

CH 7 PARAMETER ESTIMATION

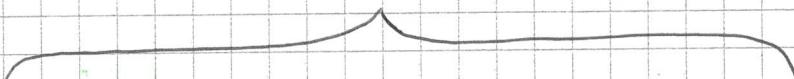
(1)

ENVIRONMENT:

- probability distribution $p(x, y)$
- unknown parameters
- data coming from training set D are i.i.d.
- D can be divided in subset $D_i \rightarrow$ classes
- for new example $x \rightarrow$ compute posterior probability

$$p(y_i | x, D) = \frac{p(x | y_i, D) p(y_i | D)}{p(x | D)}$$

\hookrightarrow trying to find parameters that describes p from D



MAXIMUM LIKELIHOOD

θ : fixed but unknown
Analysis a posteriori

$$\begin{aligned} \theta^* &= \arg \max_{\theta} p(D | \theta) \\ &= \arg \max_{\theta} \prod p(x_j | \theta) \end{aligned}$$

\hookrightarrow maximize the probability of the observed examples D_i of being part of D .

$$\log \rightarrow \nabla_{\theta} \sum_{j=1}^n \log(p(x_j | \theta)) = 0$$

\hookrightarrow Global / Local max

BAYESIAN

θ : random, known prior dist.
 \hookrightarrow update values while getting new classified data

$$p(x | y_i, D_i) = \int_{\theta} p(x_i, \theta | y_i, D_i) d\theta$$

$$\hookrightarrow p(x | D) = \int_y \int_{\theta} p(x | \theta) p(\theta | D) d\theta$$

\downarrow
Prediction of the new samples - (integral over every possible par.)

- $p(x | \theta)$: easily computed, we have the type of distribution and its parameters (e.g. Gaussian)

$$- p(\theta | D) = \frac{p(D | \theta) \cdot p(\theta)}{p(D)} \xrightarrow{\text{constant}} \text{indep. from } \theta$$

If we want the most likely value for θ we compute $p(D | \theta) p(\theta)$.

If we want the final probability we have to consider also $p(D)$

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

\downarrow Marginalization

\downarrow $p(D | \theta)$ come mai \rightarrow (integral)
la probabilità di D per tutti i valori di θ , pesata dalla probabilità a posteriori di θ

SUFFICIENT STATISTICS

Any function on a set of samples D is a sufficient statistic : $s = \Phi(D)$
If $P(D|s, \theta) = P(D|s)$ \rightarrow contiene tutto l'info su D X ottenere la stessa prob.
If θ is a random variable, a sufficient stat. s contains all the relevant information about D :

$$P(\theta|D, s) = \frac{P(D|\theta, s) P(\theta|s)}{P(D|s)}$$

Almeno i campioni non
contengono informazioni
sulla stima di θ

For a gaussian $s = \text{mean, covariance}$

CONJUGATE PRIORS

→ ESEMPI BERNOULLI and MULTINOMIAL (77-81)

Given likelihood $p(x|\theta)$ and prior $p(\theta)$, then $p(\theta)$ is a conjugate prior for $p(x|\theta)$ if the posterior distribution $p(\theta|x)$ is in the same family as the prior $p(\theta)$

CH.8 BAYESIAN NETWORKS

Bayesian Network → Directed Graphical model

- connections have direction
- node = variable
- edge = directed dependency

INDEPENDENCE ASSUMPTION: $I_B(G) = \{ \forall i: x_i \perp \text{NonDescendants}(x_i) \mid \text{Parents}(x_i) \}$

Representation of probabilistic relationships: $I_B(G) \subseteq I(p)$, set of independence assertions holding in p

The reverse is not necessarily true

To model a joint probability with m variables, we have to consider 2^m possible configurations.

↳ Try to break down into pieces → p factorizes according to G if:

STRUCTURE prob. due to causal
i.e. relations do not depend on due to existing
nella distribuzione p

$$p(x_1, \dots, x_m) = \prod_{i=1}^m p(x_i \mid \text{Parents}(x_i))$$

- ① IF G is an I-map for p , then p factorizes according to G
- ② IF p factorizes according to G , then G is an I-map for p .

PROOF ① Imap ⇒ Factorization

1. If p is an Imap satisfies $\{ \forall i: x_i \perp \text{NonDescendants}(x_i) \mid \text{Parents}(x_i) \}$
2. order variables in topological order: $x_i \rightarrow x_j \Rightarrow i < j$
3. decompose joint prob. using CHAIN RULE: $p(x_1, \dots, x_m) = \prod p(x_i \mid x_1, \dots, x_{i-1})$
4. local independences imply that for each x_i : $p(x_i \mid x_1, \dots, x_{i-1}) = p(x_i \mid \text{Parents}(x_i))$

PROOF ② Factorization ⇒ Imap

1. If p factorizes according to G , the joint probability can be written as $\prod p(x_i \mid \text{Parents}(x_i))$
2. Considering x_m , By product and sum rule: $p(x_m \mid x_1, \dots, x_{m-1}) = \frac{p(x_1, \dots, x_m)}{\sum_{x_m} p(x_1, \dots, x_m)}$
3. Applying factorization: $\frac{\prod p(x_i \mid \text{Parents}(x_i))}{\sum_{x_m} \prod p(x_i \mid \text{Parents}(x_i))} = p(x_m \mid \text{Parents}(x_m))$

\Downarrow
FACTORIZATION \Leftrightarrow IMAP

BAYESIAN NETWORK is a pair (G, p) where p factorizes over G and it is represented as a set of conditional probability distribution (CPD) associated with the nodes of G

$$p(x_1, \dots, x_m) = \prod_{i=1}^m p(x_i \mid \text{Parents}(x_i))$$

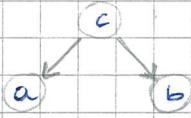
CONDITIONAL INDEPENDENCE: a,b are independent ($a \perp b \mid \emptyset$), if: $p(a, b) = p(a)p(b)$
 a,b are conditionally independent given c ($a \perp b \mid c$) if:
 $p(a, b \mid c) = p(a \mid c)p(b \mid c)$

Independence assumptions can be verified by repeated applications of sum and product rules. Graphical models allow to directly verify them through the d-separation criterion.

