x) = 1 - P(\omega_1 | x)$ .
- The error of an optimal binary classifier can be bounded as:

$$\begin{aligned} P(\text{Err}) &= \int_x \min\{P(\omega_1 | x), P(\omega_2 | x)\} p(x) dx \\ &= \int_x \min\{P(\omega_1 | x), 1 - P(\omega_1 | x)\} p(x) dx \\ &\leq \int_x 0.5 p(x) dx \quad \text{ε 0.5} \\ &= 0.5 \end{aligned}$$

i.e. the classifier predicts the correct class at least 50% of the time.

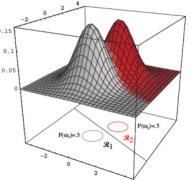
- Note that, unlike an empirical evaluation, this result is general and independent on the data distributions.

## Recap: Bayes Optimal Classifier

**Recap:**

- Assume our data is generated for each class  $\omega_j$  according to the multivariate Gaussian distribution  $p(\mathbf{x} | \omega_j) = \mathcal{N}(\boldsymbol{\mu}_j, \Sigma)$  and with class priors  $P(\omega_j)$ . The Bayes optimal classifier is derived as

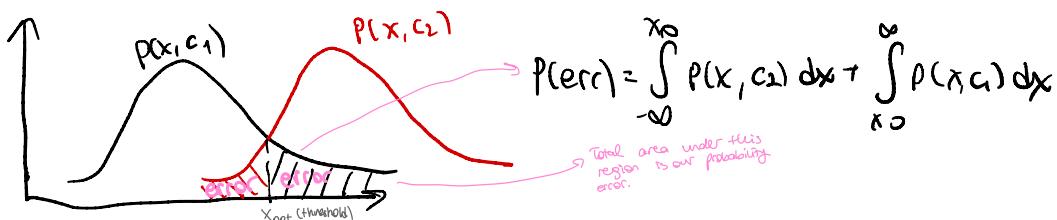
$$\begin{aligned} &\arg \max_j \{P(\omega_j | \mathbf{x})\} \\ &= \arg \max_j \{\log p(\mathbf{x} | \omega_j) + \log P(\omega_j)\} \\ &= \arg \max_j \{\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_j^\top - \frac{1}{2} \boldsymbol{\mu}_j^\top \Sigma^{-1} \boldsymbol{\mu}_j + \log P(\omega_j)\} \end{aligned}$$



- Given our generative assumptions, there is no better classifier than the one above.
- However, in practice, we don't know these distributions and only have the data.



- best classifier if I know the data properties, when we know everything about data generation
- conditions: data Gaussian distributed & both classes same cov ⇒ then Bayes best classifier
- Limitations: real life we don't know everything about data (prior, density function...)
- $\Sigma$  very hard to calculate



- Bayes error:
    - lowest possible error rate that any classifier could achieve on a given problem.
    - Bayes error all possible classifiers: linear or nonlinear
    - It's about the theoretical minimum error rate achievable by any classifier, given perfect knowledge of the data distribution.
- 

iid  $\hat{=}$  independent and identically distributed

→ Independent:  $p(xy) = p(x) \cdot p(y)$

→ Identically: all the random variables in the set have the same distribution  
Same mean, variance, PDF (prob. density function)

$$\sum p(x) = 1$$


---

Harmonic mean:  $\min(a,b) \leq \frac{2}{\frac{1}{a} + \frac{1}{b}} = \frac{2ab}{a+b}$

# Lecture 2: Parameter Estimation

## Learning $p(x | \omega_j)$ : Model-Based

we know it's Gaussian, but we don't know the parameters  $(\mu_j, \Sigma_j)$

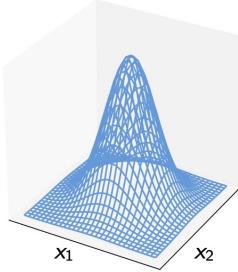
Idea:

- Assume that  $p(x | \omega_j)$  is a known parametric form, e.g.  $\mathcal{N}(\mu_j, \Sigma_j)$  where  $\mu_j, \Sigma_j$  are the parameters of the distribution.
- Estimate the parameters of the distribution that best fit the few observations we have. that explains the generated data the best

Advantage:

- With the model-based approach, we need to estimate a finite and typically small number of model parameters and not the whole data distribution. much smaller and not exponential

# parameters = # dimensions



## Max Likelihood, Simple Gaussian Case

Assume the data density is modeled as a univariate Gaussian with unit variance and unknown mean  $\theta$ . For a given data point  $x_k$ , the density function can be written as:

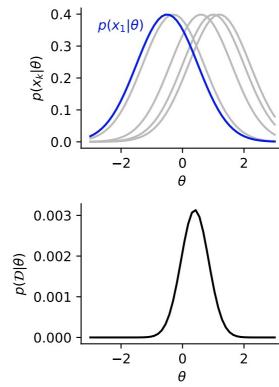
$$p(x_k | \theta) = \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2}(x_k - \theta)^2 \right]$$

Using the independence assumption, the joint density becomes:

$$p(\mathcal{D} | \theta) = \prod_{k=1}^N \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2}(x_k - \theta)^2 \right]$$

Question: How to find  $\theta$  that maximizes the function  $p(\mathcal{D} | \theta)$ .

we're looking for  $\theta$  that best represent our dataset  $\mathcal{D}$ .



## Max Likelihood, Simple Gaussian Case

Observation: The function  $p(\mathcal{D} | \theta)$  is not concave with  $\theta$ .

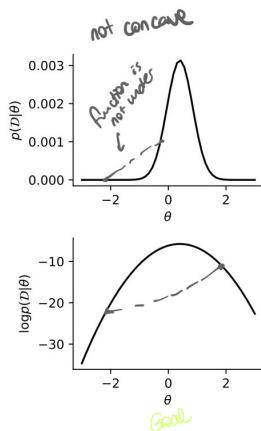
Idea: Transform the function in a way that

- doesn't change its argmax and
- makes it concave.

Applying the logarithm ensures the two properties above:

$$\begin{aligned} \log p(\mathcal{D} | \theta) &= \log \prod_{k=1}^N \frac{1}{\sqrt{2\pi}} \exp \left[ -\frac{1}{2}(x_k - \theta)^2 \right] \\ &= \sum_{k=1}^N \left[ -\frac{1}{2} \log(2\pi) - \frac{1}{2}(x_k - \theta)^2 \right] \end{aligned}$$

Cst. Concave



## Maximum Likelihood Estimation

Goal: Let  $\{p(x | \theta), \theta \in \Theta\}$  be a set of density functions (e.g. Gaussian), where  $\theta$  denotes a parameter vector (e.g. mean / covariance). We would like to find the parameter  $\theta$  that is the most likely with respect to the data.

Maximum Likelihood (ML) Approach:

- Assume that we have a dataset  $\mathcal{D} = (x_1, \dots, x_N)$ .
- Assume that each example  $x_k \in \mathbb{R}^d$  in the dataset has been generated independently and from the same density function  $p(x | \theta)$ .
- In that case, the joint density function can be written as:

$$p(\mathcal{D} | \theta) = \prod_{k=1}^N p(x_k | \theta)$$

product of all the densities

We also call this quantity the likelihood of  $\theta$  w.r.t. the dataset  $\mathcal{D}$ .

- The best parameter is then given by  $\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} | \theta)$ .

## Finding the Maximum of a Function

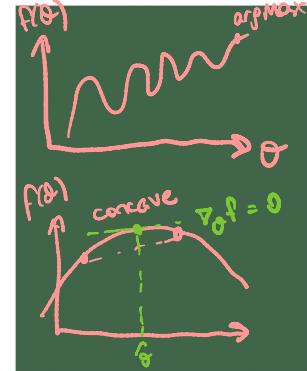
For some function  $f$  of interest (e.g. the data likelihood) we would like to compute:

$$\hat{\theta} = \arg \max_{\theta} f(\theta)$$

When the function to optimize is continuously differentiable and concave, the maximum is found at the point where the gradient is zero.

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \frac{\partial f}{\partial \theta_h} \end{bmatrix} = \mathbf{0}$$

Concave check: 2nd derivative  $\neq 0 \rightarrow \frac{\partial^2}{\partial \theta^2} p(\mathcal{D} | \theta) \neq 0$



## Max Likelihood, Simple Gaussian Case

Having found the log-likelihood w.r.t.  $\mathcal{D}$  to be

$$\log p(\mathcal{D} | \theta) = \sum_{k=1}^N \left[ -\frac{1}{2} \log(2\pi) - \frac{1}{2}(x_k - \theta)^2 \right],$$

Cst.  $\cancel{\log(2\pi)}$

the best parameter  $\hat{\theta}$  can then be found by solving  $\nabla_{\theta} \log P(\mathcal{D} | \theta) = 0$ .

$$\begin{aligned} \nabla_{\theta} \log P(\mathcal{D} | \theta) &= \sum_{k=1}^N -\frac{1}{2} \cdot 2(x_k - \theta) \cdot (-1) = \sum_{k=1}^N x_k - \theta = 0 \\ \Rightarrow \sum_{k=1}^N x_k &= \sum_{k=1}^N \theta = N \cdot \theta \\ \hat{\theta} &= \frac{1}{N} \sum_{k=1}^N x_k \end{aligned}$$

empirical mean

choose  $\theta$  that best explains the data

so  $\theta = \bar{x}$

## Max Likelihood, Multivariate Case

extension of Gaussian to multivariate

$\Sigma^{-1}$  is PSD  $\rightarrow \Sigma^{-1}$  neg. def.  $\rightarrow$  concave function

The log-likelihood of a multivariate Gaussian distribution w.r.t.  $\mathcal{D}$  is given by

$$\log p(\mathcal{D} | \mu) = \sum_{k=1}^N -\frac{1}{2} \log [(2\pi)^d \det(\Sigma)] - \frac{1}{2} (\mathbf{x}_k - \mu)^\top \Sigma^{-1} (\mathbf{x}_k - \mu)$$

does not depend on  $\mu$

Question: Assuming  $\Sigma$  is fixed, what is the optimal parameter vector  $\mu$ ?

$$\nabla_{\mu} \log p(\mathcal{D} | \mu) = \sum_{k=1}^N \mathbf{0} + \Sigma^{-1} (\mathbf{x}_k - \mu) = 0$$

emp. mean of data

$$\Sigma^{-1} \Sigma^{-1} \mathbf{x}_k = N \Sigma^{-1} \mu \rightarrow \hat{\mu} = \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k$$

$$(-\frac{1}{2} \cdot \Sigma^{-1} (\mathbf{x}_k - \mu)^2)' = -\Sigma^{-1} (\mathbf{x}_k - \mu) \cdot (-1)$$

shift  $\mu$  to  $\bar{x}$  to maximally explain the data

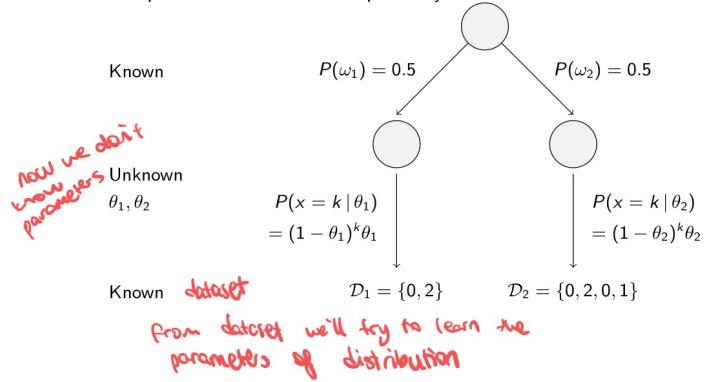
## Building a Classifier with ML

## Building a Classifier with ML

① Learn  $\theta_1, \theta_2$  given we have data

② Use  $\theta_1, \theta_2$  in Bayes to build the max. accurate classifier

Consider a labeled dataset containing two examples for the first class, and four examples for the second class. Points are in  $\mathbb{N}_0$  and we assume they are generated from classes  $\omega_1, \omega_2$  following geometric probability distributions of unknown parameters  $\theta_1$  and  $\theta_2$  respectively.



## Building a Classifier with ML

## Building a Classifier with ML

Recall:

$$\begin{aligned} P(\omega_1) &= 0.5 \\ P(\omega_2) &= 0.5 \end{aligned}$$

$$\begin{aligned} P(x = k | \theta_1) &= (1 - \theta_1)^k \theta_1 \\ P(x = k | \theta_2) &= (1 - \theta_2)^k \theta_2 \end{aligned}$$

$$\begin{aligned} \mathcal{D}_1 &= \{0, 2\} \\ \mathcal{D}_2 &= \{0, 2, 0, 1\} \end{aligned}$$

gradient:

$$\begin{aligned} P(D_1 | \theta_1) &= (1 - \theta_1)^2 \theta_1 = (1 - \theta_1)^2 \cdot \theta_1 \\ &= \theta_1 \cdot (1 - \theta_1)^2 \cdot \theta_1 \\ P(x = 0 | \theta_1) \cdot P(x = 2 | \theta_1) &= \theta_1^2 (1 - \theta_1)^2 \quad \text{not convex, apply log} \\ \log P(D_1 | \theta_1) &= 2 \log \theta_1 + 2 \log (1 - \theta_1) \\ \log (1 - \theta_1)^2 &= \frac{2}{\theta_1} + \frac{2}{(1 - \theta_1)} \cdot (-1) \stackrel{!}{=} 0 \\ \theta_1 &= 1/2 \end{aligned}$$

our ML estimator for the parameters of first class ( $\omega_1$ ).

## Bayes decision theory+ML example after finding the full model, apply Bayes

what do you want to classify?

Assume a new datapoint  $x=1$  is coming.

Recall:

$$\begin{aligned} P(\omega_1) &= 0.5 \\ P(\omega_2) &= 0.5 \end{aligned}$$

$$\begin{aligned} P(x = k | \theta_1) &= (1 - \theta_1)^k \theta_1 \\ P(x = k | \theta_2) &= (1 - \theta_2)^k \theta_2 \end{aligned}$$

$$\mathcal{D}_1 = \{0, 2\}$$

$$\mathcal{D}_2 = \{0, 2, 0, 1\}$$

$$\begin{aligned} \underset{j}{\operatorname{argmax}} P(\omega_j | x) &= \underset{j}{\operatorname{argmax}} P(x | \omega_j) P(\omega_j) \\ \omega_1: P(x=1 | \omega_1) \cdot P(\omega_1) &= (1 - \theta_1)^1 \theta_1 \cdot P(\omega_1) \\ &= \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = 0.125 \\ \omega_2: P(x=1 | \omega_2) \cdot P(\omega_2) &= (1 - \theta_2)^1 \theta_2 \cdot P(\omega_2) \\ &= (1 - \theta_2)^1 \theta_2 \cdot 0.5 = 0.124 \end{aligned}$$

outcome of classifier

$\Rightarrow$  decide (argmax)  $\omega_1$

Recall:

$$P(\omega_1) = 0.5$$

$$P(\omega_2) = 0.5$$

$$P(x = k | \theta_1) = (1 - \theta_1)^k \theta_1$$

$$P(x = k | \theta_2) = (1 - \theta_2)^k \theta_2$$

$$\mathcal{D}_1 = \{0, 2\}$$

$$\mathcal{D}_2 = \{0, 2, 0, 1\}$$

$$P(D_2 | \theta_2) = \theta_2 \cdot (1 - \theta_2)^2 \theta_2 \cdot \theta_2 \cdot (1 - \theta_2) \theta_2$$

$$\begin{aligned} D_2 &= \{0, 1, 2, 2\} \\ &= \theta_2^4 (1 - \theta_2)^3 \end{aligned}$$

$$\log P(D_2 | \theta_2) = 4 \log \theta_2 + 3 \log (1 - \theta_2)$$

$$\rightarrow 4\theta_2 - 3(1 - \theta_2) = 0$$

$$4\theta_2 + 3\theta_2 = 3$$

$$7\theta_2 = 3$$

$$\theta_2 = 3/7$$

## Bayes Parameter Estimation

p(θ), so in bayes param. est. user needs to specify what user believes θ is before observing any data

Assuming some dataset  $\mathcal{D} = (x_1, \dots, x_N)$ . We would like to model this dataset using some probability function  $p(\mathcal{D} | \theta)$  with  $\theta$  some unknown parameter that needs to be learned.

**ML parameter estimation:** Choose the parameter  $\theta$  that maximises the data likelihood:

$$\hat{\theta} = \arg \max_{\theta} p(\mathcal{D} | \theta)$$

**Bayes parameter estimation:** Instead of learning a fixed estimate, build a posterior distribution over this parameter using the Bayes theorem:

$$p(\theta | \mathcal{D}) = \frac{p(\mathcal{D} | \theta) p(\theta)}{p(\mathcal{D})} \rightarrow \text{prior distribution over new parameter } \theta$$

find a distribution over all parameters

$$= \frac{p(\mathcal{D} | \theta) p(\theta)}{\int_{\theta} p(\mathcal{D} | \theta) p(\theta) d\theta}$$

This approach requires to define a prior distribution  $p(\theta)$  which models our initial belief about the parameter  $\theta$  before observing the data.

## From ML to Bayes Classifiers

How to make a classifier from  $P(\theta|D)$ ?

**ML classifier:** Class posteriors are given by:

$$P(\omega_i | x) = \frac{p(x | \omega_i; \hat{\theta}_i) P(\omega_i)}{\sum_i p(x | \omega_i; \hat{\theta}_i) P(\omega_i)}$$

where  $\hat{\theta}_i$  is the maximum likelihood parameter for the distribution of class  $\omega_i$ .

**Bayes classifier:** Class posteriors are computed by bypassing the intermediate computation of the parameters  $\hat{\theta}_i$ , and instead conditioning directly on the data:

$$\begin{aligned} P(\omega_i | x, D) &= \frac{p(x | \omega_i, D) P(\omega_i | D)}{\sum_i p(x | \omega_i, D) P(\omega_i | D)} \\ &= \frac{p(x | \omega_i, D_i) P(\omega_i)}{\sum_i p(x | \omega_i, D_i) P(\omega_i)} \quad \text{class priors are known} \end{aligned}$$

condition the prior to dataset  
Goal: Connect this posterior  $\leftrightarrow$  the posterior of data,  $P(\theta|D)$ .

## Bayes Classifiers (cont.)

The terms of the class posterior:

$$P(\omega_i | x, D) = \frac{p(x | \omega_i, D_i) P(\omega_i)}{\sum_i p(x | \omega_i, D_i) P(\omega_i)}$$

can be expressed with model parameters as:

$$\begin{aligned} p(x | \omega_i, D_i) &= \int p(x | \theta_i, \omega_i, D_i) p(\theta_i | \omega_i, D_i) d\theta_i \\ &= \int p(x | \theta_i) p(\theta_i | D_i) d\theta_i \end{aligned}$$

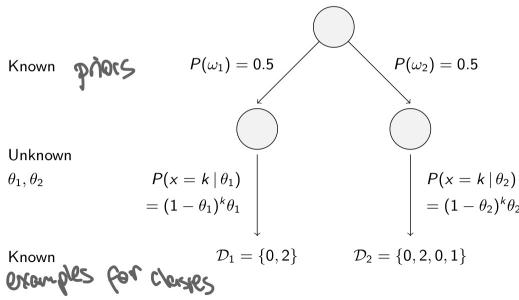
and

$$p(\theta_i | D_i) = \frac{p(D_i | \theta_i) p(\theta_i)}{\int p(D_i | \theta_i) p(\theta_i) d\theta_i}. \quad (1)$$

are our Bayes parameter estimates.

## Building a Classifier with Bayes

Recall: we consider the following data generating process



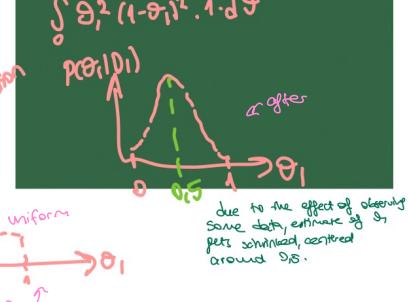
## Building a Classifier with Bayes

① Estimate a  $\theta_i$ :

$$\begin{aligned} p(\theta_i | D_i) &= \frac{p(D_i | \theta_i) p(\theta_i)}{\int p(D_i | \theta_i) p(\theta_i) d\theta_i} \quad \text{posterior} \\ &= \frac{\theta_i^2 (1-\theta_i)^2 \cdot 1}{\int \theta_i^2 (1-\theta_i)^2 \cdot 1 \cdot d\theta_i} \approx 20\theta_i^2 (1-\theta_i)^2 \end{aligned}$$

... and we further set:

$$\begin{aligned} p(\theta_1) &\sim U(0, 1) \\ p(\theta_2) &\sim U(0, 1) \end{aligned}$$



19/26

## Building a Classifier with Bayes

Recall:

$$\begin{aligned} P(\omega_1) &= 0.5 \\ P(\omega_2) &= 0.5 \\ P(x = k | \theta_1) &= (1 - \theta_1)^k \theta_1 \\ P(x = k | \theta_2) &= (1 - \theta_2)^k \theta_2 \\ D_1 &= \{0, 2\} \\ D_2 &= \{0, 2, 0, 1\} \end{aligned}$$

... and we further set:

$$\begin{aligned} \text{prior distribution} & \text{uniform:} \\ p(\theta_1) &\sim U(0, 1) \quad p(\theta_1)=1 \\ p(\theta_2) &\sim U(0, 1) \quad p(\theta_2)=2 \\ & \text{we choose these} \end{aligned}$$

$$\begin{aligned} p(\theta_2 | D_2) &= \frac{p(D_2 | \theta_2) \cdot p(\theta_2)}{\int p(D_2 | \theta_2) \cdot p(\theta_2) d\theta_2} \\ &= \frac{\theta_2^4 (1-\theta_2)^3 \cdot 1}{\int \theta_2^4 (1-\theta_2)^3 \cdot 1 \cdot d\theta_2} \\ &= 280 \cdot \theta_2^4 (1-\theta_2)^3 \\ p(\theta_2 | D_2) & \text{graph} \end{aligned}$$

find an case  $x=1$

Recall:

$$\begin{aligned} P(\omega_1) &= 0.5 \\ P(\omega_2) &= 0.5 \\ P(x = k | \theta_1) &= (1 - \theta_1)^k \theta_1 \\ P(x = k | \theta_2) &= (1 - \theta_2)^k \theta_2 \\ D_1 &= \{0, 2\} \\ D_2 &= \{0, 2, 0, 1\} \end{aligned}$$

... and we further set:

$$\begin{aligned} p(\theta_1) &\sim U(0, 1) \\ p(\theta_2) &\sim U(0, 1) \end{aligned}$$

① estimate post. dist. of each class with parameters  
② express the prob. of new datapoint coming from each class

$$\begin{aligned} p(x | \omega_1, D_1) &= \int p(x | \theta_1) \cdot p(\theta_1 | D_1) d\theta_1 \\ &= \int (1-\theta_1) \cdot \theta_1 \cdot 30 \cdot \theta_1^2 (1-\theta_1)^2 d\theta_1 \\ &= \int (1-\theta_1)^3 \cdot \theta_1^3 \cdot 30 d\theta_1 \\ &= 3/14 = 0.214 \\ p(x | \omega_2, D_2) &= \int p(x | \theta_2) \cdot p(\theta_2 | D_2) d\theta_2 \\ &= \int (1-\theta_2)^4 \cdot \theta_2^4 \cdot 280 d\theta_2 = \frac{2}{9} \approx 0.222 \end{aligned}$$

## Building a Classifier with Bayes

③ the decision

Recall:

$$P(\omega_1) = 0.5$$

$$P(\omega_2) = 0.5$$

$$P(x = k | \theta_1) = (1 - \theta_1)^k \theta_1$$

$$P(x = k | \theta_2) = (1 - \theta_2)^k \theta_2$$

$$\theta_1 = \{0, 2\}$$

$$\theta_2 = \{0, 2, 0, 1\}$$

... and we further set:

$$p(\theta_1) \sim U(0, 1)$$

$$p(\theta_2) \sim U(0, 1)$$

$$\underset{j}{\operatorname{argmax}} p(x | \omega_j, \theta_j) \cdot p(\omega_j)$$

$\omega_1: 0.214, 0.5$

$\omega_2: 0.222, 0.5$

decide for  $\omega_2$

ML said decide for  $\omega_1$ , but Bayes said decide for  $\omega_2$ . Why?

① Because in Bayes we introduced a new modelling set:  $P(\theta_1) \sim U(0, 1)$   
 $P(\theta_2) \sim U(0, 1)$

With this we influenced our initial estimate of parameters

② More observed data means more shrinkage the distribution will have.

Dataset 2 had more data, so its estimate is better.  
 → more confident as a classifier.

## ML vs. Bayes:

Observations:

- ▶ ML and Bayes classifiers do not always produce the same decisions
- ▶ Bayes classifiers are influenced by the prior distribution and are consequently less sensitive to the data.
- ▶ Bayes classifiers will tend to favor the outcome that is supported by a larger amount of data.

• Bayes posterior is influenced by the prior, if little data is available.  
 • Bayes also provides a confidence estimate, ML just gives a number, without any interval of confidence.

- ▶ In practice, parameters of the class-conditioned distributions are not known and they must be inferred from some finite dataset.
- ▶ Two main approaches for learning these parameters: (1) maximum likelihood estimation and (2) Bayesian inference.
- ▶ Bayesian inference is often more difficult than maximum likelihood estimation (both analytically and computationally), because it requires integration.
- ▶ Bayesian inference readily incorporates interesting functionalities (inclusion of priors, construction of confidence estimates).

ML-Estimation Steps:

① Objective Function

② transform objective function to concave with log, so argmax doesn't change

$$\log \Pi = \log + \log -$$

$$\begin{aligned} \textcircled{1} \quad p(\mathcal{D} | \theta) &= \prod_{k=1}^N \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x_k - \theta)^2\right] \\ \Rightarrow \log p(\mathcal{D} | \theta) &= \log \prod_{k=1}^N \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{1}{2}(x_k - \theta)^2\right] \\ &= \sum_{k=1}^N \left[ \log\left(\frac{1}{\sqrt{2\pi}}\right) + \log\left(\exp\left[-\frac{1}{2}(x_k - \theta)^2\right]\right) \right] \\ &= \sum_{k=1}^N \left[ -\frac{1}{2}\log(2\pi) - \frac{1}{2}(x_k - \theta)^2 \right] \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \theta} \log p(\mathcal{D} | \theta) &= \frac{\partial}{\partial \theta} \sum_{k=1}^N \left[ -\frac{1}{2}\log(2\pi) - \frac{1}{2}(x_k - \theta)^2 \right] \stackrel{!}{=} 0 \\ &= \sum_{k=1}^N -\frac{1}{2} \cdot 2(x_k - \theta) \cdot (-1) = 0 \\ \Leftrightarrow \sum_{k=1}^N x_k &= N\theta \\ \Rightarrow \hat{\theta} &= \frac{1}{N} \sum_{k=1}^N x_k \end{aligned}$$

optimal  $\theta$  of our ML  
is the mean of  
our data  
for multivariate Gaussians  
as well.

Var is fix

Consider the simple data density  $p(x | \mu) = \mathcal{N}(\mu, \sigma^2)$  with unknown parameter  $\mu$ . We would like to compare the ML and Bayes approaches to estimate the parameter  $\mu$  from some dataset  $\mathcal{D} = \{x_1, \dots, x_n\}$ .

ML approach:

- ▶ The maximum likelihood estimate is given by  $\hat{\mu} = \frac{1}{n} \sum_{k=1}^n x_k$ , i.e. the empirical mean (cf. previous slides).

Bayes approach:

- ▶ Assuming some prior distribution  $p(\mu) = \mathcal{N}(\mu_0, \sigma_0^2)$ , the posterior distribution can be computed as:

$$p(\mu | \mathcal{D}) = \frac{p(\mathcal{D} | \mu)p(\mu)}{p(\mathcal{D})} = \alpha \prod_{k=1}^n p(x_k | \mu)p(\mu)$$

where  $\alpha$  is a normalizing factor.

→ more data, more our Bayes will be closer to Norm. Likelihood

• Our goal is to find suitable parameters in real world  
 → few datapoints:  $\mu$  closer to  $\mu_0$  (prior)  
 • Bayes is more expensive, but it also gives confident estimates.

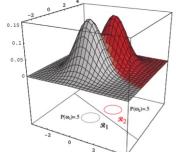
## Recap: Parameter Estimation

Example of estimator:

- ▶ Maximum likelihood estimator:

$$\mu = \frac{1}{N} \sum_{k=1}^N x_k$$

$$\widehat{\Sigma} = \frac{1}{N} \sum_{k=1}^N (x_k - \mu)(x_k - \mu)^\top$$



Problem:

- ▶ The covariance matrix (and its inverse) may be difficult to estimate.
- ▶ We make an assumption about the data (e.g. Gaussian-distributed) which may not correspond to reality.

How to turn ML to Bayes?