D - SEPARATION

TAIL - TO - TAIL

c is tail-to-tail to the path $a \rightarrow b$



$$P(a,b,c) = P(a|c) P(b|c) P(c)$$

a, b are not independent: $P(a,b) = \sum_c P(a|c) P(b|c) P(c) \neq P(a) P(b)$

a, b are conditional independent given c : $P(a,b|c) = P(a|c) P(b|c)$

HEAD - TO - TAIL

c is in the middle of a chain

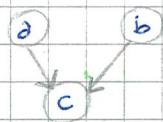


$$P(a,b,c) = P(b|c) P(c|a) P(a)$$

a, b are not independent: $P(a,b) = \sum_c P(b|c) P(c|a) \neq P(a) P(b)$

a, b are cond. ind. given c : $P(a,b|c) = P(b|c) P(c|a) P(a) = P(c) P(b|c) P(a|c)$

HEAD - TO - HEAD



$$P(a,b,c) = P(c|a,b) P(a) P(b)$$

a, b are independent: $P(a,b) = \sum_c P(c|a,b) P(a) P(b) = P(a) P(b)$

a, b are NOT condit. ind. given c

REMARK: In real world applications, we typically observe the effects rather than the causes. We have to apply the Bayes theorem in order to reason about the probability of the causes.

GENERAL HEAD - TO - HEAD

+ di 3 variabili

- let a descendant of a node x be any node which can be reached from x with a path following the direction of the arrows.
- A head-to-head node c unblocks the dependency path between its parents; if either itself or any of its descendants receives evidence.

For every descendant, though probably the further the descendant the less influent is the indirect connection.

GENERAL D - SEPARATION

- Given a generic Bayesian network
- Given A, B, C arbitrary nonintersecting nodes
- d -separated by C if all paths from any node in A to any node in B are blocked.
- A path is blocked if it includes at least s.t. either:
 - the arrows on the path meet tail-to-tail or head-to-tail at the node and it is in C or \rightarrow Non c' flusso di info
 - the arrows on the path meet head-to-head at the node and neither it nor any of its descendants is in C . \rightarrow Non c' flusso di info una volta che si è cancellato?

D - SEPARATION IMPLIES CONDITIONAL INDEPENDENCE

\hookrightarrow The sets A and B are ind. given C ($A \perp B | C$) if they're d -separated by C

BN INDEPENDENCE REVIEWED

(3)

A BN structure encodes a set of **LOCAL** independence assumptions:

$$I_e(G) = \{ \forall x_i \mid \text{NonDescendants } x_i \mid \text{Parents } x_i \}$$

i.e. Every node is cond. ind. by its non-descendants, given its parents

A BN structure G encodes a set of **GLOBAL** (Markov) independence assumptions:

$$I(G) = \{ (A \perp B \mid C), \text{dsep}(A; B \mid C) \}$$

→ "A e B sono a-separati da C"

BN EQUIVALENCE CLASSES

Different BN structures can encode the exact same set of independence assumptions.

Two BN structures G and G' are **I-equivalent** if $I(G) = I(G')$.

The space of BN structures over X is partitioned into a set of mutually exclusive and exhaustive **I-equivalence classes**.

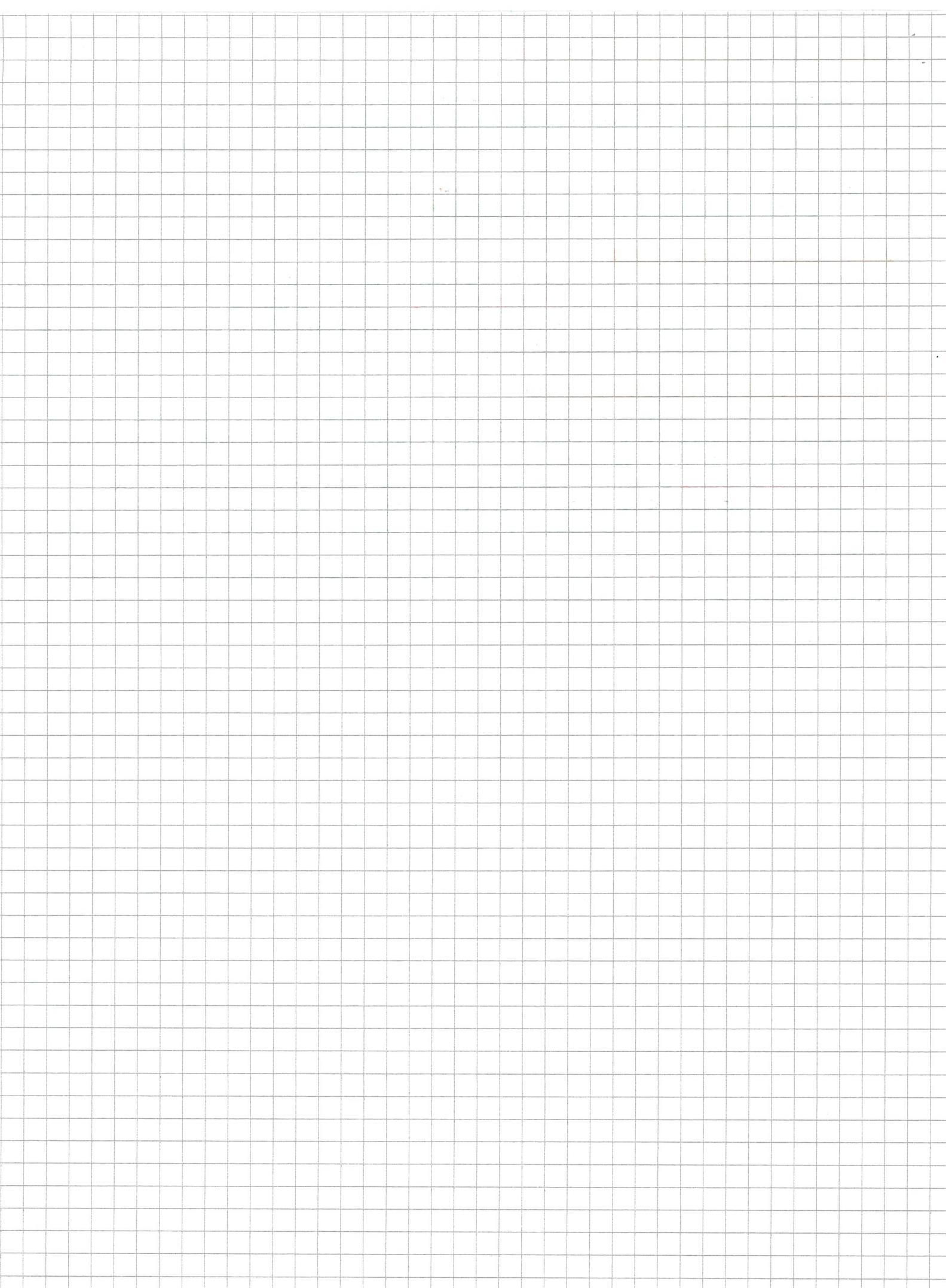
I-MAPS VS DISTRIBUTIONS

For a structure G to be an I-map for P , it does not need to encode all its independences. As a consequence a fully connected graph is an I-map for any P defined over its variables.

MINIMAL I-MAP: a min. I-map for P is an I-map G which can't be "reduced" into a $G' \subset G$ that is also an I-map for P . In other words it is no more possible to remove edges from G without introducing independences which don't hold in P .

PERFECT MAPS P-MAPS: A structure G is a perfect map (P-map) for P if it captures all its independences: $I(G) = I(P)$

PAG. 103: Suggestion for building a BN



CH 10 LEARNING BN

(4)

We have a dataset of examples $D = \{x(1), \dots, x(N)\}$. Each example $x(i)$ is a configuration for all (complete data) or some (incomplete data) variables in the model (BN).

TASK: estimate parameters of the model (CPD) from data.

SIMPLEST APPROACH: learning parameters maximizing the likelihood of data:

$$\theta^{\max} = \operatorname{argmax}_{\theta} P(D|\theta) = \operatorname{argmax}_{\theta} L(D, \theta)$$

since data are i.i.d. given θ : $P(D|\theta) = \prod_{i=1}^N P(x(i)|\theta)$

$$\text{Factorization of BN: } P(D|\theta) = \prod_{i=1}^N \prod_{j=1}^m P(x_j(i)|P_{aj}(i), \theta)$$

$P(x_j(i)|P_{aj}(i))$ depends only on $\theta_{x_j|P_{aj}}$ and not on other parameters.

$$P(D|\theta) = \prod_{i=1}^N \prod_{j=1}^m P(x_j(i)|P_{aj}(i), \theta_{x_j|P_{aj}}) \quad \text{ML ESTIMATION WITH COMPLETE DATA}$$

Parameters for each CPD can be estimated independently

$$\theta_{x_j|P_{aj}} = \operatorname{argmax}_{\theta_{x_j|P_{aj}}} \prod_{i=1}^N P(x_j(i)|P_{aj}(i), \theta_{x_j|P_{aj}}) = L(\theta_{x_j|P_{aj}}, D)$$

↳ we have to compute a maximization of this kind for every node.

→ \cup : PARENTS

A discrete CPD $P(X_i|U)$ can be represented by a table, for example:

| | | | X_2 | | | Replacing $P(x(i) P_{aj}(i))$ with $\theta_{x(i) U(i)}$, the local likelihood of a single CPD becomes: | | |
|-------|---|----------------|----------------|----------------|---|---|----------------|----------------|
| | | | R | G | B | T | $\theta_{T R}$ | $\theta_{T G}$ |
| X_1 | T | $\theta_{T R}$ | $\theta_{T G}$ | $\theta_{T B}$ | | F | $\theta_{F R}$ | $\theta_{F G}$ |
| | F | $\theta_{F R}$ | $\theta_{F G}$ | $\theta_{F B}$ | | | | |

$L(\theta_{X_i|P_a}, D) = \prod_{u \in \text{Val}(U)} \left[\prod_{x \in \text{Val}(X_i)} \theta_{x|u}^{N_{u,x}} \right]$

↓
conf. of parents

Parameters of every column can be estimated independently.

For each multinomial distribution (column), zeroing the gradient of ML:

$$\theta_{x|u}^{\max} = \frac{N_{u,x}}{\sum_x N_{u,x}} \quad \begin{cases} \text{ML parameters are the fraction of times in} \\ \text{which the specific configuration was obs.} \end{cases}$$

in data.

ML estimation tends to overfit training set → configuration not appearing in the training set will receive zero probability.

↳ combining ML with prior probability → Maximum-a-posteriori estimate:
 $\theta^{\max} = \operatorname{argmax}_{\theta} P(D|\theta) P(\theta)$ we want to maximize posterior, not likelihood.

The conjugate prior for a multinomial distribution is a Dirichlet distribution with parameters $\alpha_{x|u}$ for each possible value of x . The resulting maximum-a-posteriori estimate is:

$$\theta_{x|u}^{\max} = \frac{N_{u,x} + \alpha_{x|u}}{\sum_x (N_{u,x} + \alpha_{x|u})} \quad \begin{matrix} \text{conf. observed} \\ \text{prior} \rightarrow \text{parameters of Dirichlet distribution} \end{matrix}$$

Prior is like having observed $\alpha_{x|u}$ imaginary samples with configuration $X=x$, $U=u$

ML ESTIMATION, INCOMPLETE DATA

Some of the examples miss evidence on some of the variables.