1. calculate  $P(D|\theta)$  from ML : a.  $P(D|\theta)$  likelihood fct.

b. find  $\hat{\theta}$ : function  $\rightarrow \log \rightarrow$  derivative

2.  $P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{\int_0^1 P(D|\theta) \cdot P(\theta) d\theta}$

$$\theta \in [0, 1]$$

take the integral

---

What is the likelihood of  $\theta$ ?  $P(D|\theta) = \prod_{k=1}^N P(x_k|\theta)$  likelihood of  $\theta$

where  $D$  contains  $N$  samples,  $x_1, \dots, x_N$

# Lecture 3: Fisher Discriminant

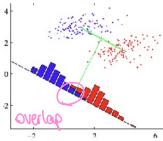
(LDA)  $\Sigma^{-1}$

## Derivation of Mean Separation

- The constrained optimization problem (subject to the constraint  $\|\mathbf{w}\| = 1$ ) can be developed as:

$$\arg \max_{\mathbf{w}} |\mu_2(\mathbf{w}) - \mu_1(\mathbf{w})|$$

$$\begin{aligned} &= \arg \max_{\mathbf{w}} \left| \frac{1}{N_2} \sum_{k \in C_2} z_k(\mathbf{w}) - \frac{1}{N_1} \sum_{k \in C_1} z_k(\mathbf{w}) \right| \\ &= \arg \max_{\mathbf{w}} \left| \frac{1}{N_2} \sum_{k \in C_2} \mathbf{w}^T \mathbf{x}_k - \frac{1}{N_1} \sum_{k \in C_1} \mathbf{w}^T \mathbf{x}_k \right| \\ &= \arg \max_{\mathbf{w}} \left| \mathbf{w}^T \left( \frac{1}{N_2} \sum_{k \in C_2} \mathbf{x}_k \right) - \mathbf{w}^T \left( \frac{1}{N_1} \sum_{k \in C_1} \mathbf{x}_k \right) \right| \\ &= \arg \max_{\mathbf{w}} \left| \mathbf{w}^T (\mu_2 - \mu_1) \right| \end{aligned}$$



where  $\mu_2$  and  $\mu_1$  are the means in input space.

- The best vector  $\mathbf{w}$  is the one that aligns with  $(\mu_2 - \mu_1)$ , i.e.  $\mathbf{w} = (\mu_2 - \mu_1) / \|\mu_2 - \mu_1\|$ .

## Linear Discriminant Analysis



**Goal:** Find a (normal vector of a linear decision boundary)  $\mathbf{w} \in \mathbb{R}^D$  that  
Maximizes mean class difference

$$(\mathbf{w}^T \mathbf{w}_o - \mathbf{w}^T \mathbf{w}_\Delta)^2 = \mathbf{w}^T \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)(\mathbf{w}_o - \mathbf{w}_\Delta)^T}_{S_B \text{ -- "between class scatter"}} \mathbf{w} \quad (2)$$

cov-matrix

Minimizes variance in each class

$$\begin{aligned} &\frac{1}{N_o} \sum_{i=1}^{N_o} (\mathbf{w}^T (\mathbf{x}_{oi} - \mathbf{w}_o))^2 + \frac{1}{N_\Delta} \sum_{j=1}^{N_\Delta} (\mathbf{w}^T (\mathbf{x}_{oj} - \mathbf{w}_\Delta))^2 \\ &= \mathbf{w}^T \underbrace{\left( \frac{1}{N_o} \sum_{i=1}^{N_o} (\mathbf{x}_{oi} - \mathbf{w}_o)(\mathbf{x}_{oi} - \mathbf{w}_o)^T + \frac{1}{N_\Delta} \sum_{j=1}^{N_\Delta} (\mathbf{x}_{oj} - \mathbf{w}_\Delta)(\mathbf{x}_{oj} - \mathbf{w}_\Delta)^T \right)}_{S_W \text{ -- "within class scatter"}} \mathbf{w} \end{aligned}$$

## Fisher LDA

### Algorithm



**Computes:** Normal vector  $\mathbf{w}$  of decision hyperplane for binary classification

**Input:** Data  $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_i \in \mathbb{R}^D, y_i \in \{-1, +1\}$ .

① Compute class mean vectors

$$\mathbf{w}_{-1} = \sum_{i \in \mathcal{Y}_{-1}} \mathbf{x}_i$$

$$\mathbf{w}_{+1} = \sum_{j \in \mathcal{Y}_{+1}} \mathbf{x}_j$$

② Compute *between-class covariance matrix*

$$S_B = (\mathbf{w}_{-1} - \mathbf{w}_{+1})(\mathbf{w}_{-1} - \mathbf{w}_{+1})^T$$

③ Compute *within-class covariance matrices*

$$S_W = \sum_{i \in \mathcal{Y}_{-1}} (\mathbf{x}_i - \mathbf{w}_{-1})(\mathbf{x}_i - \mathbf{w}_{-1})^T$$

$$+ \sum_{j \in \mathcal{Y}_{+1}} (\mathbf{x}_j - \mathbf{w}_{+1})(\mathbf{x}_j - \mathbf{w}_{+1})^T$$

④ Compute eigenvalue decomposition

$$S_B \mathbf{w} = S_W \mathbf{w} \lambda \rightarrow \text{generalized EV}$$

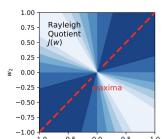
**Output:**  $\mathbf{w}$

$$\text{Fisher Discriminant: } J(\mathbf{w}) = \frac{(\mu_2(\mathbf{w}) - \mu_1(\mathbf{w}))^2}{S_1(\mathbf{w}) + S_2(\mathbf{w})}$$

$$\Rightarrow J(\mathbf{w}) \approx \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

A solution that maximizes the Rayleigh quotient  $J(\mathbf{w})$  can be obtained by first observing that  $\nabla_{\mathbf{w} \neq 0} : J(\mathbf{w}) = J(\mathbf{w})$  and searching for the particular solution for which the denominator is exactly one. This can be stated as the constrained optimization problem

$$\arg \max_{\mathbf{w}} \mathbf{w}^T S_B \mathbf{w} \text{ s.t. } \mathbf{w}^T S_W \mathbf{w} = 1$$



Optimization problem:

$$\arg \max_{\mathbf{w}} \mathbf{w}^T S_B \mathbf{w} \quad \text{s.t. } \mathbf{w}^T S_W \mathbf{w} = 1$$

pay off 1 obtain

$$\mathbf{w}^T S_W \mathbf{w} - 1 = 0 \rightarrow \mathbf{w}^T S_W \mathbf{w} = 1$$

Step 1: Lagrange

$$L(\mathbf{w}) = \mathbf{w}^T S_B \mathbf{w} + \lambda (1 - \mathbf{w}^T S_W \mathbf{w})$$

Step 2:  $\nabla L(\mathbf{w}, \lambda) = 0 \rightarrow S_B \mathbf{w} = \lambda S_W \mathbf{w}$

generalized EV-problem

Limitations of mean separation:

- class overlap in projected space
- if we rotate the projection a few degrees, then we can achieve a better classifier
- Making means distant may not be enough to induce class separability in projected space

## Fisher Discriminant

Idea:

- In addition to maximizing the separation between class means in projected space, also consider to reduce the within-class variance.



$$\mu_1 = \frac{1}{|C_1|} \sum_{k \in C_1} z_k \quad \mu_2 = \frac{1}{|C_2|} \sum_{k \in C_2} z_k$$

$$S_1 = \sum_{k \in C_1} (z_k - \mu_1)^2 \quad S_2 = \sum_{k \in C_2} (z_k - \mu_2)^2$$

- Maximizing distance between means while minimizing within-class variance can be formulated as:

$$\arg \max_{\mathbf{w}} \frac{(\mu_2(\mathbf{w}) - \mu_1(\mathbf{w}))^2}{S_1(\mathbf{w}) + S_2(\mathbf{w})} \quad \begin{matrix} \text{between means max} \\ \text{within variance min} \end{matrix}$$

How to optimize Fisher criterion:

$$\left( \frac{\mathbf{w}}{\sqrt{w}} \right)' = \frac{\mathbf{w}' - \mathbf{w}'}{\sqrt{w^2}}$$

$$\frac{(\mathbf{w}^T S_W \mathbf{w}) S_B \mathbf{w} - (\mathbf{w}^T S_B \mathbf{w}) S_W \mathbf{w}}{(\mathbf{w}^T S_W \mathbf{w})^2 \text{ cst.}} = 0$$

$$(\mathbf{w}^T S_W \mathbf{w}) S_B \mathbf{w} = (\mathbf{w}^T S_B \mathbf{w}) S_W \mathbf{w}$$

$$S_W \mathbf{w} = S_B \mathbf{w} \underbrace{\frac{\mathbf{w}^T S_W \mathbf{w}}{\mathbf{w}^T S_B \mathbf{w}}}_{\lambda}$$

$$S_W \mathbf{w} = S_B \mathbf{w} \lambda$$

$$\mathbf{w} \perp S_W^{-1} (\mathbf{w}_o - \mathbf{w}_\Delta)$$

↳ mixed covariance

$$\text{within class: } S_1(\mathbf{w}) = \sum (z_k - \mu_j)^2 = \sum (\mathbf{w}^T \mathbf{x}_k - \mathbf{w}^T \mu_j)^2$$

$$\text{variance} = \sum (\mathbf{w}^T \mathbf{x}_k - \mathbf{w}^T \mu_j)(\mathbf{w}^T \mathbf{x}_k - \mathbf{w}^T \mu_j)$$

$$= \sum (\mathbf{w}^T \mathbf{x}_k \mathbf{w}^T \mathbf{x}_k - \mathbf{w}^T \mathbf{x}_k \mathbf{w}^T \mu_j - \mathbf{w}^T \mu_j \mathbf{w}^T \mathbf{x}_k + \mathbf{w}^T \mu_j \mathbf{w}^T \mu_j)$$

$$= \mathbf{w}^T (2 \mathbf{x}_k \mathbf{x}_k^T) \mathbf{w} - \mathbf{w}^T (2 \mathbf{x}_k \mathbf{\mu}_j^T) \mathbf{w} + \sum \mathbf{w}^T \mu_j \mu_j^T \mathbf{w}$$

$$= \underbrace{\mathbf{w}^T \sum (x_k - \mu_j)(x_k - \mu_j)^T \mathbf{w}}_{= S_{ij}} = S_{ij}$$

$$S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$$

Between class scatter

$$S_W = S_1 + S_2$$

Within class

"Rayleigh Coefficient"

Constrained formulation

$$\arg \max_{\mathbf{w}} \mathbf{w}^T S_B \mathbf{w} \quad \text{s.t. } \mathbf{w}^T S_W \mathbf{w} = 1$$

maximum

$$\mathbf{w}^T S_B \mathbf{w} = \lambda \mathbf{w}^T S_W \mathbf{w}$$



$$\begin{aligned}
 & \underset{\mathbf{w}, b, \mathbf{M}}{\text{argmin}} \frac{1}{N} \sum_{k=1}^N \max(0, \mathbf{M} - (\mathbf{w}^\top \mathbf{x}_k + b) t_k) + \frac{1}{\mathbf{M}} \\
 & = \underset{\mathbf{w}, b, \mathbf{M}}{\text{argmin}} \frac{1}{N} \sum_{k=1}^N \max(0, \mathbf{M} - (\mathbf{w}^\top \mathbf{x}_k + b) t_k) + \frac{\|\mathbf{w}\|}{\mathbf{M}} \\
 & = \underset{\mathbf{w}, b, \mathbf{M}}{\text{argmin}} \frac{1}{N} \sum_{k=1}^N \max(0, 1 - (\frac{\mathbf{w}^\top \mathbf{x}_k + b}{\mathbf{M}}) t_k) + \frac{\|\mathbf{w}\|^2}{\mathbf{M}} \\
 & = \underset{\mathbf{w}, b, \mathbf{M}}{\text{argmin}} \frac{1}{N} \sum_{k=1}^N \max(0, 1 - (\mathbf{w}^\top \mathbf{x}_k + b) t_k) + \frac{\|\mathbf{w}\|^2}{\mathbf{M}} \\
 & = \boxed{\underset{\mathbf{w}, b}{\text{argmin}} \frac{1}{N} \sum_{k=1}^N \max(0, 1 - (\mathbf{w}^\top \mathbf{x}_k + b) t_k) + \frac{\|\mathbf{w}\|^2}{\mathbf{M}}}
 \end{aligned}$$

## Recap: The Perceptron Algorithm

Consider the linear model  $y = \mathbf{w}^\top \mathbf{x} + b$ , where  $\text{sign}(y)$  gives the classification decision. Let  $t \in \{-1, +1\}$  be the binary class label associated to  $x$ .

### Algorithm

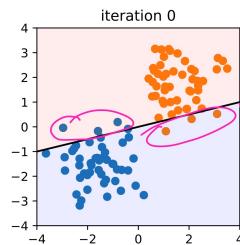
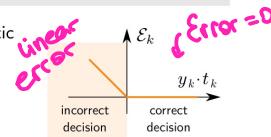
- Iterate over all data points  $\{(\mathbf{x}_k, t_k); k = 1 \dots, N\}$  (multiple times).
  - Compute  $y_k = \mathbf{w}^\top \mathbf{x}_k + b$
  - If  $x_k$  is correctly classified (i.e.  $\text{sign}(y_k) = t_k$ ), continue.
  - If  $x_k$  is wrongly classified (i.e.  $\text{sign}(y_k) \neq t_k$ ), apply:

$$\mathbf{w} \leftarrow \mathbf{w} + \gamma \cdot \mathbf{x}_k t_k \quad \text{and} \quad b \leftarrow b + \gamma \cdot t_k$$

- Stop once all examples are correctly classified.

The perceptron can also be seen as a stochastic gradient descent of the error function

$$E(\mathbf{w}, b) = \frac{1}{N} \sum_{k=1}^N \max(0, -y_k t_k)$$



→ hyperplane will rotate until all data are correctly classified  
the term  $y_k t_k$  changes.

→ change the hyperplane by rotating if data on the wrong side

→ adapting in high dimensions with  $\mathbf{x}$

→  $y_k t_k$  will be negative if data misclassified  
 $t = -1$  if data wrong classified  
 $-y_k t_k = 1 \rightarrow$  if data misclassified

→ Error fct. takes the max, so the error deviation will be big.

→ If data correct classified:  $t = 1$   
 $-y_k t_k = -1$

our  $\max(0, -y_k t_k)$  will choose 0  
so our Error function won't get anything.

- Margin: if margin between the data is large, then the hyperplane doesn't need to be that precise
- Small margin: perceptron runs long time
- Thinner the margin, longer the iterations

## Recap: Mean Separation Criterion

- We want to learn a projection of the data  $\mathbf{z}_k = \mathbf{w}^\top \mathbf{x}_k$ , with  $\|\mathbf{w}\| = 1$  such that the means of classes in projected space are as distant as possible.

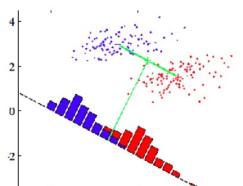
- First, we compute the means in projected space for the two classes

$$\mu_1 = \frac{1}{N_1} \sum_{k \in C_1} \mathbf{z}_k \quad \mu_2 = \frac{1}{N_2} \sum_{k \in C_2} \mathbf{z}_k$$

- Then we find  $\mathbf{w}$  that maximizes the difference of means, i.e. we express the means as a function of  $\mathbf{w}$  and pose the optimization problem:

$$\arg \max_{\mathbf{w}} |\mu_2(\mathbf{w}) - \mu_1(\mathbf{w})| \quad \text{with} \quad \|\mathbf{w}\| = 1$$

Take your data  $\mathbf{x}_k$   
apply projection on  $\mathbf{w}$   
Constraint as:  $\|\mathbf{w}\|=1$   
Goal: nice linear separation



Means as far away as possible.  
Max the separation  
 $\mu_1$  is a function of  $\mathbf{w}^\top \mathbf{x}_k$  with objectives  
 $\mu_2$  is also function of  $\mathbf{w}^\top \mathbf{x}_k$

## Summary

- Practice, train a classifier directly rather than learning the distributions
- Max mean separation and Fisher Discriminant: we only need to estimate the mean & cov of data without estimating full distributions.
- Perceptron (& its large-margin extension) focuses on the decision boundary, which leads to higher classification accuracy on general tasks.

→ Perceptron may not work if:

- ① data not linearly separable
- ② margin is very thin, takes a lot of time when fine tuning

→ If Perceptron converges, then data 100% linearly separable

→ if there's a solution the perceptron will find it in a finite number of steps.

Fisher LDA: → finds a linear combination of features that separates two or more classes

Fisher objective function:  $J(w) = \frac{w^T S_B w}{w^T S_w w}$

| max  
S<sub>B</sub>: Between class variance  
S<sub>w</sub>: within class variance

→ The max criterion assures that the classes are as separated as possible in the projected space, to get easier and more accurate classification.

- PCA seeks directions that are efficient for representation
- Fisher LDA seeks directions that are efficient for discrimination.

• Goal of LDA: move the db around and try to find a line that the projected samples are well separated.

Weight of Fisher LDA:  $w = S_w^{-1} (m_1 - m_2)$

When is the Fisher LDA minimized?  $\max_w \frac{w^T S_B w}{w^T S_w w}$

and its min if:  $w^T S_B w \stackrel{!}{=} 0$

$$S_B = (m_2 - m_1)(m_2 - m_1)^T \rightarrow w^T (m_2 - m_1)(m_2 - m_1)^T w \stackrel{!}{=} 0 \rightarrow \text{so } w \perp (m_2 - m_1)$$

$$\langle w, (m_2 - m_1) \rangle \stackrel{!}{=} 0 \rightarrow [w_1 \ w_2] \begin{bmatrix} m_1 \\ m_2 \end{bmatrix} \stackrel{!}{=} 0 \rightarrow w_1 m_1 + w_2 m_2 \stackrel{!}{=} 0$$

## Lecture 4: PCA

→ Curse of dimensionality

→ PCA

→ Lagrange

Goal: Reduce the dimension to the leading variance elements.

- min. noise & max. variance
- Lagrange multipliers to find PCA solutions
- PCA solutions  $\hat{=}$  EV's of the data covariance

### → Dimensionality Reduction : Curse of Dimensionality

• curse of dimensionality increases the models complexity in high-dimensional spaces, making it possibly overfitting.

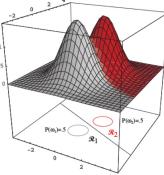
→ if  $d > N$ : ill-defined Many dimensions but not enough training data

→ if  $d \leq N$ : ill-conditioned / numerically unstable

### Example of ill-conditioning: Bayes Classifier with estimated Cov

▶ Assume our data is generated for each class  $\omega_j$  according to the multivariate Gaussian distribution  $p(\mathbf{x}|\omega_j) = \mathcal{N}(\boldsymbol{\mu}_j, \Sigma)$  and with class priors  $P(\omega_j)$ . The Bayes optimal classifier is derived as

$$\begin{aligned} & \arg \max_j \{P(\omega_j|\mathbf{x})\} \\ & = \arg \max_j \{\log p(\mathbf{x}|\omega_j) + \log P(\omega_j)\} \\ & = \arg \max_j \{\mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_j - \frac{1}{2} \boldsymbol{\mu}_j^\top \Sigma^{-1} \boldsymbol{\mu}_j + \log P(\omega_j)\} \end{aligned}$$



#### Example of estimator:

- ▶ Maximum likelihood estimator:

$$\hat{\Sigma} = \frac{1}{N} \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$$

where  $\bar{\mathbf{X}}$  is a matrix of size  $d \times N$  containing the centered data.

#### Problem:

- ▶ The matrix  $\hat{\Sigma}^{-1}$  used in the classifier only exists when  $d \leq N$  and the matrix  $\bar{\mathbf{X}}$  has full rank.
- ▶ The matrix  $\hat{\Sigma}^{-1}$  used in the classifier is typically accurate only when  $d \ll N$ . *need to have lot of data*

#### Solutions:

- ▶ Reduce the dimensionality of the data (**today**)
- ▶ Regularize the model (later)

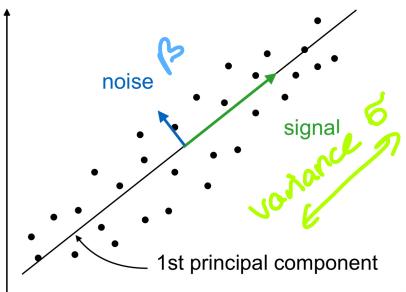
- because we don't know true  $\Sigma$ , we estimate  $\hat{\Sigma}$
- Goal: a good  $\hat{\Sigma}$  estimator
- Inverted Covariate  $\Sigma^{-1}$  is expensive to compute
- in practice: we don't know  $\mu$  or  $\Sigma$ . We just know the datapoints
- Reduce d and project data there

### f: Dimension T + space volume T available data becomes sparse

- ▶ The amount of data needed for a reliable result (e.g. no overfitting) often grows exponentially with the dimensionality.
- ▶ Here again, this can be addressed by reducing the dimensionality of the data (**today**) or regularizing (later).

## Principal Component Analysis

• PCA maps the data from a high dimension to a low dimension



Which line fits data best?  
The line  $w$  that minimizes the noise and maximizes the signal

$$\min \text{noise} \quad \max \text{signal}$$

# PCA as noise minimization 1

- Given a dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ , find a one-dimensional subspace, so that the data projected on that subspace has minimum distance to the original data.
- The reconstruction model is given by:

$$\hat{\mathbf{x}} = \mathbf{w}\mathbf{w}^\top \mathbf{x}$$

*estimate from subspace*

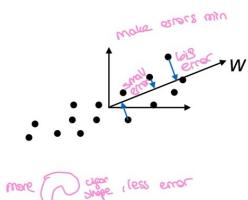
with  $\|\mathbf{w}\| = 1$ .

- Minimizing the noise (i.e. reconstruction error) can be written as:

$$\arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N \|\mathbf{x}_k - \hat{\mathbf{x}}_k\|^2 \right]$$

*error distances*

*original x - subspace x  
distance*

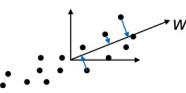


0 loss if: entire d is projected  
1 reconstruction error if: d-1 is projected

$\mathbf{w}\mathbf{w}^\top$  gives the noise and its norm is 1.  $\rightarrow \|\mathbf{w}\| = 1$

## From Min Noise to Max Signal

$$\begin{aligned} & \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N \|\mathbf{x}_k - \mathbf{w}\mathbf{w}^\top \mathbf{x}_k\|^2 \right] \\ &= \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{x}_k - \mathbf{w}\mathbf{w}^\top \mathbf{x}_k)^\top (\mathbf{x}_k - \mathbf{w}\mathbf{w}^\top \mathbf{x}_k) \right] \\ &= \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N \underbrace{\mathbf{x}_k^\top \mathbf{x}_k}_{\text{const. cut data's norm stays the same}} - 2\mathbf{x}_k^\top \mathbf{w}\mathbf{w}^\top \mathbf{x}_k + (\mathbf{w}\mathbf{w}^\top \mathbf{x}_k)^\top (\mathbf{w}\mathbf{w}^\top \mathbf{x}_k) \right] \\ &= \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N -2(\mathbf{x}_k^\top \mathbf{w})^2 + \mathbf{x}_k^\top \mathbf{w} \underbrace{\mathbf{w}^\top \mathbf{w}}_{\text{norm}=1} \mathbf{w}^\top \mathbf{x}_k \right] \\ &= \arg \min_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N -2(\mathbf{x}_k^\top \mathbf{w})^2 + (\mathbf{x}_k^\top \mathbf{w})^2 \right] \quad \text{argmin } (\cdot) = \text{argmax } (+) \\ &= \arg \max_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{x}_k^\top \mathbf{w})^2 \right] \quad \text{Max Variance} \end{aligned}$$



*⚠ Data needs to be centered*

# PCA as signal maximisation 2

- Given a dataset  $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^d$ . Find a one-dimensional subspace, so that the data projected on that subspace has maximum variance.
- The projection model is given by:

$$h = \mathbf{w}^\top \mathbf{x}$$

with  $\|\mathbf{w}\| = 1$ .

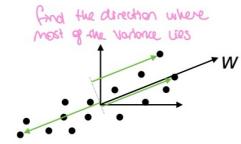
- Maximizing the signal (i.e. variance of projected data) can be written as:

$$\arg \max_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{w}^\top \mathbf{x}_k - \mathbb{E}[\mathbf{w}^\top \mathbf{x}])^2 \right]$$

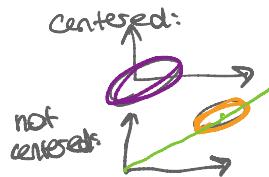
and assuming centered data, i.e.  $\mathbb{E}[\mathbf{x}] = \mathbf{0}$ , it simplifies to:

$$\arg \max_{\mathbf{w}} \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{w}^\top \mathbf{x}_k)^2 \right]$$

*Variance*  
*scalar*  
*how far we are apart?*



→ Projection of PCA won't be correct if we don't center:



- Gaussian
- Linear
- Centered

→ PCA is not robust to outliers, so remove outliers before doing PCA.

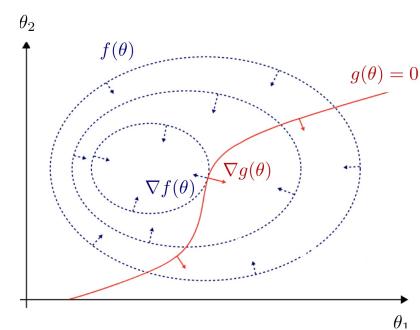
How to compute Principal Components?

Lagrange general framework:

$$\textcircled{1} \quad L(\theta, \lambda) = f(\theta) + \lambda g(\theta)$$

$$\textcircled{2} \quad \nabla L(\theta, \lambda) = 0$$

Lagrange



## Lagrange 2D Example:

Consider the constrained optimization problem

$$\arg \max_{\theta} [1 - (\theta_1^2 + \theta_2^2)] \quad \text{s.t.} \quad \theta_1 + \theta_2 = 1$$

- **Step 1:** Build the Lagrangian

$$\mathcal{L}(\theta, \lambda) = 1 - (\theta_1^2 + \theta_2^2) + \lambda \cdot (\theta_1 + \theta_2 - 1)$$

- **Step 2:** Set gradient of Lagrangian to zero:

$$\begin{aligned}\nabla_{\theta_1} \mathcal{L}(\theta, \lambda) &= 0 \Rightarrow 2\theta_1 = \lambda \\ \nabla_{\theta_2} \mathcal{L}(\theta, \lambda) &= 0 \Rightarrow 2\theta_2 = \lambda \\ \nabla_{\lambda} \mathcal{L}(\theta, \lambda) &= 0 \Rightarrow \theta_1 + \theta_2 = 1\end{aligned}$$

which gives us the solution:

$$\theta = (0.5, 0.5)$$

→ The 1st PC is the direction that the data's projection has the highest variance.

Solving EV-problem:  $\Sigma u = \lambda u$  gives the highest EV  $\lambda$  will give the direction.

Magnitude of  $\lambda$  represents the variance of the data projected onto the corresponding EV  $u$  so choosing the highest  $\lambda$  maximizes this variance.

## Solution PCA

projection problem  $\rightarrow \arg \max_w \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{x}_k^\top \mathbf{w})^2 \right] \quad \text{s.t.} \quad \|\mathbf{w}\| = 1$

can be rewritten as:

$$\begin{aligned}&= \arg \max_w \left[ \frac{1}{N} \sum_{k=1}^N (\mathbf{w}^\top \mathbf{x}_k)(\mathbf{x}_k^\top \mathbf{w}) \right] \quad \text{s.t.} \quad \|\mathbf{w}\|^2 = 1 \\ &= \arg \max_w \left[ \mathbf{w}^\top \underbrace{\left( \frac{1}{N} \sum_{k=1}^N \mathbf{x}_k \mathbf{x}_k^\top \right)}_{\widehat{\Sigma}} \mathbf{w} \right] \quad \text{s.t.} \quad \|\mathbf{w}\|^2 = 1\end{aligned}$$

where  $\widehat{\Sigma}$  is the covariance matrix (recall: the data is assumed to be centered). The covariance matrix does not depend on  $\mathbf{w}$  and can be precomputed.

Starting from our rewritten PCA optimization problem

$$\arg \max_w [\mathbf{w}^\top \widehat{\Sigma} \mathbf{w}] \quad \text{s.t.} \quad \|\mathbf{w}\|^2 = 1$$

we look for a solution by applying the method of Lagrange multipliers.

- **Step 1:** Build the Lagrangian

$$\mathcal{L}(\mathbf{w}, \lambda) = \mathbf{w}^\top \widehat{\Sigma} \mathbf{w} + \lambda \cdot (1 - \|\mathbf{w}\|^2)$$

- **Step 2:** Set gradient of Lagrangian to zero:

$$\begin{aligned}\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, \lambda) &= \mathbf{0} \Rightarrow \widehat{\Sigma} \mathbf{w} = \lambda \mathbf{w} \quad \text{EV-Problem} \\ \nabla_{\lambda} \mathcal{L}(\mathbf{w}, \lambda) &= 0 \Rightarrow \|\mathbf{w}\|^2 = 1\end{aligned}$$

! PCA solution is an eigenvector of  $\widehat{\Sigma}$

To find PC's:

- estimate cov matrices
- EV decomposition of cov matrix
- 1st EV = PCA solution

if you want more than 1 component:

- 1) Compute the 1st component
- 2) remove from data
- 3) recompute

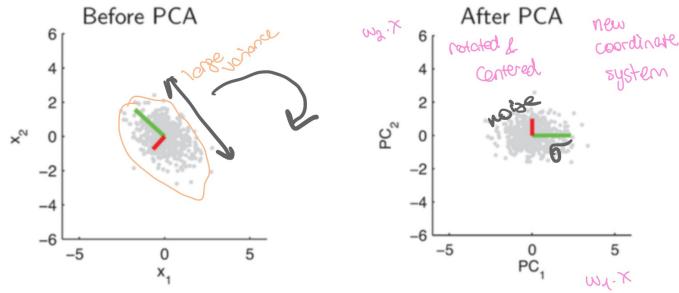
### Finding multiple PCA components

```
for i = 1 to h do
     $\mathbf{w}_i \leftarrow \text{PC}_i(\bar{X})$  1)
     $\bar{X} \leftarrow \bar{X} - \mathbf{w}_i \mathbf{w}_i^\top \bar{X}$  2)
end for
```

→ Resulting components  $\mathbf{w}_1, \dots, \mathbf{w}_k$  are the eigenvectors of  $\widehat{\Sigma}$

→ it's enough to only find the leading EV's of  $\widehat{\Sigma}$ .