↳ USE OF EXPECTATION-MAXIMIZATION

- Initialize parameters
- Fill-in missing data inferring them using current parameters by solving inference problem to get expected counts.
- Compute parameters maximizing likelihood of expected counts.
↳ Iterate procedure to improve the quality of parameters (until converge)

FORMALIZING:

number of training samples

E-STEP $E_p(x_{1:D}, \theta) [N_{ijk}] = \sum_{l=1}^n p(x_i(l) = x_k, p_{il}(l) = p_{ij} | X_l, \theta)$

- X_l : what we know about the l^{th} example (variable state)
- N_{ijk} : count where example X_i takes value x_k and parents p_{il} of the example X_i takes value p_{ij} .
- IF $x_i(l)$ and $p_{il}(l)$ are observed. For X_l , it's either 0 or 1.

H-STEP $\theta^* = \arg \max_{\theta} P(D_c | \theta)$, D_c is complete dataset

$$\theta_{ijk}^* = \frac{E_p(x_{1:D}, \theta) [N_{ijk}]}{\sum_{k=1}^K E_p(x_{1:D}, \theta) [N_{ijk}]}$$

↳ "completato con i valori attesi"

ML est. can be replaced by maximum a-posteriori (MAP) estimation giving:

$$\theta_{ijk}^* = \frac{d_{ijk} + E_p(x_{1:D}, \theta) [N_{ijk}]}{\sum_{k=1}^K (d_{ijk} + E_p(x_{1:D}, \theta) [N_{ijk}])} \rightarrow \text{we have a new estimate of the par. } \theta.$$

LEARNING STRUCTURE OF GRAPHICAL MODELS

If we have domain knowledge we could try to construct the graphical model. This is not always possible. Typically, we know some relationships, but not everything.

1. **CONSTRAINT-BASED**: test conditional independences on the data and construct a model satisfying them. Based on results of tests, we add connection to network.
2. **SCORE-BASED**: assign a probabilistic score to each possible structure, define a search procedure looking for the structure maximize the score.
3. **MODEL-AVERAGING**: Treat structure as random variable. We assign a prior probability to each structure, and average prediction over all possible structures weighted by their probabilities.

Quando si costruisce una BN MAP assegna la probabilità di una config. Altro non; qualsiasi futuro avviene con quella configurazione sarà probabilità zero

CH. 11 NAIVE BAYES

(5)

PROBLEM: CLASSIFICATION → we have a set of inputs } supervised learning
compute the output

IDEA: Probabilistic algorithm for classification.

$$\hookrightarrow \text{compute argmax given data: } y^* = \arg \max_{y \in Y} P(y | \vec{x}) = \arg \max_{y \in Y} \frac{P(\vec{x}|y)P(y)}{P(\vec{x})}$$

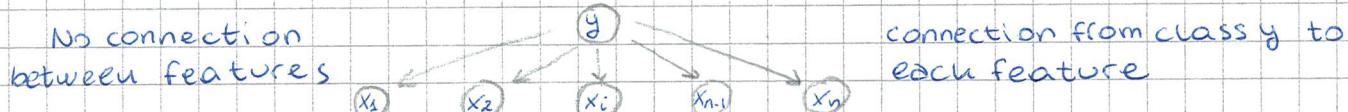
$P(\vec{x})$ is the input (a constant) so we can consider only $P(\vec{x}|y)P(y)$

PROBLEM → TO compute $P(\vec{x}|y) = P(x_1, \dots, x_n | y)$ we need 2^n configurations

SOLUTION → Naive Bayes: assume that the different features (x_1, \dots, x_n) of the input are independent given the class.