→ Cov Matrix  $\widehat{\Sigma}$  = size of EV's



- ▶ PCA rotates data into new coordinate system with the directions of largest variance being the new coordinate axes.

## Stable PCA Computation via SVD

### Singular Value Decomposition (SVD)

A singular value decomposition factorizes a matrix  $M$  as  $U\Lambda V^T = M$  where

- ▶  $U$  contains the eigenvectors of  $MM^T$
- ▶  $V$  contains the eigenvectors of  $M^T M$
- ▶ The square roots of the eigenvalues of  $MM^T$  are on the diagonal of  $\Lambda$ .

**Observation:** *normalization*

- ▶ Setting  $M = \frac{1}{\sqrt{N}}\tilde{X}$  and computing SVD results in
  - ▶ a matrix  $U = (\mathbf{w}_1 | \dots | \mathbf{w}_d)$  that contains the eigenvectors of  $\hat{\Sigma}$ , i.e. the principal components, and
  - ▶ a matrix  $\Lambda = \text{diag}(\sqrt{\lambda_1}, \dots, \sqrt{\lambda_d})$  (but applying the square root does not change the ordering of eigenvalues).
- ▶ SVD is computationally faster and more stable than covariance computation + eigenvalue decomposition. → *Vector*

### PCA Computation via SVD, limitations

- ▶ The SVD has computational complexity of  $\mathcal{O}(\min(N^2 d, d^2 N))$ .
- ▶ When  $d \approx N$ , the complexity is roughly as bad as the that of computing the eigenvectors of  $\hat{\Sigma}$ , that is,  $\mathcal{O}(d^3)$ .
- ▶ This makes the computation of principal components using SVD prohibitive for large datasets.

**Note:**

- ▶ In practice, we often need **only the first few principal components**, whereas SVD computes all of them.

**Question:**

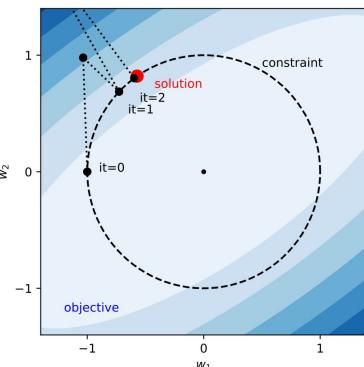
- ▶ Can we design a faster procedure that only extracts the first few principal components? **Power iteration**

## Power Iteration

Find the first component by applying the *iteration*

$$\begin{aligned}\mathbf{v}^{(t)} &\leftarrow \hat{\Sigma} \mathbf{w}^{(t-1)} \\ \mathbf{w}^{(t)} &\leftarrow \frac{\mathbf{v}^{(t)}}{\|\mathbf{v}^{(t)}\|}\end{aligned}$$

$T$  times, starting from a random (or any) unit-norm vector  $\mathbf{w}^{(0)}$ .



**Questions:**

- ▶ Does it always converge? **yes**
- ▶ How fast it converges? **exponentially fast**

• starts with a random vector  $\mathbf{w}^{(0)} \in \mathbb{R}^d$  and iteratively applies the parameter update:

$$\mathbf{w}^{(t+1)} = \frac{\hat{\Sigma} \mathbf{w}^{(t)}}{\|\hat{\Sigma} \mathbf{w}^{(t)}\|}$$

until some convergence is met

→ POWIT only gives us the **leading EV**  
→ If we want more PCA components, we need to compute them iteratively

**Naive approach (cf. slide 17): Recompute covariance at each step**

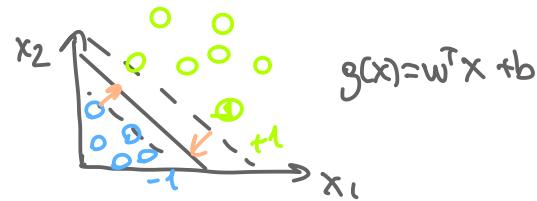
```
for i = 1 to h do
   $\hat{\Sigma} = \frac{1}{N} \tilde{X} \tilde{X}^T$ 
   $\mathbf{w}_i \leftarrow \text{POWIT}(\hat{\Sigma})$ 
   $\tilde{X} \leftarrow \tilde{X} - \mathbf{w}_i \mathbf{w}_i^T \tilde{X}$ 
end for
```

**Better approach: Work directly in covariance space** **more stable**

```
 $\hat{\Sigma} = \frac{1}{N} \tilde{X} \tilde{X}^T$ 
for i = 1 to h do
   $\mathbf{w}_i \leftarrow \text{POWIT}(\hat{\Sigma})$ 
   $\hat{\Sigma} \leftarrow \hat{\Sigma} - \mathbf{w}_i \mathbf{w}_i^T \hat{\Sigma}$ 
end for
```

# Lecture 5 : SVM (Support Vector Machines) Supervised classification

- hyperplane is  $d-1$  dimensions where  $n$  is our features  $(x_1, x_2, \dots, x_d)$   $\rightarrow \langle w, x \rangle + b = 0$
- support vectors are the closest ones to our hyperplane
- if we remove a SV, our hyperplane may shift
- margin: the distance our hyperplane is to the closest SV
- choose the hyperplane that:
  - separates data the most
  - largest margin



- distance between hyperplane and margin:  $\frac{\|w\|^2 + b\|}\|w\| = \frac{1+1}{\|w\|} = \frac{1}{\|w\|}$
- max total margin to get good data separation  
to max margin you need to  $\min \|w\|$
- SVM's can be used for linearly separable and non-linearly separable data -
- with appropriate nonlinear mapping  $\phi(\cdot)$  to a sufficiently high dimension, data can always be separated by a hyperplane.

## → Linear Soft Margin SVM

$$\text{Primal: } \min_{w, b, \zeta} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N \zeta_i \quad \text{s.t. } \forall_{i=1}^N: y_i \cdot (w^T \phi(x_i) + b) \geq 1 - \zeta_i \quad \text{and } \zeta_i \geq 0$$

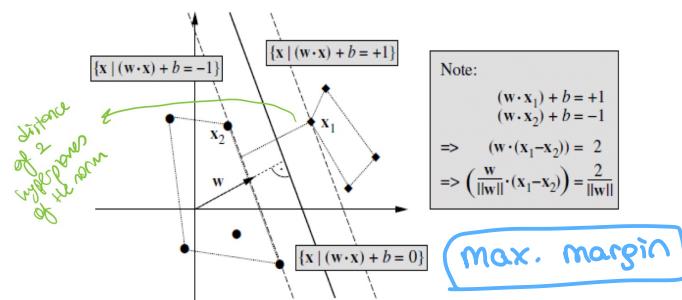
$\|\cdot\|$ : euclidian norm  
 $\phi$ : feature map

$w \in \mathbb{R}^d$   
 $b \in \mathbb{R}$   
 $x_i \in \mathbb{R}^d$   
 $y_i \in \{-1, 1\}$  : labeled data points

Once the soft margin SVM is learned, prediction of any data point  $x \in \mathbb{R}^d$  is given by the function:

$$f(x) = \text{sign}(w^T \phi(x) + b)$$

## → Linear Hyperplane classifier



- hyperplane  $y = \text{sgn}(w \cdot x + b)$  in canonical form if  $\min_{x_i \in X} |(w \cdot x_i) + b| = 1$ , i.e. scaling freedom removed.
- larger margin  $\sim 1/\|w\|$  is giving better generalization → LMC!

$$\max \|w\|^2 = \min \frac{1}{\|w\|^2}$$

Because we want new points to be correctly classified.

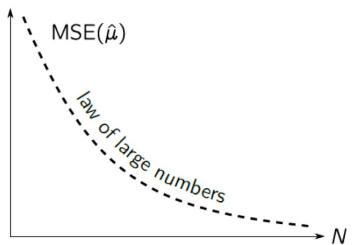
⇒ larger margin, smaller the norm, better generalization

Optimal separating hyperplane is the one that has the largest margin  $\hat{=} \text{smallest norm } \|w\|$

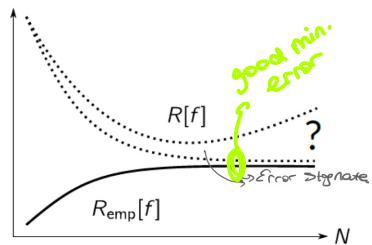
- the first datapoints we hit has norm  $\pm 1$
- margin for canonical hyperplanes is:  $\frac{1}{\|w\|}$
- If we manage to separate training data with a large margin, then we'll do well on test set.

• make the norm  $\|w\|$  small to make the margin large

Estimating a parameter



Learning a ML model



Questions:

- Given a dataset of  $N$  points, how well will my model generalize?
- How quickly  $R[f]$  and  $R_{\text{emp}}[f]$  converge with  $N$ ?  $1/N$  OR  $1/\sqrt{N}$
- Does it actually converge? there's a limit

- We don't want to overparametrize because it might lead to overfit
- We don't know the prob. distribution of our data, so we don't have the true risk  $R[f]$  instead, we have an empirical risk minimization
- We want a small generalisation (test) error

$$R[f] \leq R_{\text{emp}} + \text{const.} \quad \text{upper bound } R[f]$$

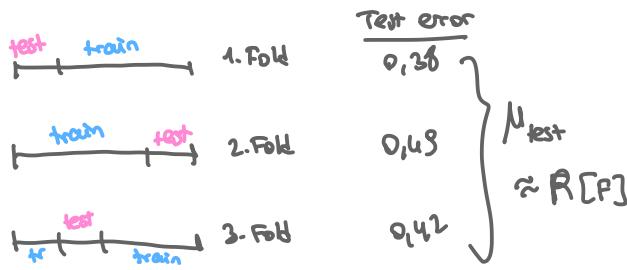
If we have our term, then we add the VC constant and say that our true risk must be  $\leq$  to that.

→ To predict  $R[f]$

Cross validation

dataset into train test

k-Fold CV:



Error Bound upper bound to  $R[f]$

## Empirical Risk Minimization (ERM)

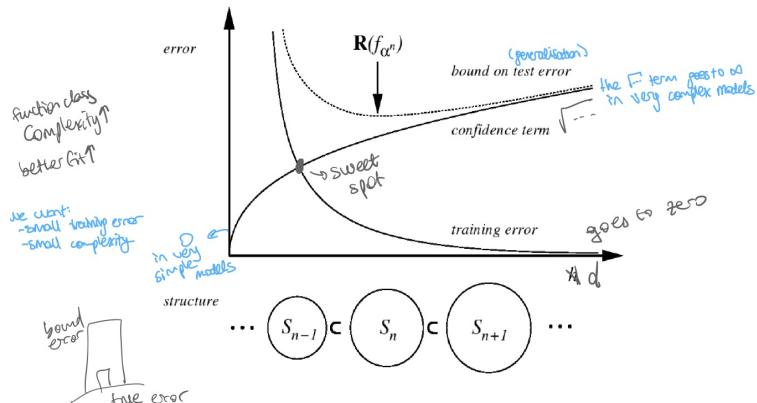
super complex function class:  
can fit anything, emp. risk  
will be zero, because the  
complexity of the function  
is enough to predict it  
with 0 error.  
VC-Dimension however is very large

Law of Large Numbers  
 $R_{\text{emp}} \rightarrow R_f$  as  $N \rightarrow \infty$

Empirical risk will converge to true risk, but it won't be the same, we need uniform convergence

- $d = \text{VC dimension}$  for class complexity
- try to minimize the VC-Dimension to get low complexity

### Structural Risk Minimization: the picture



Learning  $f$  requires small training error and small complexity of the set  $\{f_\alpha\}$ .

Dimension ↑ :- Training error ↓  
- Bound on test error ↑

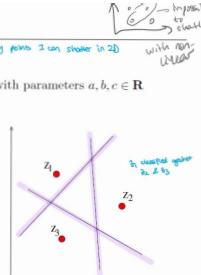
find the sweet spot

### VC Dimensions: an example

Half-spaces in  $\mathbb{R}^2$ :

$$f(x, y) = \text{sgn}(a + bx + cy), \quad \text{with parameters } a, b, c \in \mathbb{R}$$

- Clearly, we can shatter three non-collinear points.
- But we can never shatter four points.
- Hence the VC dimension is  $d = 3$
- In  $n$  dimensions: VC dimension is  $d = n + 1$



# VC Theory

## VC Theory applied to hyperplane classifiers

- Theorem (Vapnik 95): For hyperplanes in canonical form VC-dimension satisfying

$$d \leq \min\{[R^2 \|\mathbf{w}\|^2] + 1, n + 1\}.$$

w<sup>2</sup> = *w*<sup>2</sup>

Here,  $R$  is the radius of the smallest sphere containing data.  
Use  $d$  in SRM bound

$$R[f] \leq R_{emp}[f] + \sqrt{\frac{d(\log \frac{2N}{d} + 1) - \log(\eta/4)}{N}}.$$

- maximal margin = minimum  $\|\mathbf{w}\|^2 \rightarrow$  good generalization, i.e. low risk, i.e. optimize  $\min \|\mathbf{w}\|^2$
- independent of the dimensionality of the space!

Three scenarios: regression, classification & density estimation.  
Learn  $f$  from examples

$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbf{R}^N \times \mathbf{R}^M$  or  $\{\pm 1\}$ , generated from  $P(\mathbf{x}, y)$ , such that expected number of errors on test set (drawn from  $P(\mathbf{x}, y)$ ),

$$R[f] = \int \frac{1}{2} |f(\mathbf{x}) - y|^2 dP(\mathbf{x}, y),$$

is minimal (*Risk Minimization (RM)*).

**Problem:**  $P$  is unknown.  $\rightarrow$  need an *induction principle*.

*Empirical risk minimization (ERM):* replace the average over  $P(\mathbf{x}, y)$  by an average over the training sample, i.e. **minimize the training error**

$$R_{emp}[f] = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |f(\mathbf{x}_i) - y_i|^2$$

$\rightarrow$  With SLR we can upperbound the test error by the empirical error

$\rightarrow$  If hypothesis very complex, then we'll overfit.

$\rightarrow$  SVM's measure the capacity of the algorithm by their VC-Dimensions

### VC dimension: example

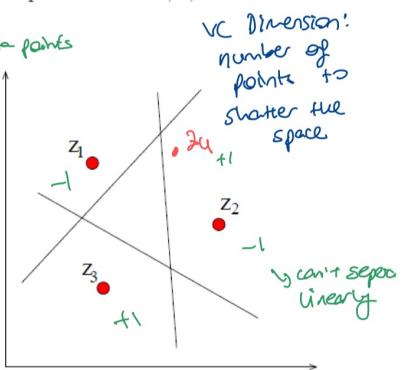
Half-spaces in  $\mathbf{R}^2$ : *if VC = 3 in 2D*  
*for 2 class problems: VC d = 3 = n+1*

$$f(x, y) = \text{sgn}(a + bx + cy), \quad \text{with parameters } a, b, c \in \mathbf{R}$$

*linear separation of data points*

- Clearly, we can shatter three non-collinear points.
- But we can never shatter four points.
- Hence the VC dimension is  $d = 3$

*Linear* • in  $n$  dimensions: VC dimension is  $d = n + 1$



$\rightarrow$  The complexity of a linear classifier is directly related to the dimensionality of the input space. Adding more features (increasing dimensionality) without sufficient data may however lead to overfitting.

### VC Dimension of a large margin



$\rightarrow$  If you can find  $\phi(\cdot)$  in higher dimensions that has:

- small training error
- low complexity,

then you can guarantee good generalization

complexity  $>$  dimensionality

## → Support Vector Approach

$$\Phi : \mathbb{R}^N \rightarrow F$$

$$x \mapsto \Phi(x)$$

where  $N \ll \dim(F)$ .

to get data  $(\Phi(x_1), y_1), \dots, (\Phi(x_N), y_N) \in F \times \mathbb{R}^M$  or  $\{\pm 1\}$ .

Learn  $\tilde{f}$  to construct  $f = \tilde{f} \circ \Phi$

→ preprocesses the data to a higher feature space

- SV-Learning: (in some cases) simpler: find  $f$ 's that make it simpler

If  $\Phi$  is chosen such that  $\{\tilde{f}\}$  allows small training error and has low complexity, then we can guarantee good generalization.

⇒ The complexity matters, not the dimensionality of the space.

## → Learning Theory

### minimal Risk Minimization (RM)

Three scenarios: regression, classification & density estimation.

Learn  $f$  from examples

$(x_1, y_1), \dots, (x_N, y_N) \in \mathbb{R}^n \times \mathbb{R}^m$  or  $\{\pm 1\}$ , generated from  $P(x, y)$ ,

such that expected number of errors on test set (drawn from  $P(x, y)$ ),

Risk: 
$$R[f] = \int \frac{1}{2} |f(x) - y|^2 dP(x, y),$$

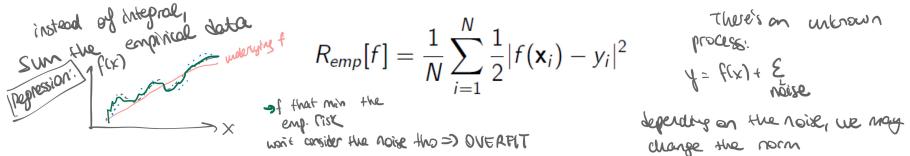
is minimal (Risk Minimization (RM)).

to get a low generalisation error

Problem: we don't know the prob. distribution  
so we look at the training error instead of generalisation error.

### empirical Risk Minimization (ERM)

Empirical risk minimization (ERM): replace the average over  $P(x, y)$  by an average over the training sample, i.e. minimize the training error



$$R_{\text{emp}}[f] = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} |f(x_i) - y_i|^2$$

There's an unknown process:  
 $y = f(x) + \varepsilon$   
noise

depending on the noise, we may change the norm

# SVM Primal

The classifier with largest margin between the positive and negative data points can be obtained by solving the minimization problem:

$$\min_{w,b} \frac{1}{2} \|w\|^2$$

$\rightarrow$  determines how flat the decision function / spaced is

subject to the constraints

$$\forall_{i=1}^N : y_i \cdot (w^\top \phi(x_i) + b) \geq 1. \quad \text{hard margin SVM}$$

Once the parameters have been optimized, the decision function is given by

$$\text{sign}(w^\top \phi(x) + b)$$

and it can be applied to classify new data.

How to optimize the decision? Lagrange had  $\hat{=}$ , but now we have  $\hat{\geq}$   
 $\Rightarrow$  Slater's condition

## Slater's Condition

data linearly separable

Consider the optimization problem

$$\min_{\theta} f(\theta)$$

$x \dots \overset{\circ}{\theta} \dots w$  you can find a  
 $\dots 0 \dots$  margin to separate

with  $f$  convex, subject to the inequality constraints

$$\forall_{i=1}^m : g_i(\theta) \leq 0$$

with  $g_i$  also convex. If there exists a point  $\theta$  in the input domain that satisfies all constraints with strict inequalities, then the solution of the optimization problem is also given by its Lagrange dual formulation, the constrained optimization problem:

$$\max_{\alpha} \left\{ \min_{\theta} \left\{ f(\theta) + \sum_{i=1}^m \alpha_i g_i(\theta) \right\} \right\} \quad \forall_{i=1}^m \alpha_i \geq 0, \quad (1)$$

where the inner minimization problem is typically solved analytically.

$\rightarrow$  1st min the Lagrange, then max over Lagrange variable  $\alpha$  where  $\alpha_i \geq 0$

$\rightarrow$  when  $g_i$  positive: constraint is not satisfied  $g_i(\theta) \leq 0$   
inner optimization

so if  $g_i(\theta) < 0$  make very large with  $\alpha$  and make it go to  $\infty$  solution,

$\rightarrow$  Unlike Lagrange, this is not a closed form but a on  $\alpha$  based dual variable solution.

## How do we apply Slater to SVM?

### SVM: from the Primal to the Dual

hard-margin

Slater's condition for SVM:

$$\exists w, b \text{ s.t. } \forall_{i=1}^N : y_i \cdot (w^\top \phi(x_i) + b) > 1$$

In other words, the two classes must be linearly separable. When  $\phi$  can be computed explicitly, this can be tested by running e.g. the perceptron algorithm.

If the Slater's conditions are verified, we inject the SVM objective in Eq. (1):

$$\max_{\alpha} \left\{ \min_{w,b} \left\{ \frac{1}{2} \|w\|^2 + \sum_{i=1}^N \alpha_i \cdot (1 - y_i \cdot (w^\top \phi(x_i) + b)) \right\} \right\} \quad \text{s.t. } \alpha \succeq 0,$$

$f(w,b) = f(\theta)$

$g(w,b) \leq 0$

$\rightarrow$  if the perceptron of the kernel converges to 0, then classes linearly separable

$\rightarrow$  or use soft margin because it doesn't request linear separability.

Roadmap for finding the dual in hard margin SVM:

- you'll find an explicit solution for  $w$ , but not for  $\alpha$ .
- for  $\alpha$ , you'll get a constraint that relates to dual variables.
- reinject  $w, b$  and get  $\alpha$

## Predicting with the Dual Solution

you get  $w$  by directly minimizing the inner optimisation problem:

Recall: we predict based on output of:

$$w^T \phi(x) + b$$

Question: how do we find the bias  $b$ ?

- ANSWER: We recall that because the model has maximized the margin between the two classes, the nearest point to the decision boundary for each class is exactly on the margin, i.e.

$$\min_i \{y_i \cdot (w^T \phi(x_i) + b)\} = 1$$

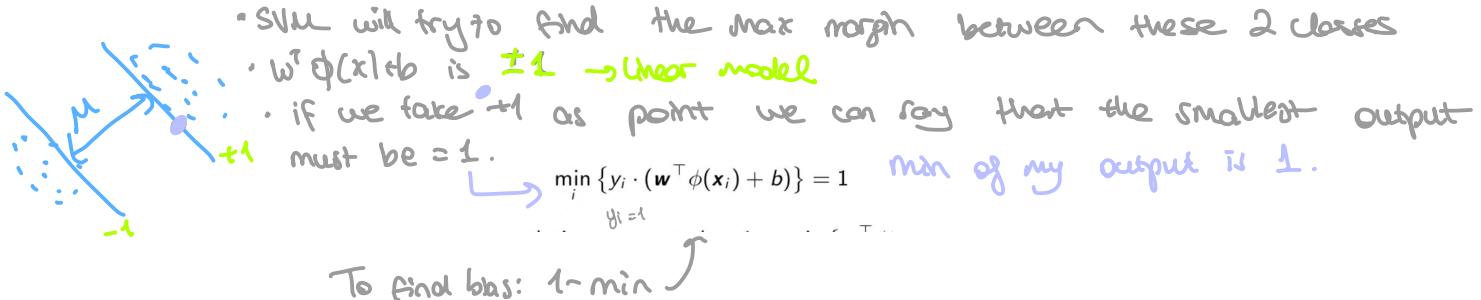
$y_i = 1$

After some manipulations, we get:  $b = 1 - \min_{i|y_i=1} \{w^T \phi(x_i) + b\}$ .

$$\max_{\alpha} \left\{ \min_{w,b} \left\{ \underbrace{\frac{1}{2} \|w\|^2}_{f(w,b)} + \sum_{i=1}^N \alpha_i \cdot (1 - y_i \cdot (w^T \phi(x_i) + b)) \right\} \right\} \quad \text{s.t. } \alpha \geq 0,$$

$f(w,b) = f(\alpha)$

$$w = \sum_i \alpha_i y_i \phi(x_i)$$



## Difference hard-Soft Margin SVM!

- there might be outliers in my class and the margin will not coincide with lowest output  $w^T \phi(x) + b$ , but it will be sth. else
- hard margin wants the output to be  $\pm 1$ , but not the case in soft margin

Question: how do we find the bias  $b$  for the soft-margin SVM?

- ANSWER: Bias can be recovered from the dual using another set of conditions, the KKT conditions.

### KKT Conditions (simplified)

Consider the optimization problem  $\min f(\theta)$ . Subject to the inequality constraints  $\forall i=1 \dots m : g_i(\theta) \leq 0$ . For this problem, the Karush-Kuhn-Tucker (KKT) conditions are:

$$\begin{aligned} \nabla f(\theta) + \sum_{i=1}^m \lambda_i \nabla g_i(\theta) &= 0 && \text{(stationarity)} \\ \forall i=1 \dots m : g_i(\theta) &\leq 0 && \text{(primal feasibility)} \\ \forall i=1 \dots m : \lambda_i &\geq 0 && \text{(dual feasibility)} \\ \lambda_i g_i(\theta) &= 0 && \text{(complementary slackness)} \end{aligned}$$

→ gradient of Lagrangian  
already given

→ useful to find bias  
primal constraint = 0  
soft SVM

- If Slater's conditions are satisfied, then the KKT conditions are satisfied at the optimum.

- The bias of the primal can be recovered from the 'complementary slackness' equation.

$$w^T \phi(x) + b - 1 \text{ is our constraint}$$

$$\alpha_i \cdot (1 - w^T \phi(x) - b) = 0$$

one of these is 0

if  $\alpha_i = 0$ : the  $(1 - w^T \phi(x) - b)$  can be anything  
if  $\alpha_i \neq 0$ : this is zero

# Lecture 6: Kernels

## Exercise Sheet 6

A kernel function  $k: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  must satisfy the *Mercer's condition*, which verifies that for any sequence of data points  $x_1, \dots, x_n \in \mathbb{R}^d$  and coefficients  $c_1, \dots, c_n \in \mathbb{R}$  the inequality

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0$$

is satisfied. If it is the case, the kernel is called a *Mercer kernel*.

Conversely, the *representer theorem* states that if  $k$  is a Mercer kernel on  $\mathbb{R}^d$ , then there exists a Hilbert space (i.e., a finite or infinite dimensional  $\mathbb{R}$ -vector space with norm and scalar product)  $\mathcal{F}$ , the so-called feature space, and a continuous map  $\varphi: \mathbb{R}^d \rightarrow \mathcal{F}$ , such that

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}} \quad \text{for all } x, x' \in \mathbb{R}^d.$$

### Kernel Trick:

Take a polynomial basis to increase the dimension in order to map the data to a higher dimension so that you get a simpler model.

## → Kernel Trick

2 polynomials

$$\begin{aligned} (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})) &= (x_1^2, \sqrt{2}x_1x_2, x_2^2)(y_1^2, \sqrt{2}y_1y_2, y_2^2)^T \\ &= (\mathbf{x} \cdot \mathbf{y})^2 \text{ requires as a scalar product in the original space} \\ &= : k(\mathbf{x}, \mathbf{y}) \text{ kernel} \end{aligned}$$

- Scalar product in (high dimensional) feature space can be computed in  $\mathbb{R}^2$ !

- works only for Mercer Kernels  $k(\mathbf{x}, \mathbf{y})$

Scalar product to not compute in  $\mathbb{R}^d$  dimension space  $\xrightarrow{\text{W} \cdot \mathbf{x}}$  scalar product based linear learning model

• take 2 polynomials

• compute their scalar product & map into a function

• map data to higher dimension to solve with a simpler model.

of kernel ≈ regularizer: if you want your function to have a certain property, then map it first to a kernel & use as a SV.

## → Mercer Kernels

[Mercer] If  $k$  is a continuous kernel of a positive integral operator on  $L_2(\mathcal{D})$  (where  $\mathcal{D}$  is some compact space),

*If this is a positive operator, then we can always use kernel trick*

$$\int f(\mathbf{x}) k(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{x} d\mathbf{y} \geq 0, \quad \text{for } f \neq 0$$

it can be expanded as

$$k(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^{N_F} \lambda_i \psi_i(\mathbf{x}) \psi_i(\mathbf{y}) \quad \text{hidden scalar product}$$

with  $\lambda_i > 0$ , and  $N_F \in \mathbb{N}$  or  $N_F = \infty$ . In that case

$$\Phi(\mathbf{x}) := \begin{pmatrix} \sqrt{\lambda_1} \psi_1(\mathbf{x}) \\ \sqrt{\lambda_2} \psi_2(\mathbf{x}) \\ \vdots \end{pmatrix} \quad \xrightarrow{\mathbf{x} \in \mathcal{D}}$$

satisfies  $(\Phi(\mathbf{x}) \cdot \Phi(\mathbf{y})) = k(\mathbf{x}, \mathbf{y})$ .