$$P(y | \vec{x}) = \prod_{i=1}^n P(x_i | y) P(y)$$

↪ consider separate table for each feature → not always accurate



TEXT CLASSIFICATION

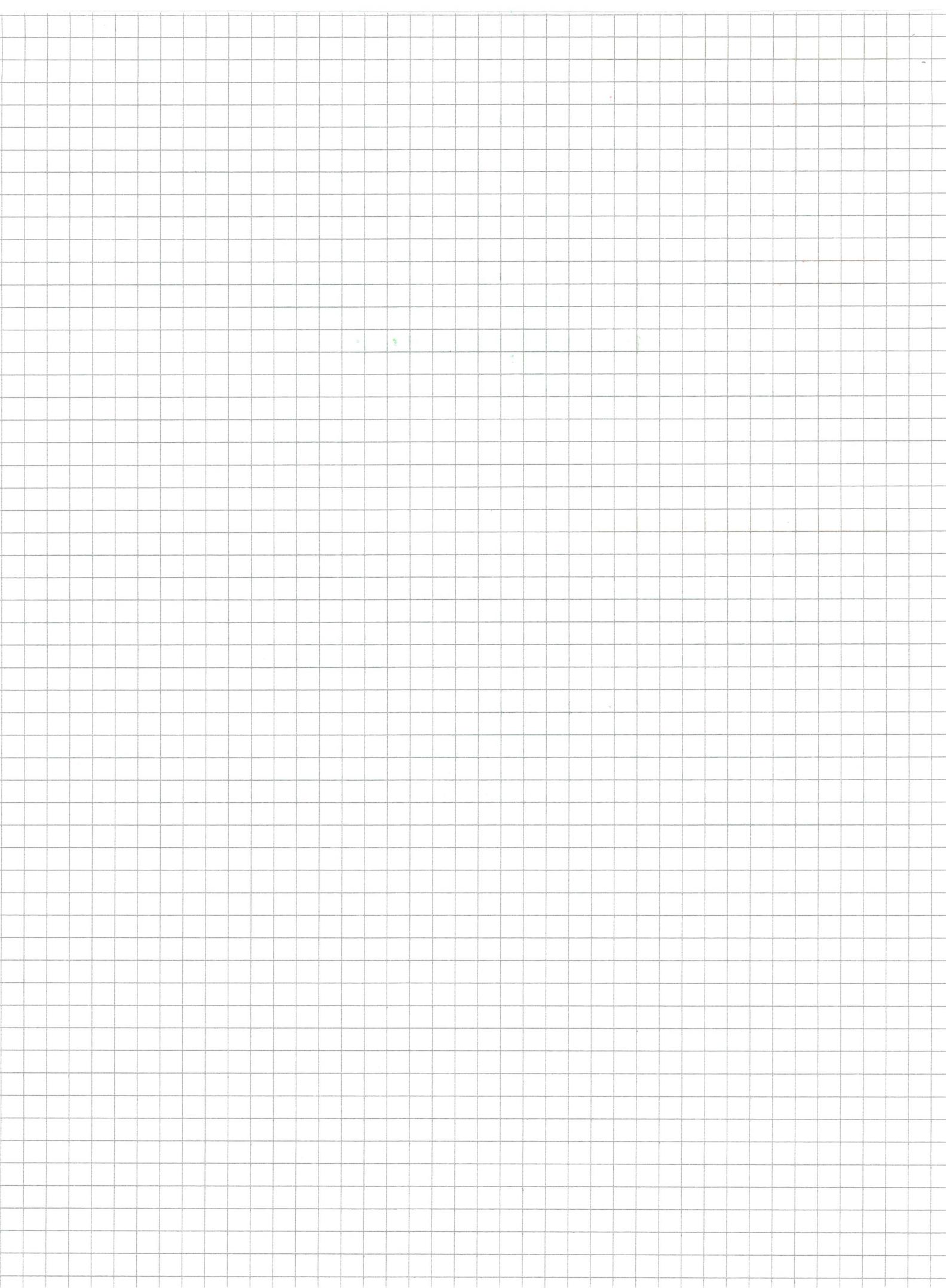
Let's consider having a piece of text x and we consider words as features

↪ $\vec{x} = (w_1, \dots, w_n)$. Probability to find a word or another is the same → ind.

To model a naive Bayes classifier is by saying that each word is modelled as independent items belonging to a greater set of words V (vocabulary). Each word can take a number of values corresponding at most with the size of the vocabulary. The Naive classifier will be:

$$\prod_{i=1}^{|V|} P(w_i | y) P(y)$$

↪ Multinomial distribution that is specific for every types of topic y , so we can get all possible words for all possible class. Each column is a multinomial distribution oversize of $|V|$ possible values.



CH. 12 LINEAR DISCRIMINANT FUNCTIONS

(6)

GENERATIVE LEARNING

Model the distribution governing the data. → (e.g. BN)

With a generative model we can generate new data (samples) according to the distribution.

DISCRIMINATIVE LEARNING

Model the discriminant function.

The purpose is to predict output given input. In classification, models the boundaries between possible classes.

- PRO: useful with complex data; learning only parameters which are interesting
- CONS: less flexible; not allow to perform inference; not possible to generate new data

A linear discriminative model is a discriminative model which use a linear function to make predictions.

The DISCRIMINANT FUNCTION is a linear combination of example features

$$f(x) = w^T x + w_0 \quad w_0 \rightarrow \text{bias / threshold}$$

In BINARY CLASSIFICATION the predicted class is obtained taking the sign of the linear function: $f(x) = \text{sign}(w^T x + w_0)$.

$f(x) = 0$ is the DECISION BOUNDARY. Weight vector w is orthogonal to the decision hyperplane. → See Proof

Functional Margin: the value of $f(x)$ for a certain point x is called Functional margin. It can be seen as a CONFIDENCE in the prediction. → The closer $f(x)$ is to zero, the smaller the confidence of the prediction.

Geometric Margin: The distance from x to the hyperplane: $r^* = \frac{|f(x)|}{\|w\|}$

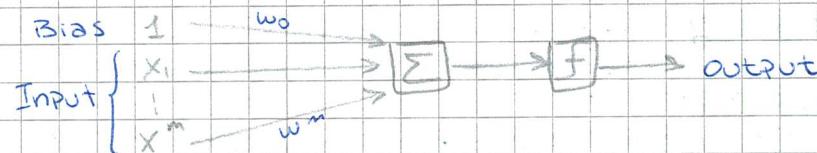
Normalized version of functional margin.

Given an arbitrary point x , it can be represented by its projection on $H(x^*)$ plus its distance to H : $x = x^* + r^* \frac{w}{\|w\|}$

PERCEPTRON

BIOLOGICAL MOTIVATION: electrochemical reactions → synapses (excite or inhibit)
→ neurons accumulates signals → exceeds threshold → signal

PERCEPTRON: Mathematical implementation of neuron structure → $f(x) = \text{sign}(w^T x + w_0)$
↳ linear combination of input features, coefficients → weights.
Result of summation is processed by an activation function giving output 0 or 1.



We can represent primitive boolean function (AND, OR, NOT), but not complex formula like XOR. We need network with 2 layers.

$$f(x) = \text{sign}(\hat{w}^T \hat{x}) \quad \hat{w} = (w_0 \quad w) \quad \hat{x} = (1 \quad x) \quad \xrightarrow{\text{Xo, bias}}$$

↳ shorter way

PARAMETER LEARNING

Function of the Parameters to be optimized \rightarrow measure of the error on the training set D \rightarrow **LOSS FUNCTION** $E(w; D)$ loss we have to pay for predicting $f(x)$ instead of y

$$E(w; D) = \sum_{(x,y) \in D} \ell(y, f(x)) \rightarrow \text{minimize } E \text{ lead to overfitting}$$

Gradient Descent common approach to minimize error function.

Progressively find a LOCAL minima of the error function.

(1) Initialize w (random)

(2) Iterate until gradient is approximately zero: $w = w - \eta \nabla E(w; D)$

$\hookrightarrow \eta$ is the **LEARNING RATE**: amount of movement at each gradient step.

\hookrightarrow Too low η implies slow converge. Too high η implies oscillations.

\hookrightarrow Techniques exists to adaptively modify η .

Perception Training Rule The error is the sum of the functional margins ($f(x)$) of incorrectly classified examples:

$$E(w; D) = \sum_{(x,y) \in D} -y f(x) s$$

$\cdot D$ is the set of current training errors $y f(x) \leq 0$

$\cdot s$ is the scaling factor to adjust the learning rate

\hookrightarrow We are taking in consideration the confidence of our predictions.

For this reason it's called **CONFIDENCE LOSS**.

If we apply gradient descent on previous equation $E(w; D)$ the amount of update is:

$$\star \quad \frac{\partial}{\partial w} E(w; D) = \eta \sum_{(x,y) \in D} y x$$

Following this approach we have to iterate over all the training set at each iteration to compute the gradient of the error for each example

\hookrightarrow SLOW CONVERGENCE

\hookrightarrow SOLUTION: **STOCHASTIC PERCEPTRON RULE**

(1) Initialize weights randomly

(2) Iterate until all examples are correctly classified

(a) For each incorrectly classified training examples (x, y) update the weight vector: $w \leftarrow w + \eta y x$

\hookrightarrow We make a gradient step for each training error. Each gradient step is very fast. At each iteration we compute the gradient on a different error function.

PERCEPTRON REGRESSION

Linear models are used also for **LINEAR REGRESSION**.

$$Xw = y \quad X: \text{training matrix} \quad \begin{cases} \text{Row} \rightarrow \text{samples} \\ \text{Column} \rightarrow \text{Features} \end{cases} \\ y: \text{output training matrix}$$

Giving solution:

$$w = X^{-1} y$$

\hookrightarrow NOT WORK X is not square

Not exact solution typically exists

We are going to minimize error function \rightarrow **MEAN SQUARED ERROR (MSE)**

$$E(w; D) = \sum_{(x,y) \in D} (y - f(x))^2 = (y - Xw)^T (y - Xw)$$

Closed form solution compute gradient and make it equal to zero $\nabla E(w; D) = 0 \Rightarrow w = (X^T X)^{-1} X^T y$ pseudo-inverse of $X^T X$

If there is an exact solution, MSE finds it, otherwise it will provide an approximation which minimizes the squared error.

With a lot of features, inverting $(X^T X)$ can be very expensive.

↳ we can compute the gradient and perform gradient descent.

For a single weight, the gradient step is computed:

For a single weight $\frac{\partial E}{\partial w_i} = \sum_{(x,y) \in D} (y - f(x)) (-x_i)$

MULTICLASS CLASSIFICATION

Solve multiclass classification by means of combination of binary classification tasks.

- ONE VS ALL positive examples \in class } classifiers needed \times ogni classe
negative examples \in class
- in decision hyperplanes, one per class
 Each classifier provides a confidence \rightarrow take class with max. confidence
- In decision boundaries $f_i(x) f_j(x)$ we don't know which class choose.
- ONE VS ONE positive examples \in class } classifiers per ogni coppia
negative examples \in other class } di classi
 $\frac{m(m-1)}{2}$ binary classifier, if we have m classes