① Stat. learning theory  
 ② risk minimisation: hard to find  
 ③ ERM: min the training error  
 ↗ not a good idea! overfit

④ SRM: bounds the pen. error from above

- small errors on training set } occurs rather  
 - simple

⑤ large margin classifier  
 - map to high-dim. in non-linear: optimal  
 - non-linear mapping: optimal generalising

## The Kernel Trick in Practice

$$\|\text{function} - \mu\|$$

### Example:

We have a dataset  $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and would like to compute the (squared) distance to the mean in feature space, i.e.

$$f(\mathbf{x}) = \left\| \underbrace{\phi(\mathbf{x})}_{\text{high-dim}} - \frac{1}{N} \sum_{i=1}^N \underbrace{\phi(\mathbf{x}_i)}_{\text{high-dim}} \right\|^2 \quad \begin{array}{l} \text{→ we need to take scalarproduct} \\ \text{of each function}-\mu \text{ terms} \end{array}$$

Applying the identity

$$\|\mathbf{a} - \mathbf{b}\|^2 = \langle \mathbf{a}, \mathbf{a} \rangle - 2\langle \mathbf{a}, \mathbf{b} \rangle + \langle \mathbf{b}, \mathbf{b} \rangle$$

we arrive at a formulation

$$f(\mathbf{x}) = \underbrace{\langle \phi(\mathbf{x}), \phi(\mathbf{x}) \rangle}_{k(\mathbf{x}, \mathbf{x})} - \frac{2}{N} \sum_{i=1}^N \underbrace{\langle \phi(\mathbf{x}), \phi(\mathbf{x}_i) \rangle}_{k(\mathbf{x}, \mathbf{x}_i)} + \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \underbrace{\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle}_{k(\mathbf{x}_i, \mathbf{x}_j)}$$

where the function (initially defined in feature space) can be expressed solely in terms of the kernel function. *Scalar product in high-dim. is gone with this trick* to avoid expensive computations.

Lagrangian  $L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y_i \langle x_i, w \rangle + b - 1]$   $\alpha_i \geq 0$

• Lagrangian must be :- maximized w.r.t.  $\alpha$   
- minimized w.r.t.  $w$  and  $b$

• Derivatives w.r.t. primal variables must vanish:  $\frac{\partial L(w, b, \alpha)}{\partial b} = 0$ ,  $\frac{\partial L(w, b, \alpha)}{\partial w} = 0$

which leads to:  $\sum_{i=1}^m \alpha_i y_i = 0$  and  $w = \sum_{i=1}^m \alpha_i y_i x_i$

• According to KKT, only the Lagrange multipliers  $\alpha_i > 0$  at the saddle point correspond to constraints.  
 $\alpha_i [y_i \langle x_i, w \rangle + b - 1] = 0$

→ The patterns  $x_i$  for which  $\alpha_i > 0$  are Support Vectors and lie exactly on the margin.

Gaussian Gram Matrix (K) each element  $K_{ij}$  represents the Gaussian kernel between observations  $x_i$  and  $x_j$ .

$$K_{ij} = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

$\sigma$ : Bandwidth of Gaussian kernel.

Properties Gram Matrix is:

- symmetric  $K_{ij} = K_{ji}$
- PSD  $u \in \mathbb{R}^N$ ,  $u^T K u \geq 0$

### Mercer kernels (most common)

Polynomial:  $k(x, y) = (x \cdot y + c)^d$

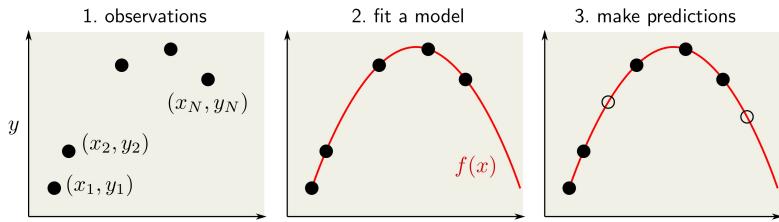
RBF:  $k(x, y) = \exp(-\|x - y\|^2 / 2\sigma^2)$

Inverse quadratic:  $k(x, y) = \frac{1}{\sqrt{\|x - y\|^2 + c^2}}$

# Lecture 7: Model Selection

→ We want to learn a model  $f(x)$  that not only predicts well, but also predicts future unseen data well.

Occam's Razor  
Poppers prediction strength  
Cross-Validation  
Bias-Variance Analysis

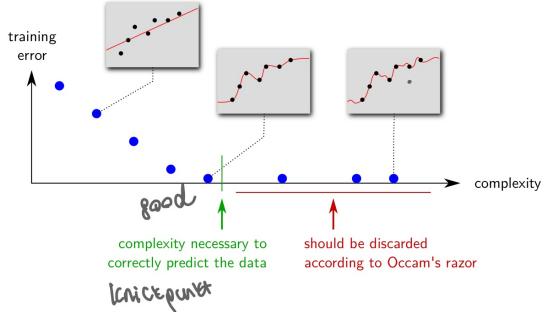


the model needs to understand the „essence“  
· perfect fit / simple / too complex is not always the best. → overfit or underfit problems

“Model complexity must not be increased beyond what is necessary to correctly predict the data.”

## → Occam's Razor

Occam's Razor for Model Selection



find the sweet spot for complexity to reduce over/underfit

choose the model that is the least complex one,  
it will: predict new (unseen) data better  
& simplicity is better

Because: generalization error bound is tighter!

## Quantifying Complexity:

### Approach 1: Counting the Parameters

Idea:

- If several models predict the data sufficiently well, prefer the one with the fewest parameters.

#### Examples of models:

Constant classifier	$g(\mathbf{x}) = \frac{C}{1}$	1 parameter
Means difference	$g(\mathbf{x}) = \mathbf{x}^T \left( \frac{\mathbf{m}_1 - \mathbf{m}_2}{2d} \right) + \frac{C}{1}$	close to NCC
PCA + Fisher	$g(\mathbf{x}) = \frac{\text{PCA}(\mathbf{x})^T}{k \cdot d} \frac{S_{W^{-1}}^{-1}}{k^2} \frac{(\mathbf{m}_1 - \mathbf{m}_2)}{2k} + \frac{C}{1}$	$k = \# \text{PCs that we project down}$

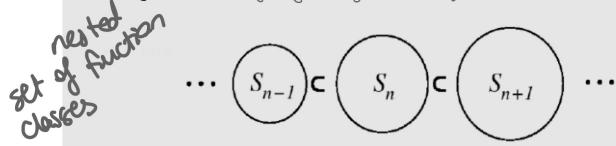
- counting model parameters
- smoothness
- SRM
- Bayesian information criterion (BIC)

## SRM

if we have the nesting, then we can start doing some generalisation bounds

#### SRM Idea:

- Structure the space of solutions into a nesting of increasingly large regions. Benefit of nesting: Bounce on generalization errors



- If two solutions fit the data, prefer the solution that also belongs to the smaller regions.

Particular choices of  $\dots \subseteq S_{n-1} \subseteq S_n \subseteq S_{n+1} \subseteq \dots$  lead to upper-bounds on the generalization error of a model (cf. next lecture).

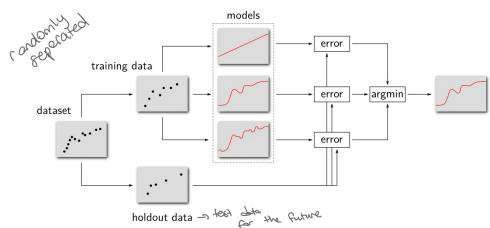
generalisation  $\approx$  training data error

→ the model with lowest generalization error is preferable.

All the points are coming from some unknown distributions. How can we know that our constructed model is good on unseen data?

## The Holdout Selection Procedure

Idea: Predict out-of-sample error by splitting the data randomly in two parts (one for training, and one for estimating the error of the model).

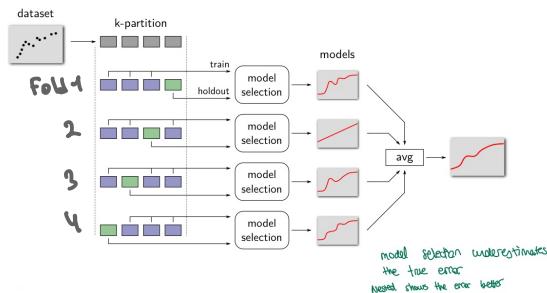


Problem: The more data we use to accurately estimate the prediction error of each model, the less data is available for training an accurate model in the first place.

## Cross-Validation ( $k$ -Fold Procedure)

splitting data to k-parts < training hold out

Idea: Retain more data for training, and make up for the lower amount of holdout data by repeating the model selection procedure over multiple data splits (and averaging the selected models).



## Holdout / Cross-Validation

**Advantages:** mimique the unseen

- The model can now be selected directly based on simulated future observations (implements Popper's principle for model selection).

**Limitations:**

- For a small number of folds  $k$ , the training data is reduced significantly, which may lead to less accurate models. For  $k$  large, the procedure becomes computationally costly.
- This technique assumes that the available data is representative of the future observations (not always true!).
  - sample selection bias: data does not properly represent the underlying distribution
  - To solve: Heckman estimator

Heckman: he only did the survey on the phone  
→ Sample Selection bias

"mimique" the unknown

- holdout data: data we haven't touched before
- we split the data to 2 parts: test & holdout
- we measure the error and take the argmin
- Problem: the more we take from the dataset to holdout, the less accurate our estimation will be.
- easiest case: lots of data

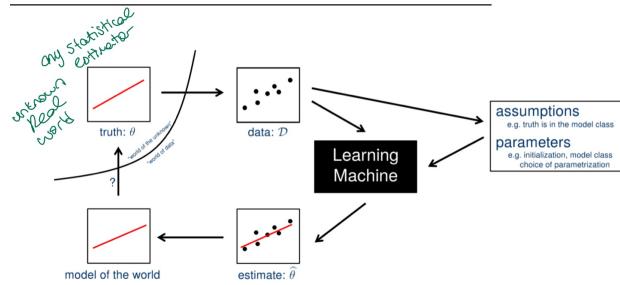
## Clever Hans effect:

Model gives the right classification for incorrect reasons: it understands the copyright as the classifier, not the horse itself.



Datasets may include nonsensical associations.

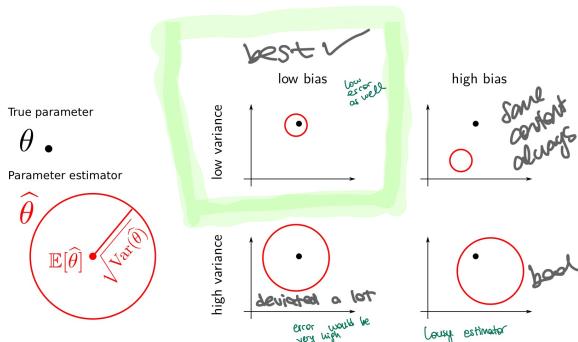
# → Bias - Variance Analysis



A good learning machine is one that produces an estimator  $\hat{\theta}$  close to the truth  $\theta$ . Closeness to the truth can be measured by some error function, e.g. the square error:

$$\text{Error}(\hat{\theta}) = (\theta - \hat{\theta})^2.$$

truth      estimator



MSE is Bias + Variance

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta} - \theta], \quad \text{Var}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2], \quad \text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

Exercise: Show that  $\text{MSE}(\hat{\theta}) = \text{Bias}(\hat{\theta})^2 + \text{Var}(\hat{\theta})$ .

MSE is our generalization error

$$\begin{aligned} \text{MSE}(\hat{\theta}) &= E[(\hat{\theta} - \theta)^2] \\ &= E[(\underbrace{\hat{\theta} - E[\hat{\theta}]}_{\text{bias}}) + \underbrace{(E[\hat{\theta}] - \theta)}_{\text{error}}]^2 \\ &= E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta}] - \theta)^2 + 2(\hat{\theta} - E[\hat{\theta}]) \cdot (E[\hat{\theta}] - \theta) \\ &= \underbrace{\text{Var}(\hat{\theta})}_{\text{low variance}} + \underbrace{\text{Bias}(\hat{\theta})^2}_{\text{high bias}} + \underbrace{0}_{\text{expectation } \hat{\theta} = \text{exp. } \theta} \end{aligned}$$

*Trick: add 0 by adding E[0]*

Parametric estimation:

►  $\theta$  is a value in  $\mathbb{R}^h$

►  $\hat{\theta}$  is a function of the data  $D = \{X_1, \dots, X_N\}$ , where  $X_i$  are random variables producing the data points.

Statistics of the estimator:

$$\text{Bias}(\hat{\theta}) = E[\hat{\theta} - \theta]$$

How far are we from the truth? (in average because dataset is sampled)

$$\text{Var}(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2]$$

How much does it deviate from the mean of overall estimates?

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

How far is our estimator from the truth?

(measures prediction error)

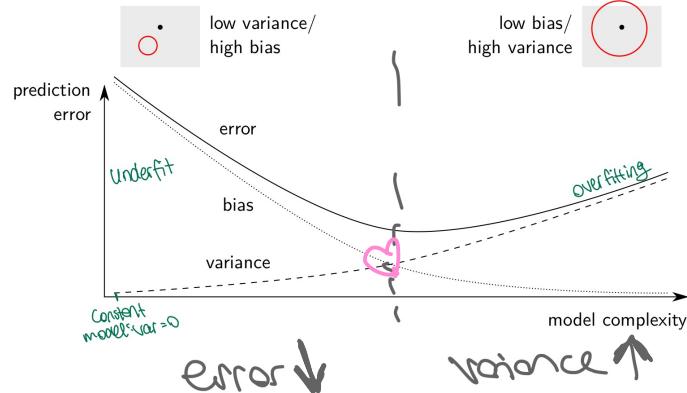
Note: for  $\theta \in \mathbb{R}^h$ , we use the notation  $\theta^2 = \theta^\top \theta$ .

product notation

How far is our model from the ground truth? - Bias - Variance will tell.

if we put too much complexity we overfit

a lot of data!  
⇒ low bias  
⇒ high variance



## → Parameters of a Gaussian:

parametric estimation:

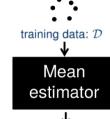
$\theta$  is a value in  $\mathbb{C}^n$  (e.g.  $\theta = (\mu, \Sigma)$  for Gaussians)

$\hat{\theta}$  is function in the data  $D = \{X_1, \dots, X_N\}$   
( $X_i$  are random variables giving back data points)

e.g. mean estimator  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$

covariance estimator  $\hat{\Sigma} = \frac{1}{N-1} (X_i - \hat{\mu})(X_i - \hat{\mu})^\top$

How the data varies around the mean?



insert  $\hat{\mu}$

Exercise: Show that the bias and variance of the mean estimator  $\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$  are given by

$$\text{Bias}(\hat{\mu}) = 0 \quad \text{and} \quad \text{Var}(\hat{\mu}) = \sigma^2/N.$$

our estimator:

$$\begin{aligned} \text{Bias}(\hat{\mu}) &= E[\hat{\mu} - \mu] \\ &= E[(\frac{1}{N} \sum_{i=1}^N X_i) - \mu] \\ &\stackrel{\text{definition of mean}}{=} (\frac{1}{N} \sum_{i=1}^N E[X_i]) - \mu \\ &= (\frac{1}{N} \sum_{i=1}^N \mu) - \mu \\ &= 0 \\ &\stackrel{\text{unbiased estimator}}{=} \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \text{Var}[\frac{1}{N} \sum_{i=1}^N X_i] \\ &= \frac{1}{N^2} \text{Var}[\sum_{i=1}^N X_i] \\ &= \frac{1}{N^2} \sum_{i=1}^N \text{Var}[X_i] \\ &= \frac{1}{N^2} \sum_{i=1}^N \sigma^2 \\ &= \frac{1}{N} \sigma^2 \\ &= \sigma^2/N \end{aligned}$$

unbiased estimator is good, because if we add more data it'll estimate more correctly.

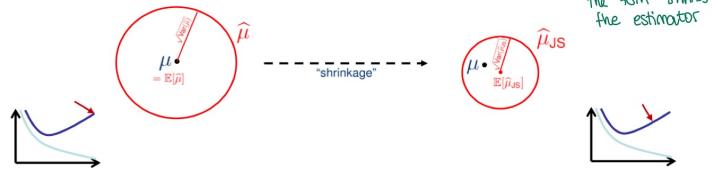
## The James-Stein Estimator

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N X_i$$

"natural" estimator

$\text{Bias}(\hat{\mu}) = 0$

$$\text{MSE}(\hat{\mu}) = \text{Var}(\hat{\mu}) = \frac{\sigma^2}{N}$$



"better than the mean"  
not unbiased anymore

$$\hat{\mu}_{\text{JS}} = \hat{\mu} - \frac{(n-2)\sigma^2}{\hat{\mu}^2} \hat{\mu}$$

James-Stein estimator

$\text{Bias}(\hat{\mu}_{\text{JS}}) > 0$

$$\text{MSE}(\hat{\mu}_{\text{JS}}) < \text{MSE}(\hat{\mu})$$

Variance is  
much  
smaller

add sth. to the diagonal  
also adds bit bias, but variance goes small

- ▶ **Occam's Razor:** Given two models that have sufficiently low training error, the simpler one should be preferred.  
*the simpler the better*
- ▶ **Measuring Complexity:** Counting the parameters can be misleading. Structured risk minimization (SRM) is more reliable in practice.
- ▶ **Popper's View:** How to make sure that a model predicts well? By testing it on out-of-sample data.  
*for few data points: leave-one out CV is a good strategy*
- ▶ **Holdout and Cross-Validation:** Common practical procedures to simulate out-of-sample prediction behavior. Has some limitations (e.g. assume that the current dataset is representative of the true distribution).
- ▶ **Bias-Variance Decomposition:** The error of a predictive model can be decomposed into bias and variance. Best models can often be interpreted as finding a good tradeoff between the two terms.

Simple models: ↓ low variance high bias ↑  
Complex models: ↑ high variance low bias ↓

## Biased Estimator

to make the estimator less affected by the sampling of data.

- A biased estimator can reduce the variance of the estimate at cost of introducing some bias.
- Introduce some bias to the model to reduce variance to get more generalized results. → Regularization: Prevent overfitting by penalizing overly complex models
- In some models, bias is deliberately introduced to achieve lower variance to get better generalisation on unseen data.
- Introducing bias to the model makes the model less sensitive (because variance is decreasing) by reducing the impact of random fluctuations in the sample.

Overfitted estimator has: high variance - low bias

High variance models are sensitive to small fluctuations in the training set, leading to overfitting.

They perform well on training data but do poor generalisation.

**Curse of dimensionality:** Number of features  $\gg$  observations

Introduce a bias (f.e. with regularization) to reduce estimation error

$$MSE = \text{bias}^2 + \text{Var}$$

## Biased Estimator

to make the estimator less affected by the sampling of data.

- A biased estimator can reduce the variance of the estimate at cost of introducing some bias.
- Introduce some bias to the model to reduce variance to get more generalized results. → Regularization: Prevent overfitting by penalizing overly complex models.
- In some models, bias is deliberately introduced to achieve lower variance to get better generalisation on unseen data.
- Introducing bias to the model makes the model less sensitive (because variance is decreasing) by reducing the impact of random fluctuations in the sample.

Overfitted estimator has: high variance - low bias

High variance models are sensitive to small fluctuations in the training set, leading to overfitting.

They perform well on training data but do poor generalisation.

**Curse of dimensionality:** Number of features  $\gg$  observations

Introduce a bias (f.e. with regularization) to reduce estimation error

$$MSE = \text{bias}^2 + \text{var}$$

# Lecture 8: Neural Networks 1

## Bayes Decision

What if we don't know the data?

parameter estimation

Theory optimal classifiers  
we know the data properties  
no better classifier

## Mean Separation / Fisher

Define classification based directly on data rather than learning data distribution

Perception  
focus on db separation of 2 classes

NN

Under + Non-Linear

Model Selection → Occams Razor: Simplicity and Generalization Evaluation  
With amount of data, can we still build a model that generalizes?

build a reduced function to feed the data; less overfit chance  
cross validation (test (holdout) data from one test set)  
training bundle done  
more accurate  
dige back  
sonra

## Perceptron (check Lecture 2)

### The Perceptron Algorithm

bias to adjust the location of the db



- Consider our linear model

$$z_k = \mathbf{w}^T \mathbf{x}_k + b \quad y_k = \text{sign}(z_k)$$

and let  $t_k$  be 1 and -1 when the true class of  $\mathbf{x}_k$  is  $\omega_1$  and  $\omega_2$  respectively.

#### Algorithm

- Iterate over  $k = 1 \dots, N$  (multiple times).
    - If  $\mathbf{x}_k$  is correctly classified ( $y_k = t_k$ ), continue.
    - If  $\mathbf{x}_k$  is wrongly classified ( $y_k \neq t_k$ ), apply:  
update rules:  
 $\mathbf{w} \leftarrow \mathbf{w} + \eta \cdot \mathbf{x}_k t_k$   
 $b \leftarrow b + \eta \cdot t_k$
- where  $\eta$  is a learning rate.

- Stop once all examples are correctly classified.

$x_1, t_1$   
 $x_2, t_2$   
 $\vdots$   
 $x_N, t_N$

1st pass  
2nd pass  
until all data points classified

→ data linear separable learn will quickly converge, else will not converge  
→ all other cases will go slow  
→ if algorithm does not converge: not linearly separable

### The Perceptron: Optimization View

The perceptron can be seen as a gradient descent of the error function

$$\mathcal{E}(\mathbf{w}, b) = \frac{1}{N} \sum_{k=1}^N \max(0, -z_k t_k)$$

GD: gradient descent rule  
decides which direction to move



Gradient:

$$\frac{\partial \mathcal{E}_k}{\partial \mathbf{w}} = \underbrace{1}_{\text{indicator}} \underbrace{\{ -z_k t_k > 0 \}}_{\text{not correctly classified}} \cdot (-x_k \cdot t_k)$$

$$\mathbf{w} = \mathbf{w} - \mu \cdot \frac{\partial \mathcal{E}_k}{\partial \mathbf{w}} = \mathbf{w} + \mu \cdot x_k t_k$$

where I have  
incorrectly  
classified  
the data

## Non-linear Classification

### Key Idea: → Kernel Trick

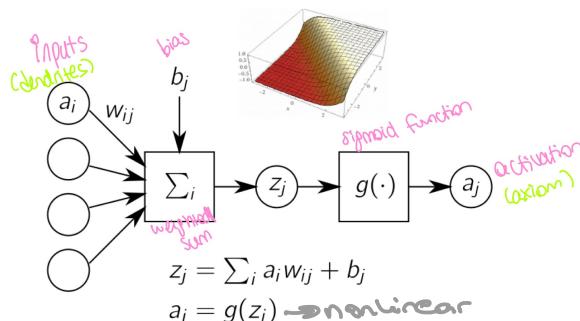
- Transform the data nonlinearly through some function  $\Phi$  before applying the linear decision function.

$$f(\mathbf{x}) = \mathbf{w}^T \Phi(\mathbf{x}) + b$$

- Example:  $\Phi(\mathbf{x}) = [x_1, x_2, x_1^2, x_2^2, x_1 x_2]$  and  $\mathbf{w} \in \mathbb{R}^5$ .

(first time talking about non-linearity in this course)

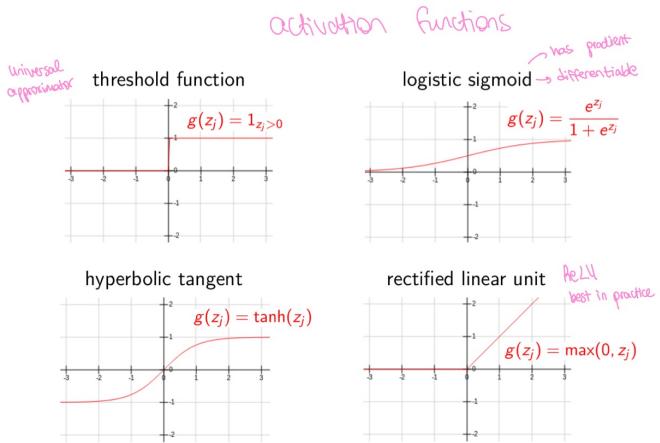
### artificial neuron:



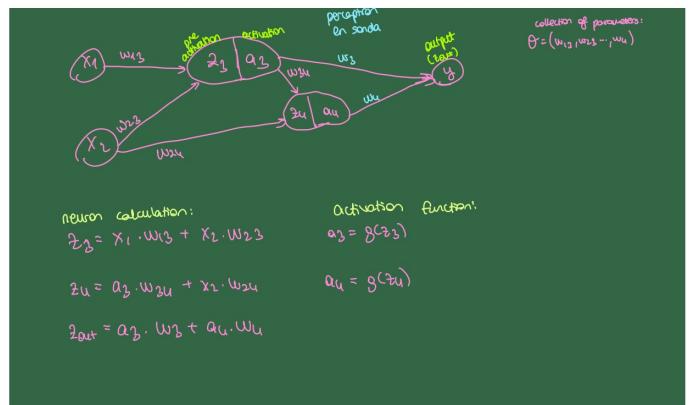
universal approx.  
theorem

- simple multivariate, non linear, differentiable
- complex nonlinear representations with neurons
- ability to learn from the data

## Examples of Nonlinear Functions



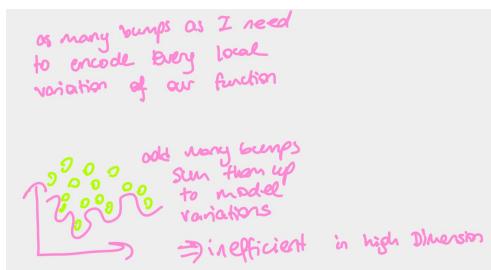
## The Forward Pass



## Universal Approx. Theorem

"With sufficiently many neurons, NN can approximate any nonlinear function."

add the bumps together, each iteration has a hidden layer



Step function:

- move hyperplane by the bias
- Scale it with w
- then add the step functions of neurons

approximation with infinite Fourier:

$$y(x_1, x_2) \approx \sum_s A_s(x_1) \cos(sx_2)$$

## Training a NN

- Use the same error function as the perceptron, but replace the perceptron output  $z$  by the neural network output  $z_{\text{out}}$ :

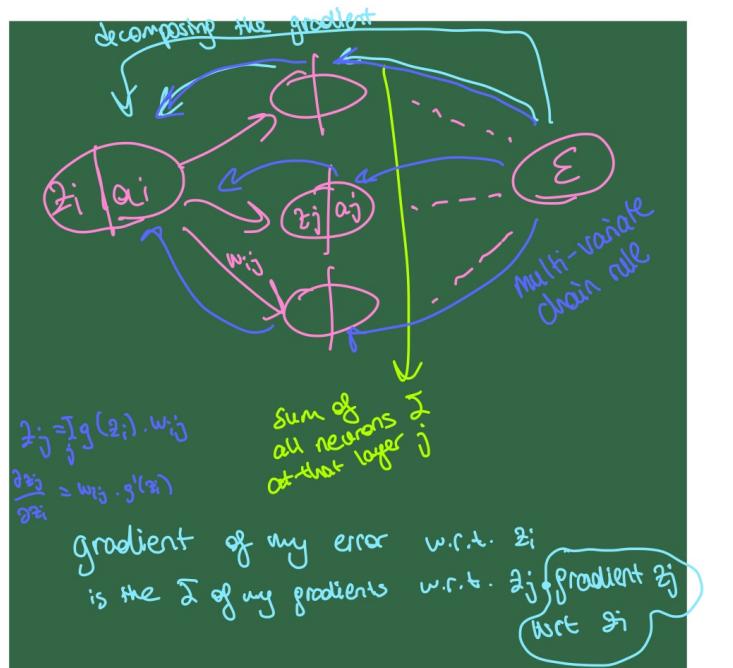
$$\mathcal{E}(\theta) = \frac{1}{N} \sum_{k=1}^N \underbrace{\max(0, -z_{\text{out}}^{(k)} t^{(k)})}_{\mathcal{E}^{(k)}(\theta)}$$

and compute the gradient of the error function w.r.t. the parameters  $\theta$  of the neural network.

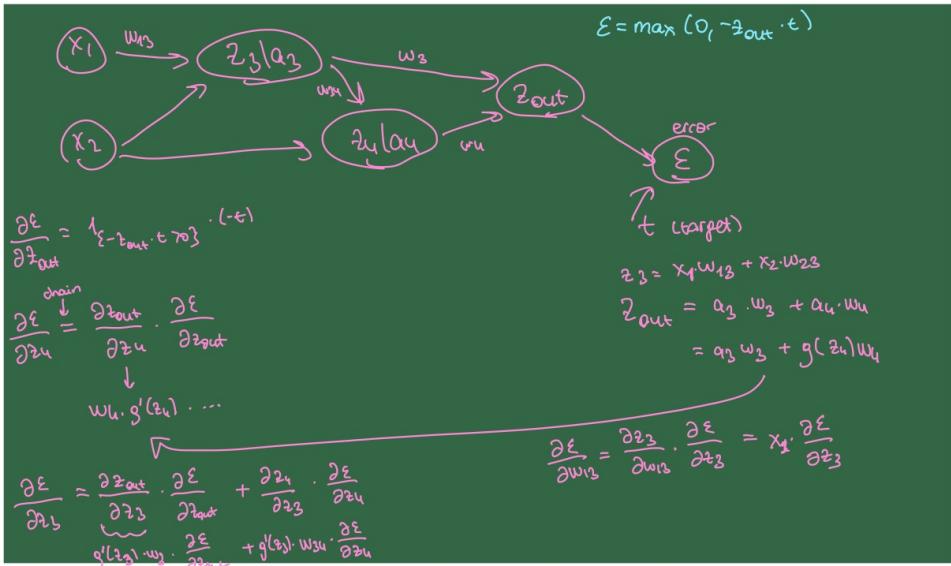
Gradient with chain rule:  $\frac{\partial \mathcal{E}}{\partial \theta_i} = \sum_j \frac{\partial z_i}{\partial z_j} \frac{\partial \mathcal{E}}{\partial z_j}$  for layers

for parameters:  $\frac{\partial \mathcal{E}}{\partial w_{ij}} = \frac{\partial z_j}{\partial w_{ij}} \cdot \frac{\partial \mathcal{E}}{\partial z_j}$

## Error Backpropagation to find Gradients



# Error Backpropagation



## Summary

- The **perceptron** and the **neural network** enable training classifiers on more complex distributions by focusing on what is critical for classification, i.e. the boundary between classes.
- The **neural network** enables learning **nonlinear** decision boundaries. This is useful when the problem is complex (most practical problems are nonlinear).
- The gradient of a **neural network** required for learning can be easily and quickly computed using the method of **error backpropagation**.
- The perceptron and the neural network do not have closed form solutions but can be trained iteratively using **gradient descent**.

► NN can learn any function, even tough not the best solution  
 → the from-NN-learned feature space can have multi-scale properties:  
 it doesn't have to be a reproducing kernel H.  
 → 1-1 mapping of NN & SVD possible.  
 choose kernel  
 ↳ representation  
 ↳ NN learns the representation

How to do backups pagination?

It will ask you  $\frac{\partial \text{output}}{\partial \text{weight of one input neuron}}$

Go to that input neuron and follow its arrows to following neurons.

Example: We need to find..  $\frac{dy}{dx}$

Wrs goes to all the blue paths

20 separate into 2 cases:  $\frac{\partial y}{\partial x_3} = \frac{\partial y}{\partial x_6} + \frac{\partial y}{\partial x_5}$  } then calculate 2 paths chain rule

$$\frac{\partial y}{\partial x_6} = \frac{\partial y}{\partial a_6} \cdot \frac{\partial a_6}{\partial z_6} \cdot \frac{\partial z_6}{\partial a_3} \cdot \frac{\partial a_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial a_1}$$

each activation function must be derivated from their  $\pi$ :

$$\frac{\partial y}{\partial x_5} = \frac{\partial y}{\partial x_5} \cdot \frac{\partial x_5}{\partial z_5} \cdot \frac{\partial z_5}{\partial x_3} \cdot \frac{\partial x_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_{11}}$$

How to find bias and w of a NN?

- give points : 1 for class A    1 for class B

$a_j = g(x \cdot w_{ij} + b_j)$  → needs to be 0 in critical points

# Lecture 9: Neural Networks 2

→ Optimizing NN: training efficiently & effectively  
→ Regularizing NN: learned model is good & generalizes well margin loss & perturbations

## optimizing the NN

### ① Initialization

weights should be initialized randomly

- Weights should be initialized *randomly* (it breaks symmetries in parameter space and reduces the risk of landing in bad local minima or getting stuck on plateaus).
- If necessary, train the network *multiple times* using different random seeds, and retain the network with the *lowest error*.

#### He-et-al<sup>1</sup> random initialization heuristic:

For neural networks with ReLU neurons, it is recommended to use:  
 $w_{ij} \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma = \sqrt{2/\#}$  input connections and where  $\mathcal{N}$  is the Gaussian distribution.

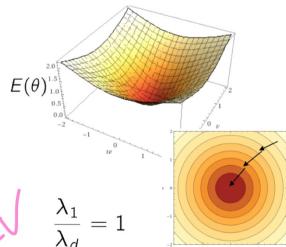
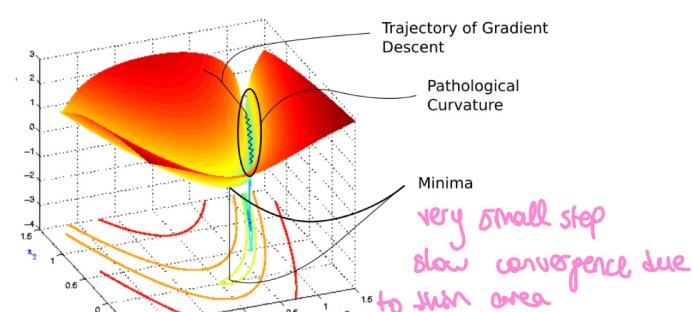
Example: For a two-layer neural network with 50 input features, 200 hidden neurons, and 1 output, weights in the first layer should be initialized using  $\sigma = \sqrt{2/50}$ , and in the second layer using  $\sigma = \sqrt{2/200}$ .

- Hessian describes how the minimum looks like
- Boulder & low conditioning number minimum is the best.
- Condition number can be reduced by centering the data before training.

- How hard is it to optimize a neural network? *pathological curvature waco*
- Are we guaranteed to converge to a good local minimum?
- How quickly do we converge to this local minimum?  
*low conditioning numbers*

### ② Neural Networks: Pathological Curvature

Another problem: The neural network may converge to some good local optimum but slowly due to *pathological curvature* of the error function.

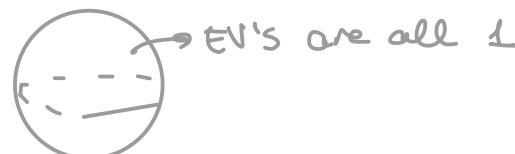


The lower the condition number, the better.

Very hard to optimize,  
therefore condition number  
very important

### ③ Hessian 2nd derivative

To understand the curvature of the error function  
→ results of Hessian = EV(?)



Condition number:  $\frac{\text{highest EV}}{\text{lowest EV}}$  ratio

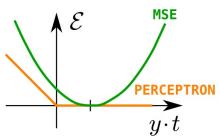
condition number ↑ (slower) the convergence ↓

⇒ lower condition number, the better.

## 4 Hessian Analysis (Linear Case, MSE)

↳ How the minimum looks like

Consider the simplest case of a homogeneous linear model  $y = \mathbf{w}^\top \mathbf{x}$ . Consider  $\mathcal{E}(\mathbf{w}) = \mathbb{E}[0.5 \cdot (1 - y \cdot t)^2]$ , the mean square error (MSE), which is easier to analyze than the perceptron loss.



The Hessian of the error function is given by:

$$H = \frac{\partial^2 \mathcal{E}}{\partial \mathbf{w}^2} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] \quad \text{data covariance}$$

i.e. the data uncentered covariance.

We center the data to get a low conditioning number:

**Condition Number:** Assuming  $\mathbf{x} \sim \mathcal{N}(\mu, \sigma^2 I)$ , the Hessian reduces to  $H = \sigma^2 I + \mu\mu^\top$ , and the condition number is

$$\frac{\lambda_1}{\lambda_d} = 1 + \frac{\|\mu\|^2}{\sigma^2}$$

(show this in the homework). In other words, the condition number can be reduced by **centering the data** before training.

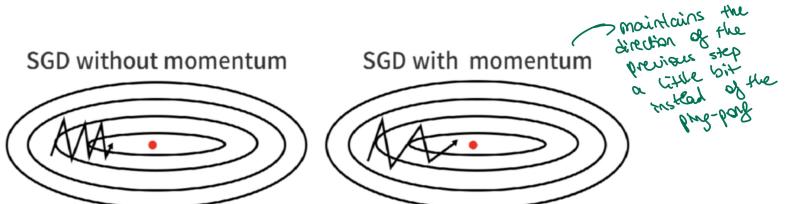
so that optimization is more smoother

19/3

→ While centering it is also ok to have a function that:  $g(z) = z$  (Urprung)  
suitable activation functions for that are:

ReLU:  $\max(0, z)$  ,  $\tanh(z)$  , centered softplus

→ Centering reduces the pathological curvature, but does not completely eliminate it.  
To make it faster: Gradient Descent with Momentum!



## 5 Improving Gradient Descent with Momentum:

**(Stochastic) Gradient Descent:**

$$\theta_{\text{new}} = \theta_{\text{old}} - \gamma \cdot \frac{\partial \mathcal{E}_k}{\partial \theta} \quad \text{with } k \sim \{1, \dots, N\}$$

where  $\gamma$  is the learning rate.

**(Stochastic) Gradient Descent + Momentum:**  $\approx$  filtering

$$\Delta_{\text{new}} = \mu \Delta_{\text{old}} + \frac{\partial \mathcal{E}_k}{\partial \theta} \quad \text{with } k \sim \{1, \dots, N\}$$

where  $\gamma$  is the learning rate and  $\mu$  is the momentum.

How to compute Hessian? Take 2nd derivative of the objective function:

$$H = \frac{\partial}{\partial w} \left( \frac{\partial J(w)}{\partial w} \right)$$

!  $\text{Cov}(x) = E(xx^T) - E(x)E(x)^T$

How to find the condition number?

$$\frac{\lambda_1}{\lambda_d} = 1 + \frac{\|M\|^2}{\delta^2}$$

1. Take 1st EV (highest):

$$\lambda_1 = \max_{\|v\|=1} v^T Hv$$

| plug in H  
|  $v^T v = 1$

$$v \perp \mu$$

If our  $H = \sigma^2 I + \mu \mu^T$ : our Eigenvalues are  $\sigma^2$  in the diagonal

# Regularizing the NN

that generalizes well to new data

## → Hinge Loss

adding a bit more to margin to penalize

So far:

- We have used the perceptron loss:

$$\mathcal{E}_k(\theta) = \max(0, -y_k t_k)$$

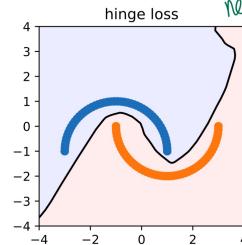
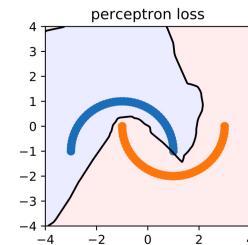
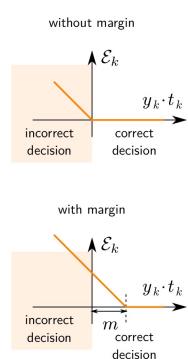
which becomes zero as soon as the current data point is being correctly classified.

Idea:

- Impose a penalty to points that lie too close to the decision boundary (i.e. add a margin), aka. the **Hinge Loss**:

$$\mathcal{E}_k(\theta) = \max(0, 1 - y_k t_k)$$

This is similar to the slack variables penalties in the soft-margin SVM formulation.



- Using the hinge loss, the decision function of the neural network moves away from the data.
- This leads to a better generalization performance.

## → How to regularize NN?

① weight decay add a term  $\gamma \|\theta\|^2$  to the objective, to make  $w$  smaller and closer to a linear model.

weight T non-linearity ↑  
try to stop early if weight  
to big, to remain in a  
linear space.

② Data perturbation anxiety

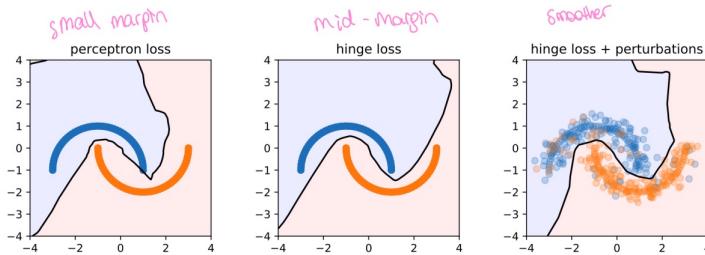
random perturbations to inputs & labels  
Gaussian noise and random flips

③ Representation perturbation Dropout method  
randomly turning the neurons on/off while training

④ Training noise

stochastic Gradient descent:  
- to regularize choose some random data (=stochastic)  
- then take a step according to these summed up  
gradients from random chosen bag of data

## Data Perturbation at Work



Observation:

- Data perturbations pull the decision boundary further away from the class manifolds, and typically further improves generalization.

St. 1

## SVMs & kernels:

- easy to optimize & regularize
- local opt for large datasets
- unique minimum
- not efficient to regularize because they are a squared algorithm

## NNs:

- can be trained on very large distributions (GPU)
- hard to optimize and regularize, we care more about generalization here.

# Lecture 10: Decision Trees and Random Forests

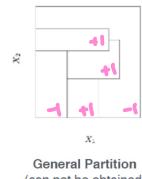
**Tree-based methods:**

- random forests segment the body to different parts and we make decisions (X-box)
- recommender systems
- banking
- medical ...

Given: Features:  $X_1, X_2, \dots, X_p$   
Target:  $Y$

Key idea:

- (1) Partition feature space into a set of rectangles.
- (2) Fit a simple model (e.g. a constant) in each one.

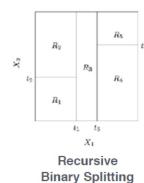
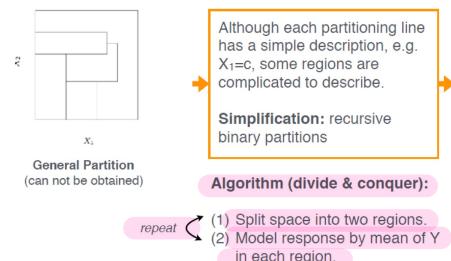


divide the space into some segments and have the very simple model predict every outcome of the rectangle

**Problem:** Dividing a large p-dimensional space to dimensions is hard

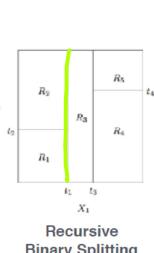
For each partition: model  $Y$  with a different constant.

very hard computation  
better to use divide & conquer



What is the best partition for  $X_i$  to have a certain classification?

## Concept of Decision Trees



① everything  $\leq t_1$  | everything  $> t_1$

Example:

- Split at  $X_1 = t_1$
- Split the region  $X_1 \leq t_1$  at  $t_2$
- etc.

② then do the same to  $t_2$

③ then to  $t_3$

compute the mean of the label for every box  $\Rightarrow$  Ensemble (bagging)

Result: Partitioning into  $R_1-R_5$

Model:  $\hat{f}(X) = \sum_{m=1}^5 c_m I\{(X_1, X_2) \in R_m\}$

Splitting is hard, can't draw the boxes in high dimensions, but we can draw a decision tree.

## Decision Trees – algorithm

Given:  $N$  observations,  $P$  features

$(x_i, y_i)$  for  $i = 1, 2, \dots, N$ , with  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ .

Wanted: splitting variables  $j$ , split points  $s$   
(to partition space into  $M$  regions)

Minimization criterion: e.g.  $\hat{c}_m = \text{ave}(y_i | x_i \in R_m)$  (sum of squares)

use in each region to find variance

Finding best binary partition is computationally infeasible (NP hard): use greedy approach

• split the area such that it gives good regression info.  
• then compute the average  
→ splitting  $p$ -dimension to  $s$  splits is too expensive  
so better to go step-by-step: **greedy approach**  
(not optimally!)

## Splitting

Given:  $N$  observations,  $P$  features

$(x_i, y_i)$  for  $i = 1, 2, \dots, N$ , with  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$ .

Consider a splitting variable  $j$  and split point  $s$  and define a pair of half planes:

$$R_1(j, s) = \{X | X_j \leq s\} \text{ and } R_2(j, s) = \{X | X_j > s\}$$

We seek  $j$  and  $s$  to solve:

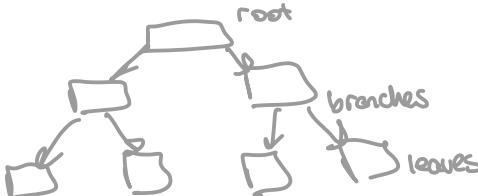
$$\min_{j, s} \left[ \min_{c_1} \sum_{x_i \in R_1(j, s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j, s)} (y_i - c_2)^2 \right]$$

### Splitting:

Inner minimization solved by:  $\hat{c}_1 = \text{ave}(y_i | x_i \in R_1(j, s))$  and  $\hat{c}_2 = \text{ave}(y_i | x_i \in R_2(j, s))$

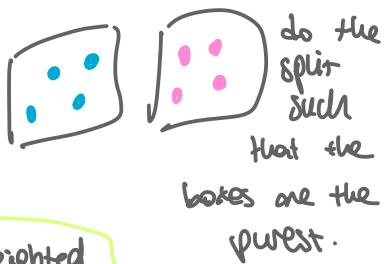
i.e. scan through all  $P$  features and determine optimal  $(j, s)$

- which variable should we split?  
**inner minimisation**



## → Splitting Metrics

To measure our splitting entropy: measure of disorder



### Information gain

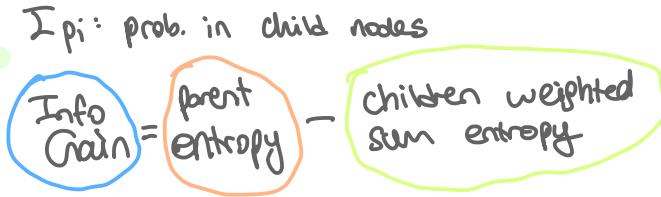
Based on the concept of entropy from information theory.

$$H(T) = I_E(p_1, p_2, \dots, p_J) = -\sum_{i=1}^J p_i \log_2 p_i$$

were  $p_1, p_2, \dots$  represent the probability of each class present in the child node that results form splitting the tree.

$$\text{Information Gain } IG(T, a) = \widehat{H(T)} - \widehat{H(T|a)}$$

Choose the split that results in the purest daughter nodes.

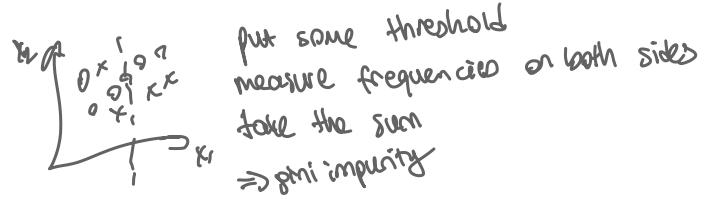


### Gini impurity

Measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset

$$GI(x) = \sum_{c=1}^C p(c)(1 - p(c)) = \sum_{c \neq \bar{c}} p(c)p(\bar{c})$$

Similar to entropy, zero if all samples belong to same class.



### Variance Reduction (often used for regression)

Variance reduction of a node is defined as the total reduction of the variance of the target variable due to the split at this node.

$$IV(N) = \frac{1}{|S|^2} \sum_{i \in S} \sum_{j \in S} \frac{1}{2} (x_i - x_j)^2 - \left( \frac{1}{|S_t|^2} \sum_{i \in S_t} \sum_{j \in S_t} \frac{1}{2} (x_i - x_j)^2 + \frac{1}{|S_f|^2} \sum_{i \in S_f} \sum_{j \in S_f} \frac{1}{2} (x_i - x_j)^2 \right)$$

were  $S$ ,  $S_t$ , and  $S_f$  are the set of pre-split sample indices for which the split test is true and false, respectively.

When to stop the split?

## Pruning

How large to grow the tree?

- too large: overfitting
- too small: not all important structure is captured

**Solution:** grow large tree, then prune using cost-complexity criterion  
(collapse any number of internal (non-leaf) nodes)

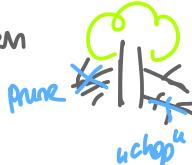
$$\begin{aligned} N_m &= \#\{x_i \in R_m\}, && \text{(training data in region)} \\ \hat{c}_m &= \frac{1}{N_m} \sum_{x_i \in R_m} y_i, && \text{(prediction in region)} \\ \text{mean} \\ Q_m(T) &= \frac{1}{N_m} \sum_{x_i \in R_m} (y_i - \hat{c}_m)^2, && \text{(squared loss)} \\ |T| & && \text{(number of terminal leaves)} \\ \alpha \geq 0 & && \text{(regularization parameter)} \end{aligned}$$

Find subtree that minimizes:  
 $C_\alpha(T) = \sum_{m=1}^{|T|} \frac{N_m Q_m(T)}{\text{cost}} + \frac{\alpha |T|}{\text{complexity}}$

a lot of splits → very large tree ⇒ overfit  
too small ⇒ not all important things are captured

⇒ PRUNE IT! cut branches or collapse them

small box size = small  $N \rightarrow MSE = 0$   
→ our cost = 0



but our tree is huge →  $|T|$  very large  
→ ↓ to balance the complexity  
→ idea: randomly remove nodes in the tree with splitting tests so that it's always true

Complexity will decrease if we start to collapse the nodes →  $|T|$  will decrease  
if we have low layers, we won't get a low training error.

too many layers → Complexity high, Cost low



Weakest Link Pruning:

1. Starting with the full tree, successively collapse the internal node that produces the smallest per node increase in  $\sum_m N_m Q_m(T)$

**Result:** Sequence of trees  $T_0, T_1, T_2, \dots, T_m$  (were  $T_m$  is a single-leaf tree)

2. Choose optimal  $\alpha$  (the optimal tree  $T_i$ ) via cross-validation.

## Decision Trees- advantages

- Interpretability
- Generalizes to higher dimensions  
*↳ no curse of dimensionality*
- Can handle mixed predictors: quantitative & qualitative
- Easily ignores redundant variables
- Handles missing data elegantly

**BUT:** decision trees are prone to overfitting (high variance)

## Bounding Generalization Error

[Haussler'88]: When the hypothesis space  $H$  is finite, then the generalization error of a consistent (i.e. correctly classifying) hypothesis  $h$  can be bounded as:

$$P(\text{error}_{\mathcal{X}}(h) > \varepsilon) \leq |H|e^{-m\varepsilon}$$

... where  $|H|$  is the size of the hypothesis space, and  $m$  is the number of iid samples in the training set.

**Observation:** generalization error grows with the number of hypotheses and decreases with the number of data points.

## Bounding Generalization Error

Limitation of Haussler's bound:

- It applies only to set of hypotheses that perfectly classify the data.  
What if we use a simple model (e.g. shallow decision tree or linear classifier?)

Another bound (derived from the Chernoff inequality) for the generalization error is given by:

$$\text{error}_{\mathcal{X}}(h) \leq \text{error}_{\mathcal{D}}(h) + \sqrt{\frac{\log |H| + \log \frac{1}{\delta}}{2m}}$$

training error ~ bias      error due to model complexity ~ variance

→ to bound the decision tree  
 variance here is the size of the hypotheses space of a decision tree

## Complexity of Decision Trees

- **Example 1:** General decision tree with  $n$  binary features (i.e. expanded until it perfectly classifies data):

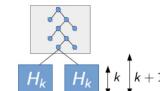
$$|H| = 2^n \quad \text{i.e. all possible lookup tables over } n \text{ binary features.}$$

- **Example 2:** Decision tree of depth 0:

$$|H| = 2 \quad H = \{\text{always labeling "0", always labeling "1"}\}$$

Solution for trees of depth  $k$ : Use induction

- Initial condition:  $|H_0| = 2$



- Recursion:  $|H_{k+1}| = n \times |H_k| \times |H_k|$

$$|H_{k+1}| = n \times |H_k| \times |H_k|$$

choice of root attribute      left tree      right tree

$\#$ : # tree depth  
 $n$ : # features

- Result  $\log_2 |H_k| = (2^k - 1)(1 + \log_2 n) + 1$

## Wrap-up

How large is the hypothesis space of decision trees?

$$\text{error}_{\mathcal{X}}(h) \leq \text{error}_{\mathcal{D}}(h) + \sqrt{\frac{\log |H| + \log \frac{1}{\delta}}{2m}}$$

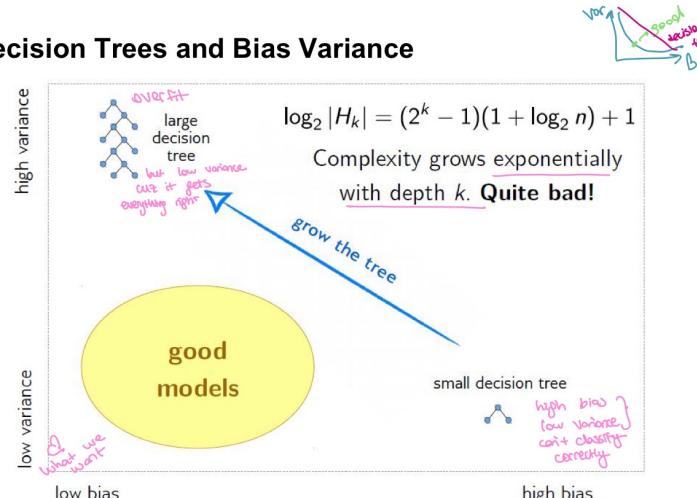
~ bias      ~ variance

**Decision tree:**

- Depth 0 (constant labeling)  $\log_2 |H_0| = 1$  - very high bias  
- low variance  
- isn't looking the data at all
- Depth  $k$   $\log_2 |H_k| = (2^k - 1)(1 + \log_2 n) + 1$
- Unrestricted  $\log_2 |H| = 2^n$  - very high bias  
- no bias

high bias / low variance  
↓  
low bias / high variance

## Decision Trees and Bias Variance



→ We need better capacity control.

**Idea:** combine decision trees with averaging, boosting.

# → Decision Tree regularization via model averaging (ensembles)

Generalization error:  $E[(y - \hat{f}(x))^2] = \text{Bias}[\hat{f}(x)]^2 + \text{Var}[\hat{f}(x)] + \sigma^2$

bias: low (if sufficiently deep)

+ variance: high (sensitive to choice of splits)

+ residual error

Consequence: decision trees often produce noisy or weak classifiers!

Idea: Combine the predictions of several randomized trees into single model.

Idea: build decision trees with some random setup, such that Var goes down, but Bias stays low.  
it needs to have  
- low training error  
- high variance needs to go away because of average

## Bagging (Bootstrap Aggregating)

Classifier  $C(\mathcal{S}, x)$   
training data  $\mathcal{S}$

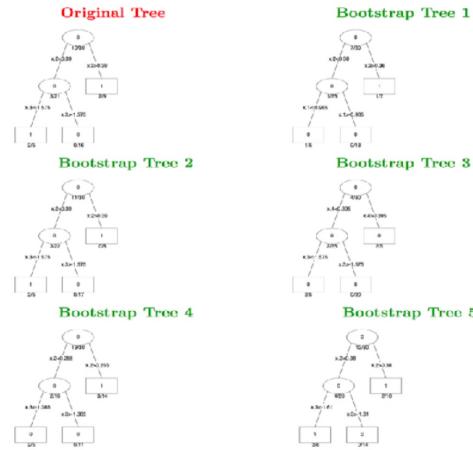
Subsample

Train  $B$  models on different dataset:

Draw  $\mathcal{S}^{*1}, \dots, \mathcal{S}^{*B}$   
samples sets of length  $N$   
(bootstrap samples: duplicates allowed!)

Prediction:

$\hat{C}_{\text{bag}}(x) = \text{Majority Vote } \{C(\mathcal{S}^{*b}, x)\}_{b=1}^B$   
(in classification; regression: average)



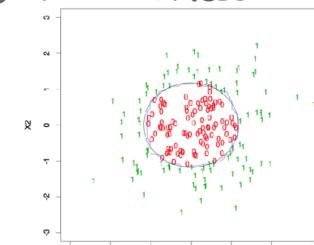
### → Limitations of bagging:

1) Datasets highly overlap:

predictions of different trees highly correlated

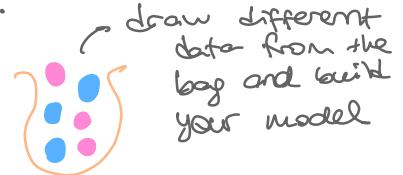
2) Variance is not reduced as much as wanted ↴

→ then: do a majority vote over all the bootstrap trees  
→ Variance is now reduced!



Improved prediction via reduced variance  
Bagging averages trees for smoother decision boundaries.

first:



# Random Forests

Random Subspace Method (Tin Kam Ho, 1998):

Randomly sample features  $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})$  without replacement

→ De-correlates estimators and decreases variance of the aggregate.

→ we take a subspace and sample different subspaces and build models for randomly sampled subspaces and add them up.

Random Forest (Breiman, 1999):

Combination of Bagging + Random Subspace Method applied to decision trees

- Complexity control through "out-of-bag" control (estimate generalization error using untrained samples)
- Random feature selection either at tree or split level

→ some samples are not used for training and we use that as the generalisation error  
→ to add randomness for variance reduction

## Random Forrests - algorithm

在 Bagging 的基础上引入了随机特征选择 (Random Feature Selection)：  
在构建每棵决策树时，除了对样本进行 Bootstrap 采样，还对特征子集进行随机采样（即每次分裂时只从部分特征中选择最佳划分）。这样可以进一步增加模型的多样性，减少相关性，提高泛化能力。

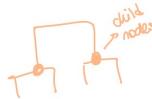
For  $b = 1$  to  $B$ :

- Draw a bootstrap sample (random sampling with replacement) of size  $N$  from the training data.
- Grow a tree  $T_b$  to the bootstrapped data until minimum node size reached:
  - Select  $m$  features at random.
  - Pick best variable/split-point among the  $m$ .
  - Split node into two daughter nodes.