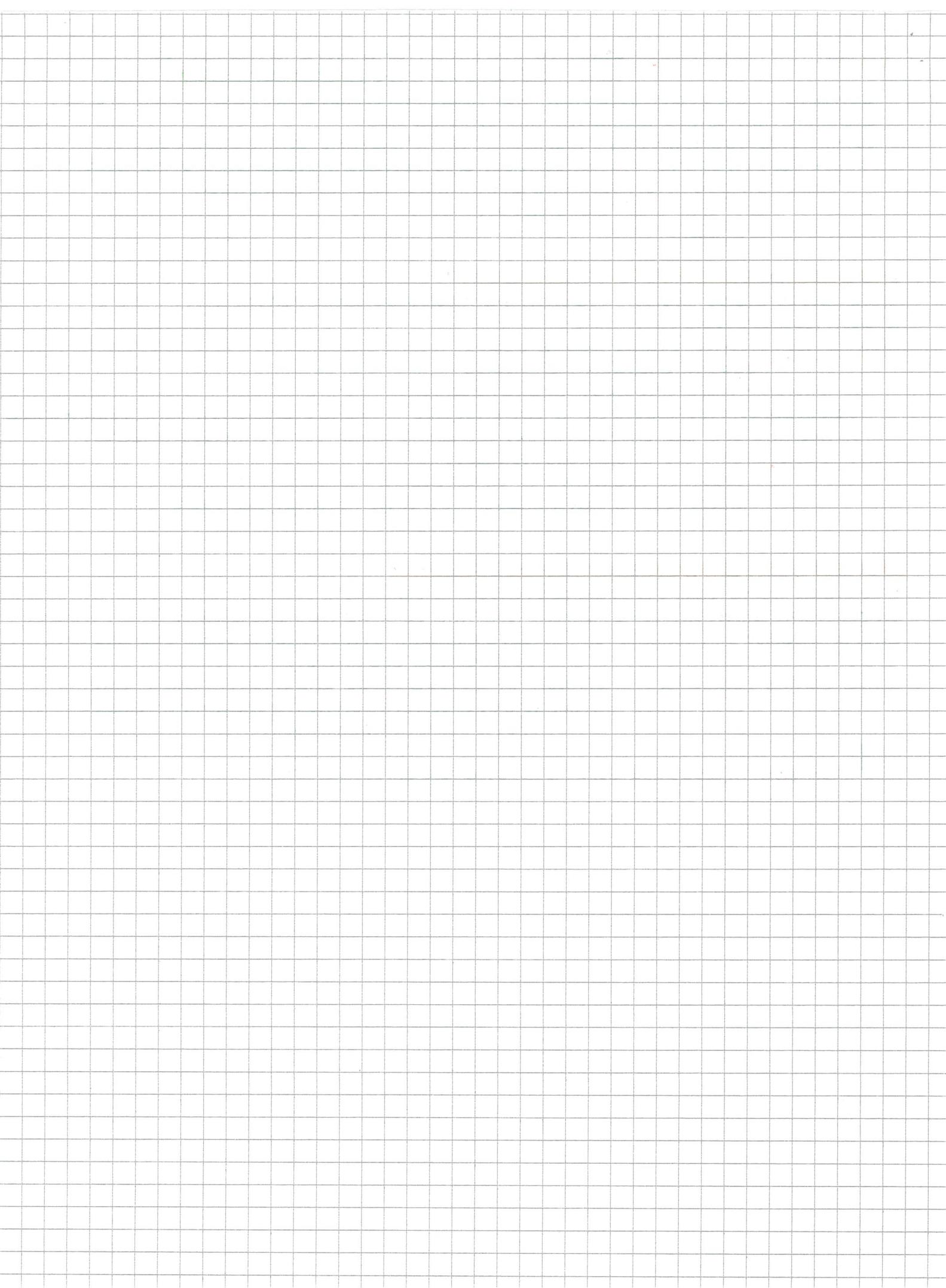
GENERATIVE LINEAR CLASSIFIER

The task is to compute the posteriori $P(y|x)$.

In gaussian distributions, linear decision boundaries are when Σ is shared among classes ($\Sigma_i = \Sigma$). \rightarrow Only under this assumption Gaussian is a linear classifier.

Naive Bayes classifier \rightarrow linear if all features are independent one with respect to the other given the class:
 $f_i(x) = P(x|y_i)P(y_i)$

↳ In both cases we can use the LOG function.



CH. 13 SUPPORT VECTOR MACHINE NO NON-LINEAR SVM / SVR REGRESSION (8)

SVM are a popular method for linear (and non-linear) classification.

- Linear classifier selecting hyperplane maximizing separation margin between classes (Large Margin Classifiers)
- Solution only depends on a small subset of training samples (Support Vectors).
- Sound generalization theory (bounds or errors based on margin).
- Can be easily extended to non linear separation (kernel machines).

Classifier margin Given a training set D , a classifier confidence margin is

$$P = \min_{(x,y) \in D} y f(x)$$

This is the minimal confidence of the classifier in a correct prediction (which corresponds to $y f(x)$ as for perceptron) and has to be positive in order to correctly separate the classes.

The Geometric Margin is the same value divided by the norm of w :

$$\frac{P}{\|w\|} = \min_{(x,y) \in D} \frac{y f(x)}{\|w\|}$$

Canonical Hyperplane there is an infinite number of equivalent hyperplanes that encodes for the same hyperplane: \rightarrow Indicate multiple solutions

$$d(w^T x + w_0) = 0 \quad d \neq 0$$

The canonical hyperplane having the confidence (classifier margin) equal to 1

$$P = \min_{(x,y) \in D} y f(x) = 1$$

and its geometric margin is:

$$\frac{P}{\|w\|} = \frac{1}{\|w\|}$$

VEDI FIG. 13.2 PAG. 149

HARD MARGIN SVM

THEOREM: Margin Error Bound

Consider a set of decision functions $f(x) = \text{sign}(w^T x)$ with $\|w\| \leq \Delta$ and $\|x\| \leq R$ for some $R, \Delta > 0$.

Let $p > 0$ and γ denote the fraction of training samples with margin smaller than $P/\|w\|$.

For all distributions P generating the data, with probability at least $1 - \delta$ (δ prob. of error) over the drawing m training patterns, and for any $p > 0$ and $\delta \in (0, 1)$, the probability that a test pattern drawn for P will be misclassified is bound from above by:

Component of all margin errors \hookrightarrow
$$(2) + \sqrt{\frac{C}{m} \left(\frac{R^2 \Delta^2}{p^2} \ln^2 m + \ln \left(\frac{1}{\delta} \right) \right)}$$
 $C = \text{universal constant}$

↳ This theorem shows a bound on the generalization error of a classifier which is trained to be an hard margin SVM maximized.

↳ The probability of test error depends on

- number of margin errors γ (samples with margin $< P/\|w\|$)
- number of training sample (error depend on $\sqrt{(\ln m)/m}$)
- size of the margin (error depends on $1/p^2$)

Learning Problem The learning objective is $\min_{w, w_0} \frac{1}{2} \|w\|^2$
subject to $y_i(w^T x_i) + w_0 \geq 1$

Hard margin is called this way because it require all the examples to have a confidence at least equal to 1. \rightarrow not always possible

↳ Quadratic optimization problem

TECHNIQUES FOR SOLVE THE PROBLEM

Karush-Kuhn-Tucker approach (KKT) convert problem into an unconstrained problem with the same solution.

Let's have a general constrained problem:

$$\min_z f(z) \text{ subject to } g_i(z) \geq 0 \quad \forall i$$

Introduce $\alpha_i \geq 0$ (Lagrange multiplier) for each constraint and rewrite

$$\min_z \max_{\alpha \geq 0} f(z) - \sum_i \alpha_i g_i(z) \quad i: \text{all possible constraints}$$

↳ Remove constraints adding it into the objective

The optimal solution(s) x^* for this problem is/are also optimal solution for the original constraint problem.

If there is at least 1 constraint that x^* doesn't satisfy this leads to a non-valid solution for the original problem.

If all constraints are satisfied, maximization over the α will set all elements of the summation to zero, so that z^* is a solution for $\min_z f(z)$.