→ ensemble of  $B$  trees where we did random sampling

→ algorithm simple but powerful

Result: ensemble of  $B$  trees.



Prediction:  $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$  (regression)

$\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$  (classification)

## Random Forest vs. Single tree:

- Train each tree on bootstrap resample of data.
- For each split: consider only a subset of  $m$  randomly selected features. (typically  $m = \sqrt{p}$  or  $\log_2 p$  were  $p$  is the total number of features)
- No pruning.
- Fit  $B$  trees this way and aggregate predictions.

### Single Trees

- + yield insight into decision rules
- + rather fast
- + easy to tune parameters
- predictions with high variance  
*don't generalize well*

### Random Forests

- + smaller prediction variance: better performance
- + easy to tune parameters
- slower
- "black box"

# Boosting

Idea: given a weak learner, run it multiple times on (reweighted) training data, then let learned classifiers vote.

On each iteration:

- Like in bagging, draw a sample of the observations from data (with replacement).
- Unlike in bagging, observations are not sampled randomly  
Higher weight observations are more likely chosen.
- Weight each training example by how incorrectly it was classified.  
Weight is higher for examples that are harder to classify.

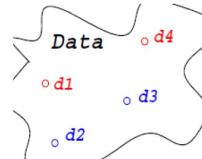
- prior distribution when drawing the data is changed  $\rightarrow$  essence of boosting
- harder to learn data will be evaluated better

## The Adaboost Algorithm

Input:  $N$  examples  $\{(x_1, y_1), \dots, (x_N, y_N)\}$

Initialize:  $d_i^{(1)} = 1/N$  for all  $i = 1 \dots N$

Do for  $t = 1, \dots, T$ ,



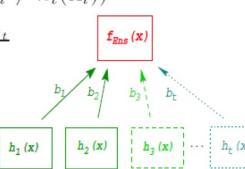
1. Train base learner according to example distribution  $d^{(t)}$  and obtain hypothesis  $h_t : x \mapsto \{\pm 1\}$ .

2. compute weighted error  $\epsilon_t = \sum_{i=1}^N d_i^{(t)} I(y_i \neq h_t(x_i))$

3. compute hypothesis weight  $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$

4. update example distribution

$$d_i^{(t+1)} = d_i^{(t)} \exp(-\alpha_t y_i h_t(x_i)) / Z_t$$



Output: final hypothesis  $f_{\text{Ens}}(x) = \sum_{t=1}^T \alpha_t h_t(x)$

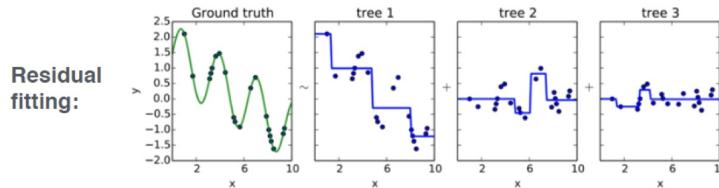
## Gradient Tree Boosting - conceptually

Idea: Iteratively improve a weak learner by correcting it:  $F_{m+1}(x) = F_m(x) + h(x)$

• all the trees are added up

A perfect correction  $h(x)$  implies:  $F_{m+1}(x) = F_m(x) + h(x) = y$   
i.e.  $h(x) = y - F_m(x)$

Hence we fit a new tree to the residual:  $y - F_m(x) \longrightarrow \frac{1}{2}(y - F(x))^2$   
(negative gradient of quadratic loss function)

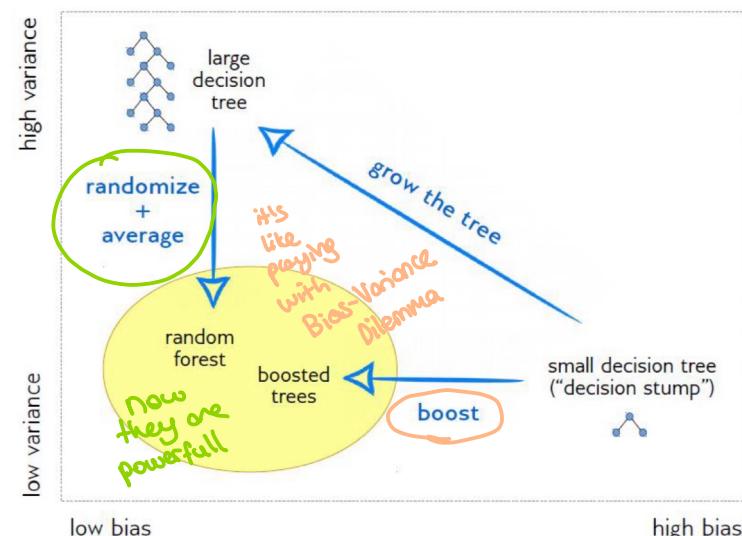


Gradient boosting is a gradient descent algorithm and can be trivially generalized with other loss functions and their gradients.

before: decision trees couldn't generalize well  
now: 2 approaches for low bias & low variance

- random forest
- boosted trees

## Decision Trees and Bias Variance



# Lecture 11: Boosting & Ensemble Learning

main goal: find  $\alpha_t$

## Ensemble Learning and Classification

- Ensemble for binary classification consists of
  - Hypotheses (basis functions)  $\{h_t(x) : t = 1, \dots, T\}$
  - \* of some hypothesis ("concept") set  
 $H = \{h \mid h(x) \mapsto \{\pm 1\}\}$  sign fn.
  - Weights  $\alpha = [\alpha_1, \dots, \alpha_T]$  ~ find this
    - \* satisfying  $\alpha_t \geq 0$
- Classification Output: weighted majority of the votes
  - $f_{\text{Ens}}(x) = \sum_{t=1}^T \alpha_t h_t(x)$  weighted sum
- How to find the hypotheses and their weights?
  - Bagging (Breiman, 1996):  $\alpha_t = 1/T$  T: # weak classifier
  - AdaBoost (Freund & Schapire, 1994)

- we wanna do binary classification
- we have  $T$  basis functions (hypotheses)
- non-linear  $\nearrow$
- for each of the  $T$  function we have a weight that is positive

→ Ensemble: weighted sum over all the hypotheses  
 $\hat{=} \quad$   
 Majority vote of the classification

Bagging: take all the functions  $h_t$  and think of them as small classifiers
 

- sum them up with equal weight:  $1/t$
- sum up things and get the majority decision

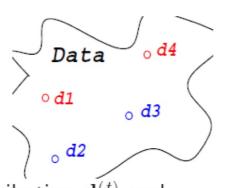
AdaBoost: reweight sum of the hypotheses in a way that you trust some more than others

## → AdaBoost Algorithm

Input:  $N$  examples  $\{(x_1, y_1), \dots, (x_N, y_N)\}$

Initialize:  $d_i^{(1)} = 1/N$  for all  $i = 1 \dots N$

Do for  $t = 1, \dots, T$ ,



1. Train base learner according to example distribution  $d^{(t)}$  and obtain hypothesis  $h_t : x \mapsto \{\pm 1\}$ .

2. compute weighted error  $\epsilon_t = \sum_{i=1}^N d_i^{(t)} I(y_i \neq h_t(x_i))$

3. compute hypothesis weight  $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$

4. update example distribution

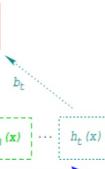
$$d_i^{(t+1)} = d_i^{(t)} \exp(-\alpha_t y_i h_t(x_i)) / Z_t$$

Output: final hypothesis  $f_{\text{Ens}}(x) = \sum_{t=1}^T \alpha_t h_t(x)$

- In each round the prob. of drawing that data is changed
- start with equal prob for all data

draw with equal prob.  
then build your model

→ only in the 1st iteration:  $d_i^{(1)} = \frac{1}{N} \begin{pmatrix} \text{prob. of} \\ \text{drawing} \end{pmatrix}$   
and after it changes



weight how for  
↓ from label  
hypothesis  
finally put:  $\alpha \cdot y \cdot h$  to  
your exp and then normalize  
with  $Z$  to make a probability  
that sums to 1.

- at 2: we see how good fit was. If  $h_t = -1$  when  $y_t = +1$ : then misclassification.  $\Rightarrow$  Error
- it counts how often the hypothesis  $h_t$  is wrong
- weight the error with the prob. of drawing  $x_i$

- Start with an even distribution
- train the model
- see where the model is wrong
- sum up the error
- depending on error amount, weight it
- look at each datapoint and look how strong your fit is

finally put:  $\alpha \cdot y \cdot h$  to  
your exp and then normalize  
with  $Z$  to make a probability  
that sums to 1.

### Exercise 2: AdaBoost as an Optimization Problem (25 + 25 P)

Consider AdaBoost for binary classification applied to some dataset  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ . The algorithm starts with uniform weighting ( $\forall_{i=1}^N : p_i^{(1)} = 1/N$ ) and performs the following iteration:

for  $t = 1 \dots T$ :

- Step 1:  $\mathcal{D}, p^{(t)} \mapsto h_t$  (learn  $t$ th weak classifier using weighting  $p^{(t)}$ )
- Step 2: weighted error:  $\epsilon_t = \mathbb{E}_{p^{(t)}}[1_{\{h_t(x) \neq y\}}]$  (compute the weighted error of the classifier)
- Step 3:  $\alpha_t = \frac{1}{2} \log \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)$  contribution ↑ (set its contribution to the boosted classifier)
- Step 4:  $\forall_{i=1}^N : p_i^{(t+1)} = Z_t^{-1} p_i^{(t)} \exp(-\alpha_t y_i h_t(x_i))$  (set a new weighting for the data)

The term  $\mathbb{E}_{p^{(t)}}[\cdot]$  denotes the expectation under the data weighting  $p^{(t)}$ , and  $Z_t$  is a normalization term. An interesting property of AdaBoost is that it can be shown to minimize some objective function

$$\text{closed form solution to minimize } G(\alpha) \rightarrow G(\alpha) = \sum_{i=1}^N \exp(-y_i f_{\alpha,t}(x_i))$$

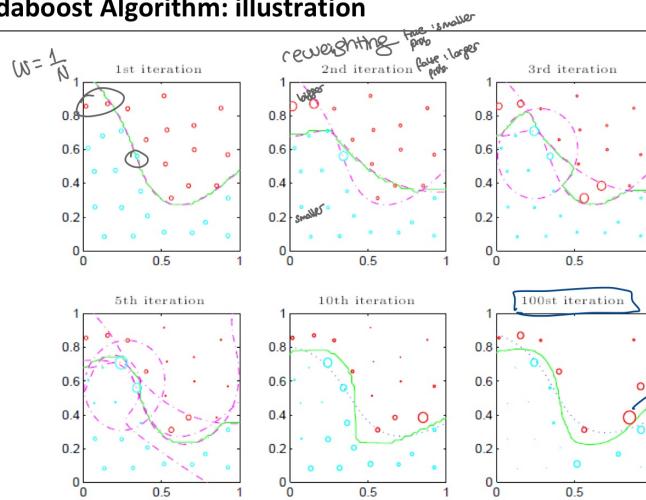
where  $f_{\alpha,t}(x) = \sum_{\tau=1}^t \alpha_\tau h_\tau(x)$  is the output score of the boosted classifier after  $t$  iterations.

$t=1 \dots N$

- wrong classified data gets higher prob. of drawing until it's correctly classified

#### AdaBoost Algorithm: illustration

- function gets completer in each iteration



- wrong classified points one bigger

- most of the data becomes irrelevant we only care about datapoints close to the db / misclassified

outliers & SVMs are important

→ The distribution changes on each iteration and we focus more and more on hard-to-get data.

#### Error Function of AdaBoost

- AdaBoost stepwise minimizes a function of

$$y_i f_{\alpha}(x_i) = y_i \sum_t \alpha_t h_t(x_i)$$

$$G(\alpha) = \sum_{i=1}^N \exp\{-y_i f_{\alpha}(x_i)\}$$

- The gradient of  $G(\alpha^{(t)})$  gives exactly the example weights used for AdaBoost:

$$\frac{\partial G(\alpha^{(t)})}{\partial f(x_i)} \sim \exp\{-y_i f_{\alpha}(x_i)\} \sim d_i^{(t+1)}$$

- The hypothesis coefficient  $\alpha_t$  is chosen, such that  $G(\alpha^{(t)})$  is minimized:

$$\alpha_t = \underset{\alpha_t \geq 0}{\operatorname{argmin}} G(\alpha^{(t)}) = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

- AdaBoost is a **gradient descent** method to minimizes  $G(\alpha)$ .

⇒ Bregman Divergences (Entropy Projections, ...)

⇒ Coordinate Descent Methods & Column Generation

- Boosting does step-wise minimization

→ we want that  $y_i$  and  $d_i h_t(x_i)$  has the same sign

- $G(d)$  is an exponential margin function.

## PAC Boosting – exponential convergence

*Upperbound on the training error:*

Theorem 1 (Schapire et al. 1997) Suppose AdaBoost generates hypotheses with weighted training errors  $\epsilon_1, \dots, \epsilon_T$ . Then we have

$$\sum_{i=1}^N \mathbb{I}(y_i \neq \text{sign}(f_{\text{Ens}}(\mathbf{x}_i))) \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t(1-\epsilon_t)} \quad \begin{matrix} \text{this bound decreases exponentially fast} \\ \text{error of ensemble} \end{matrix}$$

If  $\epsilon_t < \frac{1}{2} - \frac{1}{2}\gamma$  (for all  $t = 1, \dots, T$ ), then the training error will

decrease exponentially fast, i.e. will be zero after only

$$\frac{2 \log(N)}{\gamma^2} = \mathcal{O}(\log(N))$$

iterations.

## PAC Boosting – VC dimension of combined Hypothesis

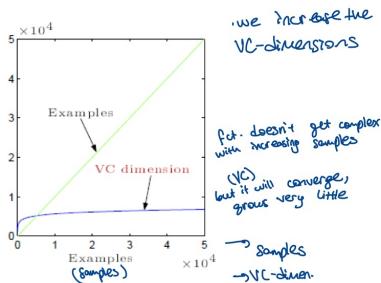
Let  $d$  be the VC dimension of the base hypothesis class  $\mathcal{H}$ .

Then the VC dimension of the class of combined functions is

$$d_{\text{Ens}}(N, \gamma) = \mathcal{O}\left(d \underbrace{\frac{\log(N)}{\gamma^2}}_{\sim T} \log\left(\frac{\log(N)}{\gamma^2}\right)\right) = \mathcal{O}(d \log(N) \log^2(N)).$$

An Example

- VC dimension  $d = 2$  (e.g. decision stumps)
- $\epsilon_t \leq 0.4 = \frac{1}{2} - \frac{1}{2}\gamma$   
 $\Rightarrow \gamma \geq 0.2$



## PAC Boosting – Digestion

- properties of weak learner imply exponential convergence to a consistent hypothesis
- Fast convergence ensures small VC dimension of the combined hypothesis
- small VC implies small deviation from the empirical risk
- for any  $\varepsilon > 0$  and  $\delta > 0$  exists a sample size  $N$ , such that with probability  $1 - \delta$  the expected risk is smaller than  $\varepsilon$
- Any weak learner can be boosted to achieve an arbitrary high accuracy! (~ strong learner)

• take an ensemble

• our error will be max 1/2  $\rightarrow$  yes/no  
our error is then smaller as:

$$\epsilon_t < \frac{1}{2} - \frac{1}{2}\gamma$$

$\gamma$  is smarter than chance

- if we do that our training error will decrease exponentially fast to zero

Weak learner: you have to be a bit smarter than luck

→ to limit overfit

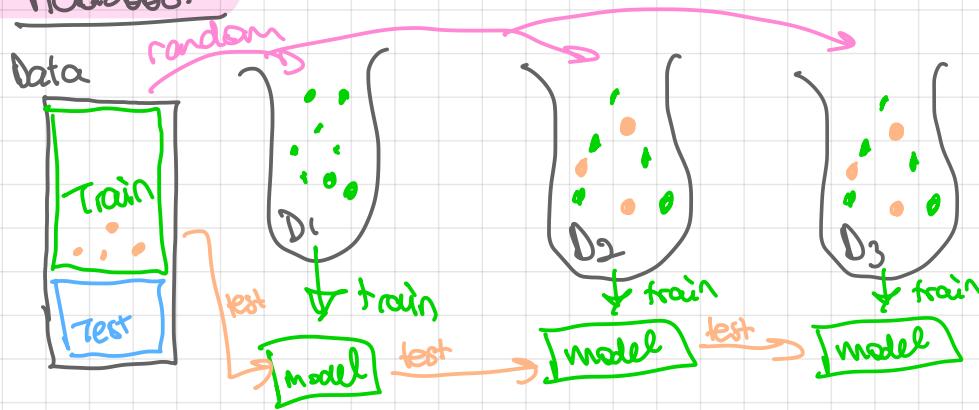
- the complexity doesn't go up even though we increase our ensemble because we are already at 0 training error.

VC describes the overall complexity by shattering.

→ because you don't need many ensembles  
→ because each learner can be very simple

boost to achieve arbitrary high accuracy

## Adaboost



1. select randomly from training data
2. take all training data and use it to test the model.
3. some points will be not well predicted  $\rightarrow$  error
4. for  $D_2$ : choose again randomly, but each instance is weighted according to the error of  $D_1$ .  
So points with significant error are more likely to be chosen

**Bagging:** choosing some subset of the data at random with replacement and create each bag in the same way.

**Boosting:** add-on to bagging, where in subsequent bags we choose those that has been modeled poorly in the overall system before.

→ both methods build an ensemble of models

**Bagging**:

- builds them parallel
- less training time
- classifiers are equally weighted

→ both reduce the variance  
but: bagging does not have bias reduction

• boosting very likely to overfit

**Boosting**:

- sequentially
- models are weighted in the ensemble based on their training performance

# → Margin Distributions

- Function set used in boosting: Convex Hull of  $\mathcal{H}$
- $$S := \left\{ f : \mathbf{x} \mapsto \sum_{h \in \mathcal{H}} \alpha_h h(\mathbf{x}) \mid \alpha_h \geq 0, \sum_{h \in \mathcal{H}} \alpha_h = 1 \right\}$$
- the  $\alpha$ 's are the parameters

- Find a hyperplane in the Feature Space spanned by the hypotheses set  $\mathcal{H} = \{h_1, h_2, \dots\}$

- Margin  $\rho$  for an example  $(\mathbf{x}_i, y_i)$  by proportional to the margin  $\approx$  SVM

$$\rho_i(\boldsymbol{\alpha}) := y_i f_{\text{Ens}}(\mathbf{x}_i) = y_i \sum_{t=1}^T \frac{\alpha_t}{\sum_t \alpha_t} h_t(\mathbf{x}_i) \xrightarrow{\substack{\text{proportional} \\ \text{of ensemble}}} \text{Normalized } \rho$$

rewriting of hypothesis

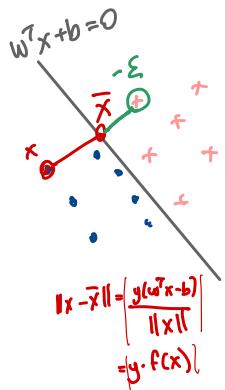
- Margin  $\varrho$  for a function  $f_{\text{Ens}}$  by  $\varrho(\boldsymbol{\alpha}) := \min_{i=1, \dots, N} \rho_i(\boldsymbol{\alpha})$

SVM: look at all datapoints find smallest distance to  $\mathbf{w}$   $\Rightarrow$  margin

Show for any datapoint that the nearest point on margin:  $\mathbf{x} - \bar{\mathbf{x}}$  is basically

$$\frac{\mathbf{y} \cdot \mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$

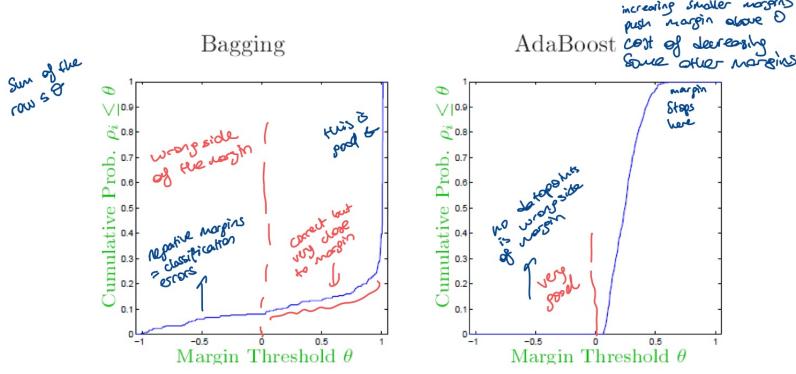
do it for every datapoint as a signed distance



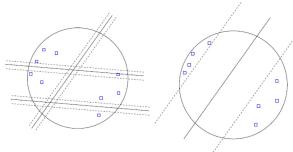
→ if the margin  $\rho_i$  is very large, then the model classified correctly  
→ margin of the whole Ensemble Function is the min one

## Margin Distributions - illustration

AdaBoost tends to increase small margins, while decreasing large margins



Small margin → many ways of shattering data  
large margin → not so many ways



## Lower Bounding the margin

Theorem 2 Suppose the base learning algorithm generates hypothesis with weighted training errors  $\epsilon_1, \dots, \epsilon_T$ . Then we have for any  $\theta$

can also be bounded by: margin is below threshold  $\theta$

$$P_Z(y f_{\text{Ens}}(\mathbf{x}) \leq \theta) \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\theta} (1-\epsilon_t)^{1+\theta}}$$

prob. that margin below threshold  $\theta$

Corollary 3 If the base learning algorithm always achieves

$\epsilon_t \leq \frac{1}{2} - \frac{1}{2}\gamma$  then AdaBoost will generate a combined hyperplane with margin at least  $\frac{1}{2}\gamma$ .

if each layer error is better than guessing ( $\frac{1}{2}$ )

$\gamma$ : fd. of our error gives a bound on the margin → use test error

→ prob. of finding a  $\theta$  that is smaller than our margin is exp. small

## SVM vs. Boosting

- SVMs

$$R[f] \leq R_{emp}[f] + O\left(\sqrt{\frac{\log(N\theta^2)}{\theta^2 N}} + \frac{\log(1/\eta)}{N}\right).$$

- Boosting

$$R[f] \leq R_{emp}^\theta[f] + O\left(\sqrt{\frac{d \log^2(\frac{N}{d})}{\theta^2 N}} + \frac{\log(1/\delta)}{N}\right)$$

depends on VC dimension of Bayes

- independent of the dimensionality of the space!

Margin  $\theta$  is part of emp. risk

→ both of them in terms of margins are the same:  
 large margin: small term at  $+O(\sqrt{\dots})$   
 small VC-Dimension of the hypothesis: small  $+O(\sqrt{\dots})$

## Boosting the limit

### What happens in the long run?

- Explicit expression for  $d_i^{(t+1)}$ :

$$d_i^{(t+1)} = \frac{\exp\{-\rho_i(\alpha^{(t)})\} \|\alpha^{(t)}\|}{\sum_{j=1}^N \exp\{-\rho_j(\alpha^{(t)})\} \|\alpha^{(t)}\|}$$

↔ Soft-Max Function with parameter  $\|\alpha^{(t)}\|_1$

- $\|\alpha\|_1$  will increase monotonically ( $\sim$  linear)

↔ the  $d$ 's concentrate on a few difficult patterns

→ Support Patterns

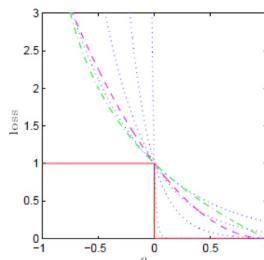
↔ Annealing Process:

$$\mathcal{G}(\alpha) = \sum \exp\{-\rho_i(\alpha)\} \|\alpha\|_1$$

→ 0/ $\infty$ -Loss approximated asymptotically

→ Barrier Optimization

→ with every iteration,  $\alpha$  will grow since vector gets longer: the 1-norm will increase monotonically.



## SVM's vs. Boosting

### Mathematical Program Formulation- SVMs

The SVM minimization of

$$\begin{aligned} & \min_{\mathbf{w} \in \mathcal{F}_\Phi, \rho \in \mathbf{R}_+} \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{subject to } y_i \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle \geq 1, \quad i = 1, \dots, N. \end{aligned} \quad (1)$$

reformulate as maximization of the margin  $\rho$

$$\begin{aligned} & \max_{\mathbf{w} \in \mathcal{F}_\Phi, \rho \in \mathbf{R}_+} \rho \\ & \text{subject to } y_i \sum_{j=1}^D w_j \Phi_j(\mathbf{x}_i) \geq \rho \quad \text{for } i = 1, \dots, N \\ & \|\mathbf{w}\|_2 = 1, \end{aligned} \quad (2)$$

where  $D = \dim(\mathcal{F})$  and  $\Phi_j$  is the  $j$ -th component of  $\Phi$  in feature space:

$$\Phi_j = P_j[\Phi]$$

- boosting also as a max margin problem
- tries to solve a linear program because:  $\|\mathbf{w}\|_1 = 1$

→ AdaBoost is very similar to SVMs: they both focus on hard to learn data.

SVM as a margin maximisation of some projected components

we'd like to have all margins  $> \rho$   
 push every individual norm to max  
 normalize  $\|\mathbf{w}\|_2$  as 1.  
 quadratic:  $\|\mathbf{w}\|_2 = 1$

### Boosting as a Mathematical Program

- master hypothesis

$$f(\mathbf{x}) = \sum_{t=1}^T \frac{w_t}{\|\mathbf{w}\|_1} h_t(\mathbf{x})$$

- base hypotheses  $h_t$  produced by the base learning algorithm.

Arc-GV solution is asymptotically the same as linear program solution, maximizing smallest margin  $\rho$ :

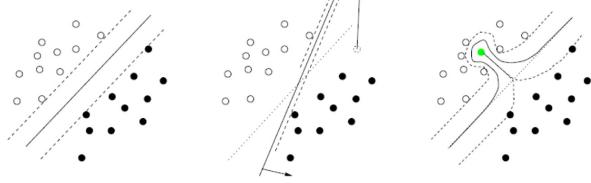
$$\begin{aligned} & \max_{\mathbf{w} \in \mathbf{R}^J, \rho \in \mathbf{R}_+} \rho \\ & \text{subject to } y_i \sum_{j=1}^J w_j h_j(\mathbf{x}_i) \geq \rho \quad \text{for } i = 1, \dots, N \\ & \|\mathbf{w}\|_1 = 1, \end{aligned} \quad (3)$$

where  $J$  is the number of hypotheses in  $\mathcal{H}$ .

# → Soft Margins

## Hard Margin Classification

no noise  
no margin



- The problem of finding a maximum margin "hyper-plane" on reliable data (left), data with outlier (middle) and a mislabeled pattern (right). The hard margin implies **noise sensitivity**.

• hard margin does not have slack variables (exceptions), so it's very sensitive to outliers

→ so we use soft margin sum to make it easier.

## AdaBoost with Soft Margins

- Define a **Soft Margin**

$$\tilde{p}_n(\alpha) = p_n(\alpha) + \zeta_n,$$

— where  $\zeta_n$  is the **amount of uncertainty** in example  $(x_n, y_n)$

- Once we have defined the **uncertainty measure**  $\zeta_n$ , we can easily get a new **regularized Boosting algorithm**.

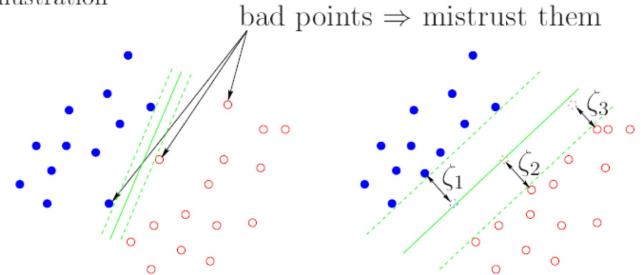
⇒ Improve the Error Function by plugging-in the **Soft Margin**

$$\tilde{G}(\alpha) = \sum_{n=1}^N \exp \{-\|\alpha\|_1 \tilde{p}_n(\alpha)\}$$

$$d_n^{t+1} = \frac{\partial \tilde{G}(\alpha)}{\partial f_\alpha(x_n)}$$

$$\alpha_t = \underset{\alpha_t \geq 0}{\operatorname{argmin}} \tilde{G}(\alpha)$$

Illustration



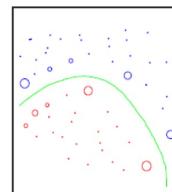
bad points ⇒ mistrust them

## Regularizing AdaBoost – Reducing the *Influence*

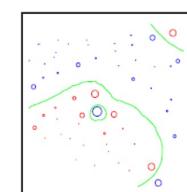
- How can we know **which** patterns are unreliable?

AdaBoost focuses on **difficult-to-learn** patterns by assigning high pattern weights  $d_n$  that we can exploit. Hence, we define the **Influence** of a pattern

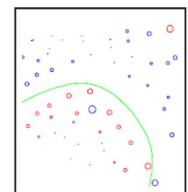
$$\mu_n^t = \sum_{r=1}^t \frac{\alpha_r}{\|\alpha\|_1} d_n^r \quad \zeta_n = C(\mu_n^t)^2$$



AdaBoost without "Noise"



AdaBoost with "Noise"



AdaBoost<sub>Reg</sub> with "Noise"

## Regularizing AdaBoost

### Positive:

- first algorithm that addresses overfitting in Boosting
- much improved results

### Negative:

- Modification on algorithmic level
- hard to analyze
  - which optimization problem is solved
  - no generalization results

### Idea:

- go back to beginning and redesign optimization problem
- use convergence results for leveraging
- apply margin bounds

- Boosting algorithms

- AdaBoost
- PAC Motivation
- Boosting with Large Margins
- Strategies for Dealing with High Dimensional Spaces
- Relations to Mathematical Programming & SVMs

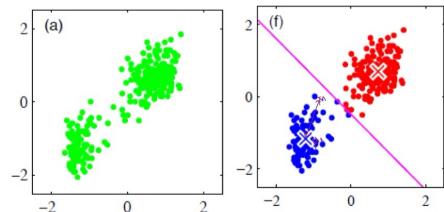
# Lecture 12: Clustering

Clustering: unsupervised  
no labels

to understand the underlying data structure  
would be nice to find a Gaussian distribution  
no labels, but we try to separate the data

Given  $N$   $d$ -dimensional datapoints  $X = \{x_1, \dots, x_N\}$  (no labels),

- how good is the cluster?
- How do we measure similarity?
  - assume data is uniform, homogeneous
  - find the structure of bunch of data points in high dimensional space



partition data into  $K$  disjoint sets  $\mathcal{S}_k$  based on similarity.

$$\forall k \forall I \mathcal{S}_k \cap \mathcal{S}_I = \emptyset, \quad \bigcup_{k=1}^K \mathcal{S}_k = X$$

$\underbrace{\text{sets shouldn't overlap}}$        $\overbrace{\text{union of our sets is the whole data}}$

## K-Means

Define clusters by minimum Euclidean distance to cluster mean.

→ Find  $\theta = \{\mathcal{S}_1, \dots, \mathcal{S}_K\}$  that minimize

$$J(\theta) = \sum_{k=1}^K \sum_{x_n \in \mathcal{S}_k} \text{squared distance of data point - respective cluster mean}$$

How much is the scatter around the mean

$$\text{where } \mu_k = \frac{1}{|\mathcal{S}_k|} \sum_{x_n \in \mathcal{S}_k} x_n$$

↓ normalizing

## Algorithm ("Expectation Maximization")

1. Choose  $K$  random points as initial cluster centers

$$\mu_1^{(0)}, \dots, \mu_K^{(0)}$$

(1) choose random  $\mu$ 's  
(2) expectation step: everything closest to the center goes to that center  
(3) assign the data points new to the cluster center & compute the  $\mu$ 's again - update mean  
(4) do it again to find local min

2. Assignment (E):

$$\mathcal{S}_k^{(t)} = \{x_n : \|x_n - \mu_k^{(t)}\|^2 \leq \|x_n - \mu_i^{(t)}\|^2 \forall i, 1 \leq i \leq K\}$$

3. Update (M):  $\mu_k^{(t+1)} = \frac{1}{|\mathcal{S}_k^{(t)}|} \sum_{x_n \in \mathcal{S}_k^{(t)}} x_n$

4. Iterate 2. and 3. until convergence to local minimum

you can't get global in from k-means

Ex: 100 datapoints, 5 clusters → very hard to solve with enumeration  $10^{67}$   
instead: do expectation maximization

- choose random  $\mu$ 's → choose  $K$  points
- find the closest center to that  $\mu$  → Assignment (E)
- compute the mean again → Update (M)
- do the assignment step again and again until local minimum is reached

- K-means clustering
- Gaussian Mixture Model
- Expectation Maximisation:
  - ML-bound for latent variables
  - lower bound on log-likelihood

Data comes from gaussian like clusters

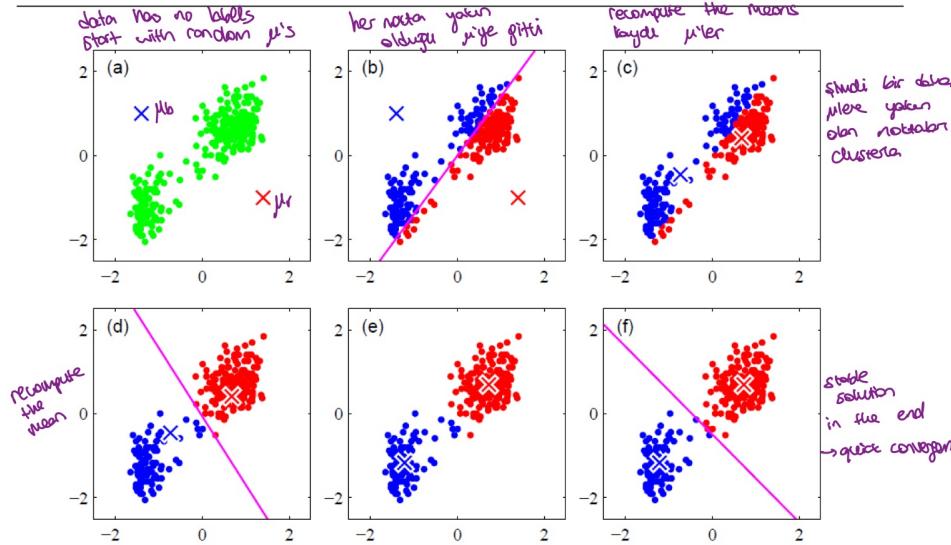
$\mu$ : center of the cluster

► How do we measure similarity?  
euclidian distance

► How do we judge similarity?  
distance all data points of that cluster  $x_n$  to their cluster center  $\mu_k$

• little variance around each cluster mean

## K-Means



assignment  
↳ recompute the mean  
until convergence  
→ very fast, if data not large.

How do you choose the number of clusters? Open problem no labels hard to model selection

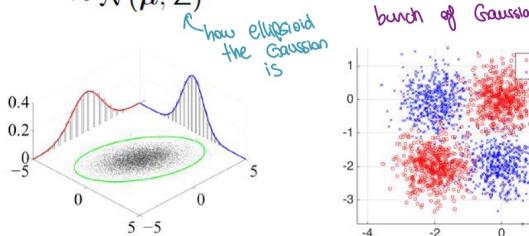
How to choose the 1st  $\mu$ ? Random, no right way

## Density Estimation

Multivariate Gaussian:

$$p(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left[ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$$

$$\sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$



→ 1 Gaussian is not enough  
idea: take bunch of Gaussians and weight them.  
scaling

Does not always model data (e.g. class-conditional densities) well.

## Gaussian Mixture Model (GMM)

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{k=1}^K \tau_k p_k(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

↑ weight of Gaussians

$$p_k(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$\tau_k$ : scaling or "prior" of  $p(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

$$\sum_{k=1}^K \tau_k = 1$$

normalize

Parameter vector  $\boldsymbol{\theta}$ :

$$\boldsymbol{\theta} = \{\tau_1, \dots, \tau_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_K\}$$

how they are scaled      means of all Gaussians      covs Gaussians

$\tau_k$  is our scaling factor

- we use GMM to do complex probability distributions
- GMM contains a parameter vector that has the means and covariances of all the Gaussians+ how they are superimposed ( $\tau_k$ )

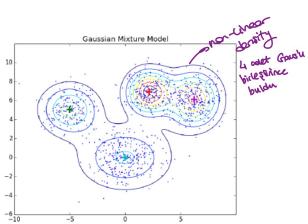
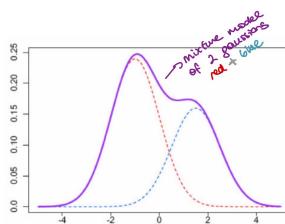
$$p(\mathbf{x}|\boldsymbol{\theta})$$

$\xrightarrow{\tau}$   $\xrightarrow{\boldsymbol{\mu}}$   $\xrightarrow{\boldsymbol{\Sigma}}$

Clustering is essentially finding & describing bunch of blobs and finding their structure.

# Gaussian Mixture models are universal density approximators.

- $\infty$  many  $P(x|\theta)$  can model anything, because it is summed up from smaller Gaussians.
- As a byproduct they provide a clustering solution.



## Fitting a GMM using EM

GMM is probabilistic & generative model

Trick: introduce auxiliary variables indicating the membership of each sample to a Gaussian

$$z_1, \dots, z_N \in \mathbb{R}^K \sim \text{Categorical}(\tau)$$

$$p(z_{nk} = 1) = \tau_k$$

$$\forall n \exists! k \quad z_{nk} = 1, z_{nj, j \neq k} = 0$$

e.g.  $z_n = (0, 0, 0, 1, 0, \dots, 0)^T$  tells the cluster  
 ↓  
 n-th data point  
 k-th cluster

if the n-th data point belongs to the k-th cluster

Note: this is also how to sample from a GMM

1. Sample  $z \sim \text{Categorical}(\tau)$  from which of the Gaussians?
2. Sample  $x \sim \mathcal{N}(\mu_k, \Sigma_k)$ , where  $z_k = 1 \rightarrow$  sampling from the respective distribution

## Fitting a GMM using EM: algorithm

1. Initialize  $t=0, \theta^{(0)} = \{\tau_1, \dots, \tau_K, \mu_1, \dots, \mu_K, \Sigma_1, \dots, \Sigma_K\}$   
 (e.g.,  $\tau_k^{(0)} = 1/K, \Sigma_k^{(0)} = I, \mu_k^{(0)} = \text{rand}$ )

2. **Expectation:** compute membership probabilities given  $\theta^{(t)}$

$$q^{(t)}(z_{nk}) := p(z_{nk} = 1 | x_n, \theta^{(t)}) \stackrel{\text{Bayes}}{=} \frac{p(x_n | z_{nk}, \theta^{(t)}) p(z_{nk}, \theta^{(t)})}{p(x_n | \theta^{(t)})}$$

$$= \frac{\tau_k^{(t)} p_k(x_n | \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{l=1}^K \tau_l^{(t)} p_l(x_n | \mu_l^{(t)}, \Sigma_l^{(t)})}$$

↓  
 lots of iterations  
 run

3. **Maximization:** update  $\theta$  given (soft) cluster assignments

$$\tau_k^{(t+1)} = \frac{1}{N} \sum_{n=1}^N q^{(t)}(z_{nk}) \quad \text{new means: } \mu_k^{(t+1)} = \frac{1}{N \tau_k^{(t+1)}} \sum_{n=1}^N q^{(t)}(z_{nk}) x_n$$

$$\Sigma_k^{(t+1)} = \frac{1}{N \tau_k^{(t+1)}} \sum_{n=1}^N q^{(t)}(z_{nk}) (x_n - \mu_k^{(t+1)}) (x_n - \mu_k^{(t+1)})^\top$$

weighted sum for the means

limit of soft assignment is hard assignment

→ with this we don't need to sample from the huge distribution but we can choose the k-th Gaussian we want

lump parameters:  $\tau, \mu, \Sigma$

variance  $\Sigma_k = \Sigma$ : unit (no ellipsoids allowed)  
 $\mu_k = \text{random}$

prob. of sampling from a certain Gaussian is  $\tau_k = 1/K \rightarrow$  uniform distribution

## Fitting a GMM using EM: algorithm

Note: also possible to use hard cluster assignments.

1. **Expectation:**

Max over all possible things

$$q^{(t)}(z_{nk}) = p(z_{nk} = 1 | x_n, \theta^{(t)})$$

Hard assignments:  
 Expectation is the max of all possible things

$$z_{nk}^{(t)} = \begin{cases} 1 & \text{if } q^{(t)}(z_{nk}) = \max_l q^{(t)}(z_{nl}) \\ 0 & \text{otherwise} \end{cases}$$

2. **Maximization:** update  $\theta$  given hard cluster assignments

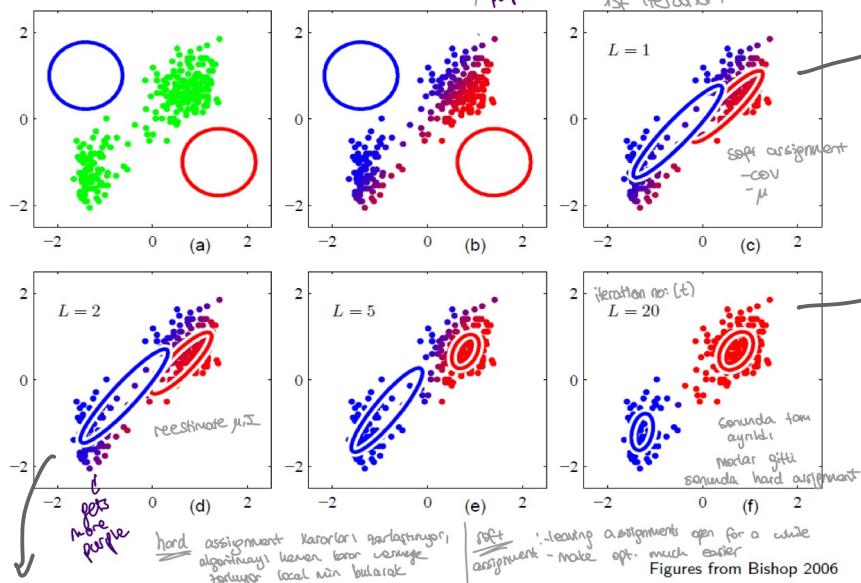
$$\tau_k^{(t+1)} = \frac{1}{N} \sum_{n=1}^N z_{nk}^{(t)} \quad \mu_k^{(t+1)} = \frac{1}{N \tau_k^{(t+1)}} \sum_{n=1}^N z_{nk}^{(t)} x_n$$

$$\Sigma_k^{(t+1)} = \frac{1}{N \tau_k^{(t+1)}} \sum_{n=1}^N z_{nk}^{(t)} (x_n - \mu_k^{(t+1)}) (x_n - \mu_k^{(t+1)})^\top$$

just summing

## Fitting a GMM using EM

soft assignment!  
middle ones can not  
decide



- signs are not that hard

compute cov and  $\mu$ 's

L: # iterations (t)

- find clusters in full color
  - solid cov-structure
  - hard assignments

Not  
hard  
blue

• classes aren't that easy to separate. Having a hard assignment will make the algorithm push to a local minimum. But if we have a soft assignment we're leaving assignments open for a while and that makes the optimisation numerically easier.

→ k-means assumes that classes have similar size. If the class sizes are not so similar, we'll get suboptimal solutions.

Assignment       $\xrightarrow{\text{latent variables inc}}$       Estimation      parameters  $(\mu, \sigma^2)$

Latent part: Where does the data belong to is unknown

first optimize randomly  
then refine

then estimate latent variables classes

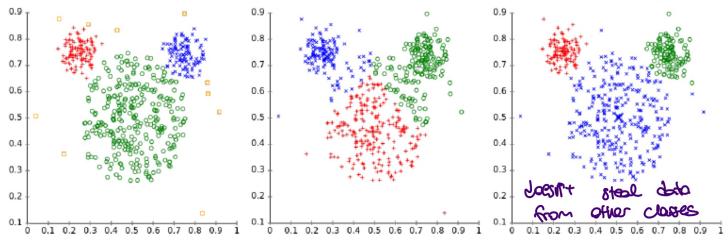
then given that the assignment was right, assign parameters to get  $\mu$ 's and  $\sigma$ 's.  
fixing the  $\mu$ 's and  $\sigma$ 's we can do again an assignment  
again do  $\mu$ 's &  $\sigma$ 's...

## GMM vs. k-means clustering

In contrast to K-means, GMM allows for

- Unequal cluster variances
  - Unequal cluster probabilities
  - Non-spherical clusters
  - Soft cluster assignment

Different cluster analysis results on "mouse" data set:  
Original Data                  k-Means Clustering                  EM Clustering



→ GrMm is more okay to different class sizes

→ K-means not ok for different class size

# → Fitting a GMM

# with Max-Likelihood

- Change  $\theta$  such that likelihood of generating  $X$  is maximum  $\Rightarrow p(X|\theta)$

## General EM Algorithm

### Maximum Likelihood for latent variable models

$z$  latent (unobserved variables)

$X$  observed data

$\theta$  model parameters

We want to maximize the likelihood of the observed data (= incomplete-data likelihood),  $L(\theta|X) = p(X|\theta)$ :

$$\hat{\theta} = \arg \max_{\theta} \log [p(X|\theta)] = \arg \max_{\theta} \log \left[ \sum_{z \in Z} p(X, z|\theta) \right].$$

Maximizing this directly is difficult because of  $\log \sum \dots$  headache

On the other hand, it is often easy to optimize the complete-data likelihood,  $L(\theta|X, z) = \log p(X, z|\theta)$ .

## EM= Expectation Maximisation

- universal tool to apply for latent variables
- $z$  (latent variables) tell us the underlying structure of data
- max log-likelihood w.r.t. model variables  
we are transparent that there's latent variables

## Trick: get rid of $\log \sum$ :

### Example: GMM

① Incomplete-data log-likelihood

$$\log [p(X|\theta)] = \sum_{n=1}^N \log \left[ \sum_{k=1}^K \tau_k p_k(x_n|\mu_k, \Sigma_k) \right]$$

② Complete-data log-likelihood

$$\log [p(X, z|\theta)] = \sum_{n=1}^N \log \left[ \sum_{k=1}^K \delta_{z_{nk}=1} \tau_k p_k(x_n|\mu_k, \Sigma_k) \right]$$

z<sub>nk</sub>: n<sup>th</sup> datapoint is in the k<sup>th</sup> cluster

→ Analytic ML estimate for each  $\theta_k = (\tau_k, \mu_k, \Sigma_k)$

**Problem:** We don't know  $z$ .

## EM Algorithm

$z$  is unknown, so estimate jointly with  $\theta$

Expectation Maximization algorithm:

- Iterate between updates of hidden variables and parameters
- **Theory:** updates are defined in a way such that  $p(X|\theta)$  increases in each step
- Guaranteed to find local maximum of  $p(X|\theta)$   
(hard to find global maximum if  $p(X|\theta)$  is non-concave)
- **Technically:** optimize a lower bound on  $p(X|\theta)$  and subsequently improve bound

Split the optimization to:  
 ↗ 1st part is taking care of unknown latent variables  
 ↘ other part " rest of parameters ( $\tau, \mu, \Sigma$ )

- update the Max-likelihood in a way that log-likelihood increased

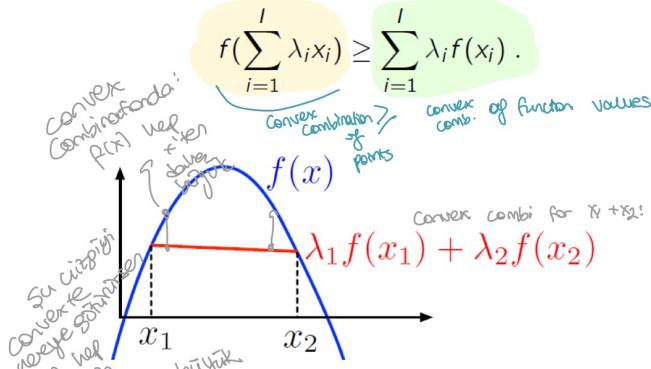
# → Jensen's Inequality

Instead of optimizing the function itself, maximize the lower bound to the function  
 → optimize a lower bound on  $p(x|\theta)$  and subsequently improve bound

## Jensen's Inequality

For any convex combination  $\lambda_1, \dots, \lambda_I$ ,  $\lambda_i \geq 0$ ,  $\sum_{i=1}^I \lambda_i = 1$

and any concave function  $f$ :



(conversely for convex  $f$ , analogous for continuous  $f$ )

## A lower bound on the log-likelihood

$\theta$ : a parameter setting

$q(z)$ : a probability mass function of choice on  $z$

$$\begin{aligned} \log p(X|\theta) &= \log \sum_z p(X, z|\theta) \quad | \cdot q(z) \\ &= \log \sum_z q(z) \left[ \frac{p(X, z|\theta)}{q(z)} \right] \quad \text{convex combination of some probabilities} \\ &\quad \rightarrow \text{Jensen's inequality} \downarrow \quad (\text{remember that } \sum_z q(z) = 1, \text{ log concave}) \\ &\quad \text{is lower bounded by:} \quad \geq \sum_z q(z) \underbrace{\log \left[ \frac{p(X, z|\theta)}{q(z)} \right]}_{\text{sum disparity, okay}} \\ &=: F(q(z), \theta) \quad \text{sum of abstract function of parameters and } q(z) \end{aligned}$$

This lower bound is much easier to optimize ( $\log \sum$  vs.  $\sum \log$ )

Jensen ile logdeki sum'ı disa atmaya caliyorum  
 lower bound最大化 amaci  
 q(z)yle sineye bi fonksiyon

## Expectation Maximization

$$\hat{\theta} = \arg \max_{\theta} \log p(X|\theta).$$

This is difficult.

**EM objective:** maximize, w.r.t.  $q$  and  $\theta$ , the lower bound

$$\hat{q}, \hat{\theta} = \arg \max_{q, \theta} F(q(z), \theta).$$

## Maximization of lower bound

$$F(q(z), \theta) = \sum_z q(z) \log \left[ \frac{p(X, z|\theta)}{q(z)} \right]$$

increase  $q(z)$   
 to get a larger function  
 change the  
 probabilities so  
 not too small  
 in proving

There are two ways to improve the lower bound:

1. **Expectation:** improve  $q(z)$  for given  $\theta$
2. Maximization: improve  $\theta$  for given  $q(z)$

# 1. How to select the $q(z)$ ?

Take the difference between what we actually want to optimize and its lower bound.  
How large is this difference? Can we minimize that?

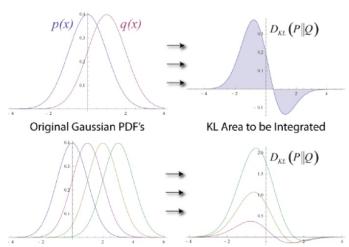
$$\text{we wanna actually optimize: } \log p(X|\theta) - \sum_z q(z) \log \left[ \frac{p(X, z|\theta)}{q(z)} \right] = \text{KL}(q(z) || p(z|x, \theta))$$

data log-likelihood

lower bound

## KL-Divergence

$$KL(P||Q) = \sum_x P(x) \log \frac{Q(x)}{P(x)} = \mathbb{E}_x \left[ \log \frac{Q(x)}{P(x)} \right] \geq 0$$



- Measures the distance between two distributions  $P$  and  $Q$
- Not a true metric (not symmetric, no triangle inequality)
- Important quantity in information theory

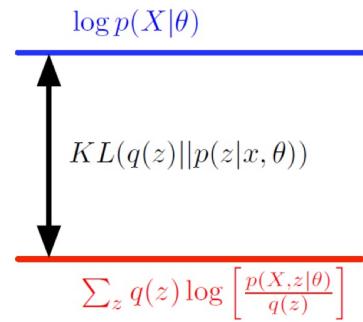
- a way of comparing prob. distributions
- positive quantity
- difference metric
- ideally we want the KL-divergence between data log likelihood and lower bound 0 so that they are the same.

### Expectation: improving $q(z)$

$$KL(q(z) || p(z|x, \theta)) = 0 \iff \underbrace{q(z)}_{\text{assignment}} = p(z|X, \theta)$$

⇒ Lower bound is strict if  $q(z) = p(z|X, \theta)$ .

**Expectation step:** set  $q^{(t)}(z) = p(z|X, \theta^{(t)})$  blue and red lines will be together



## 2. Now improve $\theta$ :

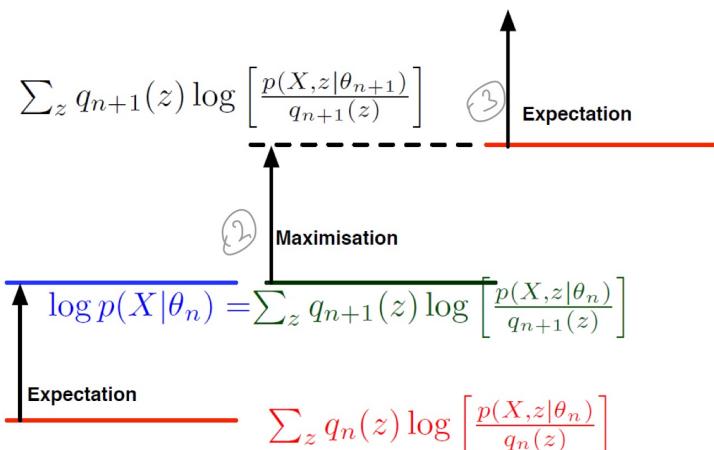
**Goal:** maximize  $F(q(z), \theta)$  w.r.t.  $\theta$ .

$$\begin{aligned} \theta^* &= \arg \max_{\theta} \sum_z q(z) \log \left[ \frac{p(X, z|\theta)}{q(z)} \right] \quad q \text{ is fix} \\ &= \arg \max_{\theta} \sum_z q(z) \log p(X, z|\theta) - \underbrace{\sum_z q(z) \log q(z)}_{\text{"Entropy" } H(q) \text{ Some constant, doesn't depend on } \theta} \\ &= \arg \max_{\theta} \sum_z q(z) \log p(X, z|\theta) \end{aligned}$$

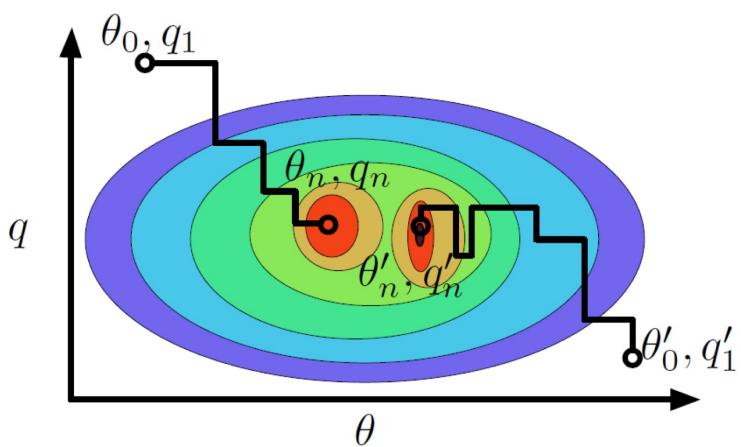
**Approach:** set gradient to zero ...

Typically, easy (analytic) solution, due to  $\sum \log$  rather than  $\log \sum$ .

### Iterative Optimization



### Convergence to local maximum



## EM summary

$z$  latent (unobserved variables)

$X$  observed data

$\theta$  model parameters

1. Initialize  $\theta^{(0)} = \text{rand}$

2. Expectation:  $q^{(t)}(z) = p(z|X, \theta^{(t)})$  assignment

3. Maximization:  $\theta^{(t+1)} = \arg \max_{\theta} \sum_z q^{(t)}(z) \log p(X, z|\theta)$   
↳ reestimating  $\mu$  and  $\Sigma$

4. Iterate until convergence local min als Ergebnis

## Summary

- EM is a "meta-algorithm" for obtaining local ML estimates
- Also applicable to maximum a-posteriori (MAP) estimation
- Particularly useful in models with latent variables  $z$ , where optimizing the incomplete-data likelihood directly is hard, but optimizing the complete-data likelihood  $p(X, z|\theta)$  is easy.

→ Alternate between estimating  $z$  and  $\theta$

- Can be applied to a GMM, but EM is not equal to a GMM
- Other applications:
  - Hidden Markov Models (Baum-Welch algorithm)
  - Missing/incomplete data
  - Only summary data observed

- Clustering basically measures similarity and judges the solution
- Optimisation of how to find the clustering solution?  
iterative, otherwise it would be expensive because combinatorial
- k-Means or Gmm
- k-Means is like a meta-algo because we optimise our error function and find the local max/min
- EM is a one big algorithm that we can apply to many models (Gmm included)

## Properties of EM

### Pro

- No stepsize/learning rate
- Each iteration improves likelihood

### Con

- "Only" local minima found
- Solution dependent on initialization
- Can be slow

**Note:** sometimes possible to use generic solvers (e.g. Newton)

### But:

- Complicated gradients, update rules  
↳ Jensen requires densities
- No improvement guarantee (e.g., Jensen requires densities)

# Lecture 13: Kernel Ridge Regression

## Classification

map vectors to output {0, 1}

predict classes

data generated from 2 classes

learn a decision fn. to map observations

Goal: predict class labels

LDA, SVM

## Regression

different values

map datapoint  $x$  to real value

price prediction, weather...

above / below a certain threshold

Goal: predict continuous function value

## Linear Regression

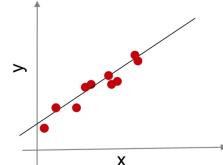
### Linear Regression

#### Problem:

Given a set of samples  $(x_1, y_1), \dots, (x_T, y_T)$  predict the outcome  $y$  for a new input  $x$

#### Model:

$$\hat{y} = w^T x$$



$$y = \underbrace{\begin{matrix} \downarrow \\ w^T \end{matrix}}_{\rightarrow} \quad \downarrow \quad X \in \mathbb{R}^{D \times T}$$

For optimal  $w$ :

minimize the loss between true and prediction:

$$\min E(w)$$

$$E(w) = \sum_{t=1}^T (y_t - w^T x_t)^2$$

Least-Square-Error

How to minimize LSE function?

$$E_{\text{LSE}}(w) = (y - w^T X)^2 = y^T y - 2y^T w^T X + w^T X X^T w$$

take derivative w.r.t.  $w$ :  $\frac{\partial E_{\text{LSE}}(w)}{\partial w} = -2X^T y + 2X^T X^T w$

set derivative to zero:  $2X^T X^T w = 2X^T y$

$$w = \underbrace{(X^T X^T)^{-1}}_{\text{Cov-matrix}} \cdot X^T y$$

→ However, linear regression is not robust to outliers

## → Ridge Regression

- has a regularizer (squared Euclidean norm of  $w$ ) with  $\gamma$  in loss function to control the complexity of the function.
- Whole idea comes from SLM: If we have a different set of solutions and if each of them correctly fit the data, then it is preferable to chose the smallest one because the smallest one is more likely to generalize well.
- RR adds a penalty term that shrinks the regression coefficients towards zero, reducing the model's complexity and thus reducing the risk of overfitting.

How to find the optimal  $w$  in RR?

$$E_{RR}(w) = (y - w^T X)^2 + \gamma \|w\|_2^2$$

→ Regularization term in RR

$$\frac{\partial E_{RR}(w)}{\partial w} = -2Xy + 2XX^Tw + \gamma 2w \stackrel{!}{=} 0$$

$$2w(\gamma + XX^T) = 2Xy$$

$$w = (XX^T + \gamma I)^{-1} \cdot Xy$$

Parameter  $\gamma$ :

- determines the influence of the regularizer:
- if  $\gamma$  big, then the values of  $w$  will go towards zero (shrinkage)
  - $\gamma = 0$ : we recover the linear model

higher  $\gamma \uparrow$ , flatter the predicted function

High values of  $\gamma$  is good for noisy high-dimensional data.

- if  $\gamma$  large: no small EVs.

# → Kernel Ridge Regression

nonlinear

- map data to a higher dimension and do the regression there if the relationship between input and output non linear:  $x \rightarrow \phi(x)$  mapping

Non-linear Ridge Regression  $\hat{y} = \mathbf{w}^\top \phi(\mathbf{x})$

$$\mathbf{w} = (\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1} \phi(\mathbf{X})\mathbf{y}$$

How to find the optimal  $\phi$ ? → Kernel Trick:

Optimal solution has to lie in the space spanned by the data.

New predictions of arbitrarily complex categorizations require **only** comparisons with previously learned examples

## the Kernel Trick

The prediction  $\hat{y}$  given a new data point  $\mathbf{x}$  can be expressed as linear combination of data points  $\mathbf{x}_i$

$$\hat{y} = f(\mathbf{x}) = \sum_{i=1}^N k_\phi(\mathbf{x}, \mathbf{x}_i) \alpha_i$$

where  $i = 1, \dots, N$  are the indices of previously seen data points.

Rewrite nonlinear ridge regression model again **in more detail**

The **kernel function**  $k$  corresponds to an inner product in feature space and measures the **similarity between two data points**

$$k_\phi(\mathbf{x}, \mathbf{x}_i) = \langle \phi(\mathbf{x}), \phi(\mathbf{x}_i) \rangle_\phi = \phi(\mathbf{x})^\top \underbrace{\phi(\mathbf{x}_i)}_{\text{seen}}$$

$$\begin{aligned} \hat{y} &= \phi(\mathbf{x})^\top \mathbf{w} = \phi(\mathbf{x})^\top (\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1} \phi(\mathbf{X})\mathbf{y} \\ &= \phi(\mathbf{x})^\top (\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1} \phi(\mathbf{X})(\mathbf{K} + \lambda I)(\mathbf{K} + \lambda I)^{-1}\mathbf{y} \\ &= \phi(\mathbf{x})^\top (\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1} \phi(\mathbf{X})(\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda I)(\mathbf{K} + \lambda I)^{-1}\mathbf{y} \\ &= \phi(\mathbf{x})^\top \underbrace{(\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1}(\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)}_I \phi(\mathbf{X})(\mathbf{K} + \lambda I)^{-1}\mathbf{y} \\ &= \underbrace{\phi(\mathbf{x})^\top \phi(\mathbf{X})}_{\mathbf{k}^*} (\underbrace{\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda I}_{\mathbf{K}})^{-1} \mathbf{y} \end{aligned} \quad (1)$$

$$\begin{aligned} \hat{y} &= \phi(\mathbf{x})^\top \mathbf{w} = \phi(\mathbf{x})^\top (\phi(\mathbf{X})\phi(\mathbf{X})^\top + \lambda I)^{-1} \phi(\mathbf{X})\mathbf{y} \\ &= \underbrace{\phi(\mathbf{x})^\top \phi(\mathbf{X})}_{\mathbf{k}^*} (\underbrace{\phi(\mathbf{X})^\top \phi(\mathbf{X}) + \lambda I}_{\mathbf{K}})^{-1} \mathbf{y} \end{aligned}$$

↑ new datapoint  
↑ matrix of all data

$$\text{apply Kernel Trick: } \hat{y} = \underbrace{\mathbf{k}^*}_{\mathbf{k}} (\mathbf{K} + \lambda I)^{-1} \mathbf{y} = \sum_{i=1}^N \underbrace{k_\phi(\mathbf{x}, \mathbf{x}_i)}_{\mathbf{k}} \alpha_i$$

$\mathbf{k}$  is now a scalar product of Kernel  $(\mathbf{x}, \mathbf{x}_i)$

Results in the kernel ridge regression model

$$\hat{y} = \sum_{i=1}^N k_\phi(\mathbf{x}, \mathbf{x}_i) \alpha_i$$

$$\text{with } \alpha = (\mathbf{K} + \lambda I)^{-1} \mathbf{y} \text{ and } \mathbf{K}_{i,j} = k_\phi(\mathbf{x}_i, \mathbf{x}_j)$$

The matrix  $\mathbf{K}$  is called **kernel matrix**;  $K_{i,j}$  denotes the similarity between data points  $\mathbf{x}_i$  and  $\mathbf{x}_j$  in the training data set.

# → Bayesian Linear Model

Apply Bayesian inference to model **supervised learning** problem:  
given samples  $(x_1, y_1), \dots, (x_T, y_T)$  make prediction for new input  $x^*$

1. Assume that the observation  $y$  arises from a **parametric model**  $f$  plus and independent, identically distributed noise term  $\varepsilon$ :

$$y = f_w(x) + \varepsilon \quad \text{with} \quad f_w(x) = x^\top w; \quad p(\varepsilon) = \mathcal{N}(0, \sigma^2)$$

2. Calculate **likelihood** of observed data:  $y = w \cdot x$

*thus likelihood is gaussian*

$$p(y | X, w) = \prod_{i=1}^n p(y_i | x_i, w) = \mathcal{N}(X^\top w, \sigma^2)$$

*data i.i.d. no correlation*

3. Define **prior**  $p(w)$  and calculate **posterior**  $p(w | X, y)$  using Bayes rule:

$$p(w) = \mathcal{N}(0, \Sigma)$$

$$p(w | X, y) = \frac{p(w)p(y | X, w)}{p(y | X)} = \frac{\text{prior} \times \text{likelihood}}{\text{marginal likelihood}}$$

$$\propto \mathcal{N}(\sigma^{-2} A^{-1} X y, A^{-1})$$

$$\text{with } A = \sigma^{-2} X X^\top + \Sigma^{-1}$$

4. Make **prediction** for new input  $x^*$ :

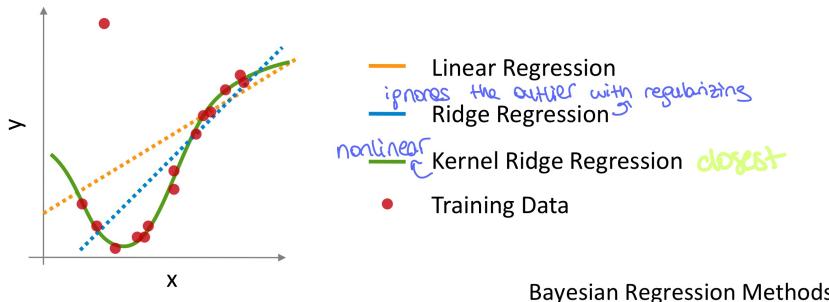
$$p(f^* | x^*, X, y) = \mathcal{N}(\sigma^{-2} x^{*\top} A^{-1} X y, x^{*\top} A^{-1} x)$$

$$\text{with } A = \sigma^{-2} X X^\top + \Sigma^{-1}$$

Note: In contrast to previous models the Bayesian prediction is not a single number but a **predictive distribution**.

The **mean** often serves as predicted value; the **variance** as an estimate of confidence.

## Overview:



Bayesian Regression Methods

- Bayesian methods **predict distribution** of outputs
- **Bayesian Linear Model**
  - Bayesian counterpart of linear regression
- **Gaussian Process**
  - Bayesian method for nonlinear regression
  - Can be viewed as inference in space of functions
  - Predictive mean equals the kernel ridge regression prediction, if a squared exponential covariance function is assumed

# → nonlinear Bayesian Model

To make the Bayesian approach nonlinear we reconsider the predicted mean

$$\bar{f}^* = \sigma^{-2} x^{*\top} A^{-1} X y \quad \text{with } A = \sigma^{-2} X X^\top + \Sigma^{-1}$$

1. Project data into feature space via function  $\Phi$

$$\bar{f}^* = \sigma^{-2} \phi(x^*)^\top A^{-1} \phi(X) y$$

with  $A = \sigma^{-2} \phi(X) \phi(X)^\top + \Sigma^{-1}$

2. Apply kernel trick

$$\bar{f}^* = \sigma^{-2} \phi(x^*)^\top A^{-1} \phi(X) y$$

$$\begin{aligned} &= \underbrace{\phi(x^*)^\top \Sigma \phi(X)}_{k^*} \underbrace{(\phi(X)^\top \Sigma \phi(X) + \sigma^2 I)^{-1} y}_{K} \\ &= k^*(K + \sigma^2 I)^{-1} y \end{aligned}$$

*Bayesian always the same result*

*Frequencies*

3. Resulting model is called **Gaussian Process**

$$\begin{aligned} \bar{f}^* &= k^*(K + \sigma^2 I)^{-1} y \\ \text{var}(f^*) &= k(x^*, x^*) - k^*(K + \sigma^2 I)^{-1} k^{*\top} \end{aligned}$$

Again a **kernel function**  $k$  measures the similarity between samples

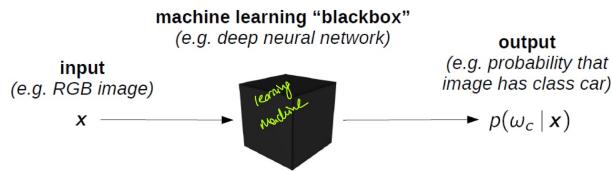
$$\begin{aligned} \text{here } k(x_i, x_j) &= \phi(x_i)^\top \Sigma \phi(x_j) \\ k^* &= k(x^*, X); \quad K = k(X, X) \end{aligned}$$

Method	biased	robust	nonlinear
Linear Regression	no	no	no
Ridge Regression	yes	yes	no
Kernel Ridge Regression	yes	yes	yes

# Lecture 14: XAI

- Explaining non-linear or single decisions is difficult.
- Single decisions are based on heatmaps, model thinks its red because of the ship

## Interpreting with class prototypes



**Goal:** Understand what the model considers as a "typical car".

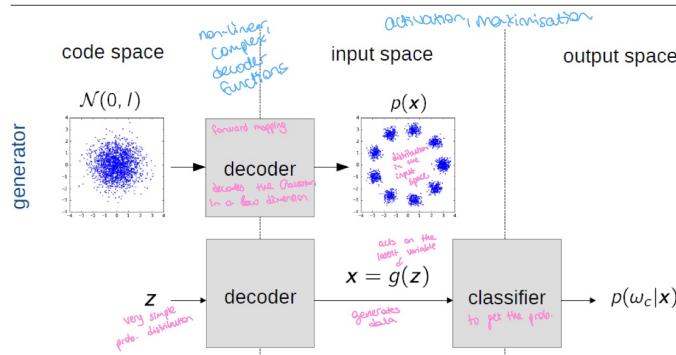
### Activation maximization framework:

Given: Who will can we make  $x$  such that it has the generality of the prototype?

$$\max_{x \in \mathcal{X}} \log p(w_c | x) - \lambda \|x\|^2$$

class prototype      prob. model scan  $X$  space find the answer that Max the prob.

## Building Prototypes using a generator



Activation maximization in code space: not in very high-dimensional  $x$  space anymore

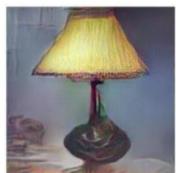
$$\max_{z \in \mathcal{Z}} \log p(w_c | g(z)) - \lambda \|z\|^2$$

was  $x$  before regenerating in code space

## Types of Interpretation

### Global Interpretation:

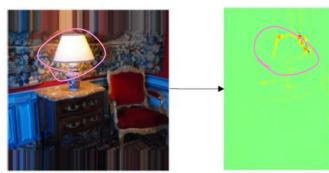
Understanding how a lamp typically looks like.



model's prototypical lamp

### Local Interpretation:

Understanding why this image of a lamp contains a lamp.



## Building more natural prototypes

$\max_{x \in \mathcal{X}} \log p(w_c | x) - \lambda \|x\|^2$  find input pattern that produces the strongest class activation

replace l2-norm by a data-dependent regularizer

$\max_{x \in \mathcal{X}} \log p(w_c | x) + \log p(x)$  we only allow density of our observation

$\max_{x \in \mathcal{X}} \log p(x | w_c)$  apply the Bayes theorem

find most likely input pattern for a given class

**Insight:** a good prototype must also depend on the data distribution.

## NN: classifier

take a latent space and generate  $x$ -space latent space very low dim.  
so it's very easy to find its max  
→ Class prototypes in latent space

## Approaches to interpretability

by the book interpretable

### Ante-hoc interpretability:

Choose a model that is readily interpretable and train it.

Example:

$$f(x) = \sum_{i=1}^d g_i(x_i)$$

contribution of  $i$ th variable

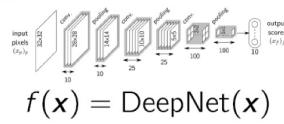
Is the model expressive enough to predict the data?

ex: decision tree upper set interpretable  
only 1st layer < 2nd layer gets complicated

### Post-hoc interpretability:

Choose a model that works well in practice, and develop a special technique to interpret it.

Example:



lets take the best model

How to determine the contribution each input variable?

# Explaining NN predictions

Explanation: "Which pixels contribute how much to the classification"  
 (what makes this image to be classified as a car)

$$f(x) = \sum_p h_p$$

Sensitivity / Saliency: "Which pixels lead to increase/decrease of prediction score when changed"  
 (what makes this image to be classified more/less as a car)

$$h_p = \left\| \frac{\partial}{\partial x_p} f(x) \right\|_\infty$$

↑ of the nonlinear  
wrt one point + x

Deconvolution: "Matching input pattern for the classified object in the image"  
 ↳ activity maximization (relation to  $f(x)$  not specified)

Each method solves a **different** problem!!!

change  
increase/decrease

*this way in a NN*

classification → explanation

*this way in a NN*

$$x_j = \sigma(\sum_i x_i w_{ij} + b_j)$$

$$r_j = f(x)$$

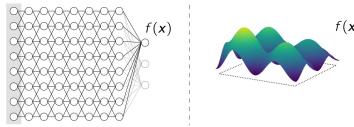
- model fully trained
- don't do backpropagation
- only to understand the model

$$r_i = x_i \sum_j \frac{w_{ij} r_j}{\sum_i x_i w_{ij}} = x_i c_i$$

$r_i$  depends on the activations **and** the weights

## → Activation Maximization

► Step 1: Think of the ML model as a function of the input



► Step 2: Explain the function  $f$  by generating a maximally activating input pattern:

$$x^* = \arg \max_x f(x, \theta)$$

**Problem:** In most cases  $f(x)$  does not have single point corresponding to the maximum.

E.g. in linear models,  $f(x) = \mathbf{w}^\top \mathbf{x} + b$ , we can keep moving the point  $\mathbf{x}$  further along the direction  $\mathbf{w}$ , and the output continues to grow).

Therefore, we would like to apply a preference for 'regular' regions of the input domain, i.e.

$$x^* = \arg \max_x [f(x) + \Omega(x)]$$

In practice, the preference can be for data points with small norm (i.e. we set  $\Omega(x) = -\lambda \|\mathbf{x}\|^2$  so that points with large norm are penalized.)

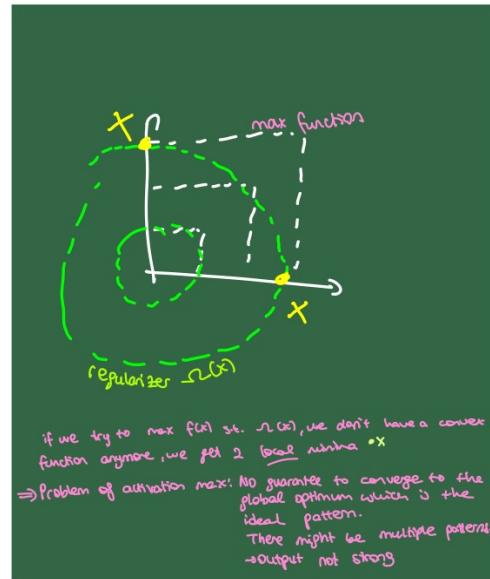
## Activation Maximization: Examples

$f(x) = \mathbf{w}^\top \mathbf{x} + b$  and  
 $\Omega(x) = -\lambda \|\mathbf{x}\|^2$

$$\begin{aligned} \frac{\partial}{\partial x} [f(x) + \Omega(x)] & \text{to find the global optimum} \\ &= \omega + (-2\lambda x) = 0 \\ x^* &= \frac{\omega}{2\lambda} \quad \text{optimal pattern } x \\ & x \uparrow \omega \uparrow \end{aligned}$$

it makes sense to find an  $x$  that also grows with  $\omega$ .  
 ↑ large  $\omega$  means strong regularization, smaller  $x$  will get.  
 → The more we regularize our activation max. solution, the more  $x$  will go to origin ( $= 0$ ).

$$f(x) = \max(x_1, x_2) \text{ and} \\ \Omega(x) = -\lambda \|\mathbf{x}\|^2$$



# Activation Maximization: Probability View

Assume the model produces a log-probability for class  $\omega_c$ :

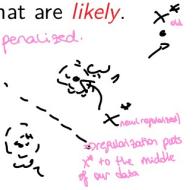
$$f(\mathbf{x}) = \log p(\omega_c | \mathbf{x})$$

The input  $\mathbf{x}^*$  that maximizes this function can be interpreted as the point where the classifier is the most *sure* about class  $\omega_c$ .

Choose the regularizer  $\Omega(\mathbf{x}) = \log p(\mathbf{x})$ , i.e. favor points that are *likely*. ↗ everything unlikely gets penalized.

The optimization problem becomes:

$$\begin{aligned}\mathbf{x}^* &= \arg \max_{\mathbf{x}} \log p(\omega_c | \mathbf{x}) + \log p(\mathbf{x}) \\ &= \arg \max_{\mathbf{x}} \log p(\mathbf{x} | \omega_c)\end{aligned}$$

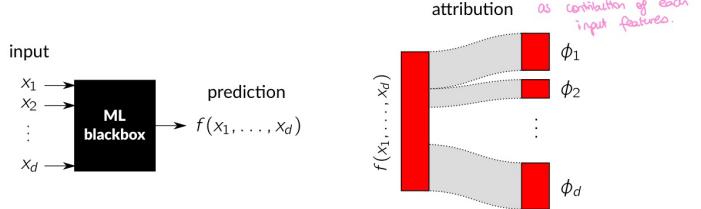


where  $\mathbf{x}^*$  can now be interpreted as the most *typical* input for class  $\omega_c$ .

## Afftribution

### Attribution of a Prediction to Input Features

To understand the prediction for a given datapoint



1. The data  $\mathbf{x} \in \mathbb{R}^d$  is fed to the ML black-box and we get a prediction  $f(\mathbf{x}) \in \mathbb{R}$ .
2. We explain the prediction by determining the contribution of each input feature.

Key property of an explanation: **conservation** ( $\sum_{i=1}^d \phi_i = f(\mathbf{x})$ ).

Sum of scores is  $f(\mathbf{x})$   
For  $i$  is a number of different input features

contribution of each input feature:

- Shapley Values
- Taylor Expansions
- LRP Method

## → Shapley Values

The Shapley values  $\phi_1, \dots, \phi_d$  measuring the contribution of each feature are:

$$\phi_i = \sum_{S: i \notin S} \frac{|\mathcal{S}|!(d-|\mathcal{S}|-1)!}{d!} [f(\mathbf{x}_{S \cup \{i\}}) - f(\mathbf{x}_S)]$$

where  $(\mathbf{x}_S)_S$  are all possible subsets of features contained in the input  $\mathbf{x}$ .

$\phi_i$ : each features contribution needs to be calculated

take the feature list as:

1.  $\emptyset$
2. other features alone
3. other features combined

$$d_5 = \frac{|\mathcal{S}|!(d-|\mathcal{S}|-1)!}{d!}$$

$$S: \text{size of condition} \quad S_\emptyset = 0 \quad S_{i,1} = 1 \quad S_{2,3} = 2$$

d: # features

each input variable is a player and function output is the profit realized by the cooperating players

### Attribution: Shapley Values

for  $\phi_i$ :  $S = \emptyset, i=1$   
Sum of all coalitions that is not or is  $i$

$$\text{Recall: } \phi_i = \sum_{S: i \notin S} \underbrace{\frac{|\mathcal{S}|!(d-|\mathcal{S}|-1)!}{d!}}_{\text{size of condition}} \underbrace{[f(\mathbf{x}_{S \cup \{i\}}) - f(\mathbf{x}_S)]}_{\Delta_S}$$

Input with feature  $i$       Input without feature  $i$

for  $\phi_1$ :  $\emptyset = 1, d=3$   
for  $\phi_2$ :  $\emptyset = 1, d=2$   
for  $\phi_3$ :  $\emptyset = 1, d=1$

**Worked-through example:** Consider the function  $f(\mathbf{x}) = x_1 \cdot (x_2 + x_3)$ . Calculate the contribution of each feature to the prediction  $f(\mathbf{1}) = 1 \cdot (1+1) = 2$ .

Feature	$\emptyset$	$x_1$	$x_2$	$x_3$	$x_1 \cup x_2$	$x_1 \cup x_3$	$x_2 \cup x_3$	$x_1 \cup x_2 \cup x_3$
Look at all possible feature sets that are not or are $i$ certain to be	$\emptyset$ empty set	$\frac{1}{3} \cdot 0$	$\emptyset$	$\emptyset$	$\frac{1}{3} \cdot 1$	$\emptyset$	$\emptyset$	$\emptyset$
feature	$\emptyset$	$\frac{1}{3} \cdot 1$	$\{x_2\}$	$\{x_3\}$	$\{x_1\}$	$\{x_1, x_2\}$	$\{x_1, x_3\}$	$\{x_2, x_3\}$
size of condition	$0$	$1$	$1$	$1$	$1$	$2$	$2$	$2$
$d$	$1$	$2$	$2$	$2$	$3$	$3$	$3$	$3$
$d_5$	$1$	$1$	$1$	$1$	$1$	$1$	$1$	$1$
$\Delta_S$	$0$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$
$\phi_i$	$0$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$\sum \phi_i$	$0$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$
$\sum \phi_i = 1$	$0$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

Shapley value satisfies conservation axiom:  
All contribution = sum of output of action  
I:  $\phi_1 = 1/2$   
 $\phi_2 = 1/2$   
 $\phi_3 = 1/2$   
also symmetry:  $x_1+x_2$  does have the same effect:  $\phi_1 = \phi_2$

$$\phi_3 = \frac{1}{2}$$

## → Taylor Expansions

- The function can be approximated locally by some Taylor expansion:

$$f(\mathbf{x}) = f(\tilde{\mathbf{x}}) + \sum_{i=1}^d \underbrace{[\nabla f(\tilde{\mathbf{x}})]_i \cdot (\mathbf{x}_i - \tilde{\mathbf{x}}_i)}_{\phi_i} + \dots$$

- First-order terms  $\phi_i$  of the expansion can serve as an explanation.
- The explanation ( $\phi_i$ ) depends on the choice of root point  $\tilde{\mathbf{x}}$ .

**Example:** Attribute the prediction  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  with  $\mathbf{x} \in \mathbb{R}^d$  on the  $d$  input features.

$$\begin{aligned}\tilde{\mathbf{x}} &= \mathbf{0} \\ \phi_i &= [\nabla f(\tilde{\mathbf{x}})]_i \cdot (\mathbf{x}_i - \tilde{\mathbf{x}}_i) \\ &= w_i \cdot (\mathbf{x}_i - \mathbf{0}) \\ &= w_i x_i\end{aligned}$$

- root point  $\tilde{\mathbf{x}} = \mathbf{0}$
- take derivative of the function
- $(\mathbf{x}_i - \text{root point}) \cdot \text{derivative}$

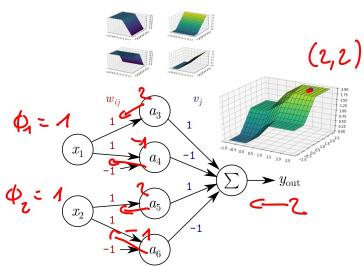
→ other you have your first order terms as explanation  
Because gradient is the same at the entire function.

**Problem:** Gradient information too localized  
it can't explain discontinuities in a function.

**Solution:** decompose the problem in multiple parts  $\Rightarrow$  LRP Method

## → LRP Method

**Example:** Consider  $y_{\text{out}}$  to be the quantity to explain.



- Step 1: Propagate on the hidden layer

$$\forall_{j=3}^6 : R_j \leftarrow \frac{a_j v_j}{\sum_{j=3}^6 a_j v_j} y_{\text{out}}$$

- Step 2: Propagate on the first layer

$$\forall_{i=1}^2 : R_i \leftarrow \sum_{j=3}^6 \frac{x_i w_{ij}}{\sum_{i=1}^2 x_i w_{ij}} R_j$$

**Note:** Other propagation rules can be engineered, and choosing appropriate propagation rules is important to ensure LRP works well for practical neural network architectures.

- to identify the contributions of each feature

$$\mathbf{x} \rightarrow \mathbf{a} \rightarrow \mathbf{R}$$

if weight: -1  $\rightarrow a \cdot w = 0$  ignore  
weight: +1  $\rightarrow$  take it ✓

## → Taylor Expansions

- The function can be approximated locally by some Taylor expansion:

$$f(\mathbf{x}) = f(\tilde{\mathbf{x}}) + \sum_{i=1}^d \underbrace{[\nabla f(\tilde{\mathbf{x}})]_i \cdot (\mathbf{x}_i - \tilde{\mathbf{x}}_i)}_{\phi_i} + \dots$$

- First-order terms  $\phi_i$  of the expansion can serve as an explanation.
- The explanation ( $\phi_i$ ) depends on the choice of root point  $\tilde{\mathbf{x}}$ .

**Example:** Attribute the prediction  $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$  with  $\mathbf{x} \in \mathbb{R}^d$  on the  $d$  input features.

$$\begin{aligned}\tilde{\mathbf{x}} &= \mathbf{0} \\ \phi_i &= [\nabla f(\tilde{\mathbf{x}})]_i \cdot (\mathbf{x}_i - \tilde{\mathbf{x}}_i) \\ &= w_i \cdot (\mathbf{x}_i - \mathbf{0}) \\ &= w_i x_i\end{aligned}$$

- root point  $\tilde{\mathbf{x}} = \mathbf{0}$
- take derivative of the function
- $(\mathbf{x}_i - \text{root point}) \cdot \text{derivative}$

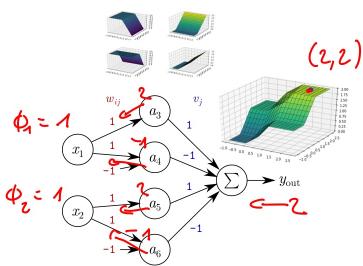
→ other you have your first order terms as explanation  
Because gradient is the same at the entire function.

**Problem:** Gradient information too localized  
it can't explain discontinuities in a function.

**Solution:** decompose the problem in multiple parts  $\Rightarrow$  LRP Method

## → LRP Method

**Example:** Consider  $y_{\text{out}}$  to be the quantity to explain.



- Step 1: Propagate on the hidden layer

$$\forall_{j=3}^6 : R_j \leftarrow \frac{a_j v_j}{\sum_{j=3}^6 a_j v_j} y_{\text{out}}$$

- Step 2: Propagate on the first layer

$$\forall_{i=1}^2 : R_i \leftarrow \sum_{j=3}^6 \frac{x_i w_{ij}}{\sum_{i=1}^2 x_i w_{ij}} R_j$$

**Note:** Other propagation rules can be engineered, and choosing appropriate propagation rules is important to ensure LRP works well for practical neural network architectures.

- to identify the contributions of each feature

$$\mathbf{x} \rightarrow \mathbf{a} \rightarrow \mathbf{R}$$

if weight: -1  $\rightarrow a \cdot w = 0$  ignore  
weight: +1  $\rightarrow$  take it ✓