Applying KKT to SVM we can formalize: $\min_{w, w_0} \frac{1}{2} \|w\|^2$ subject to $y_i(w^T x_i + w_0) \geq 1$

constraints can be included in minimization:

$$L(w, w_0, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i (y_i(w^T x_i + w_0) - 1)$$

↳ minimized respect to w, w_0 . Maximized respect to α

$$\nabla_w L = 0 \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \rightarrow \text{write primal variables in terms of dual } (\alpha)$$

$$\frac{\partial L}{\partial w_0} L = 0 \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0$$

$$\text{↳ substituting in Lagrangian: } L(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j x_i^T x_j$$

↳ to be maximized due to the dual variables α

The objective becomes: $\max_{\alpha} L(\alpha)$ with constraints $\alpha_i \geq 0 \quad i \in \{1, m\}$

$$\sum_{i=1}^m \alpha_i y_i = 0$$



Formulation in dual variables. Dual constraints are simpler.

↳ choice depends on number of variables

$$\text{DECISION FUNCTION: } f(x) = w^T x + w_0 = \underbrace{\sum_{i=1}^m \alpha_i y_i x_i^T x}_{\text{Primal}} + w_0 \quad \underbrace{\alpha_i y_i}_{\text{Dual}}$$

Why we call this method support vector? Lagrangian is a min-max problem (\min_w, \max_{α}), the result is a SADDLE POINT. Each of the conditions should be equal to zero, therefore each element in the summation should be zero. This happens in two ways:

1. α_i are equal to zero

2. $y_i(w^T x_i + w_0) - 1 = 0 \rightarrow$ confidence for example should be 1.

There are 2 hyperplanes with confidence = 1 \rightarrow pos. examples conf = 1
↳ neg examples conf = -1

These are the only points with $\alpha_i > 0 \rightarrow$ support vectors.

Decision hyperplane $f(x) = 0$ is defined only in terms of the support vectors.

Bias w_0 can be computed from KKT conditions: $w_0 = \frac{1 - y_i w^T x_i}{y_i}$ ↗ support vector

↳ It's usually averaged over all support vectors

SOFT MARGIN SVM

PROBLEM WITH HARD MARGIN: adding new points, previous hyperplane are not longer valid. \rightarrow OVERFITTING

SOLUTION: soft margin SVM \rightarrow maximize margin but with soft constraints: some exceptions of falsely classified samples are allowed.

We add a slack variable ξ_i for each sample in the dataset, then we will sum them according to a multiplicative factor C which is a hyperparameter:

$$\min_{w \in \mathbb{R}^n, w_0 \in \mathbb{R}, \xi \in \mathbb{R}^m} \frac{\|w\|^2}{2} + C \sum_{i=1}^m \xi_i$$

subject to: $y_i(w^\top x_i + w_0) \geq 1 - \xi_i, \xi_i \geq 0$

$\xi_i = 0 \rightarrow$ correct classification

$\xi_i > 0 \rightarrow$ miss-classification

Larger $\sum_{i=1}^m \xi_i$, larger the number of miss-classified examples

Collecting slack variables $C \sum_{i=1}^m \xi_i$ trying to combine margin maximization and penalty minimization.

C gives a trade-off between maximizing the margin and fitting examples.

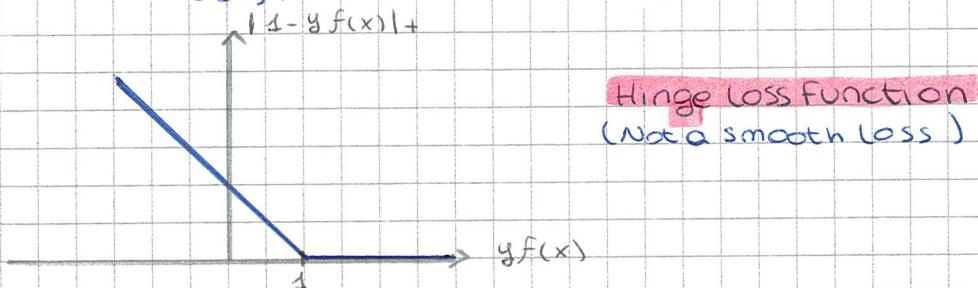
Regularization Theory: Objectives to be learned that combine a complexity term and training errors:

$$\min_{w \in \mathbb{R}^n, w_0 \in \mathbb{R}, \xi \in \mathbb{R}^m} \frac{\|w\|^2}{2} + C \sum_{i=1}^m \text{loss}(y_i, f(x_i))$$

In soft margin ξ_i is the loss of false predictions: $\xi_i = \text{loss}(y_i, f(x_i))$

ξ_i is non-negative $\rightarrow \text{loss}(y_i, f(x_i)) = |1 - y_i(w^\top x_i + w_0)| +$

\hookrightarrow PLOT:



Lagrangian For soft margin SVM

KKT formulation: $L = \frac{1}{2} \|w\|^2 + C \sum_{i=0}^m \xi_i - \sum_{i=1}^m \alpha_i (y_i(w^\top x_i + w_0) - 1 + \xi_i) - \sum_{i=0}^m \beta_i \xi_i$

Primal variable: w, w_0 Dual variables: α, β

$$\nabla_w L \Rightarrow w = \sum_{i=1}^m \alpha_i y_i x_i \quad \frac{\partial L}{\partial w_0} \Rightarrow \sum_{i=1}^m \alpha_i y_i = 0 \quad \frac{\partial L}{\partial \xi} \Rightarrow C - \alpha_i - \beta_i = 0$$

$$L(\alpha) = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j + \sum_{i=1}^m \alpha_i \Rightarrow \text{correspond exactly to hard.}$$

Dual Formulation: $\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j x_i^\top x_j$

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

At the saddle point it holds that for all i : $\alpha_i (y_i(w^\top x_i + w_0) - 1 + \xi_i) = 0$
 $\beta_i \xi_i = 0$

Support vectors are examples for which $(y_i(w^\top x_i + w_0)) \leq 1$

STOCHASTIC GRADIENT DESCENT

Compute gradient on a single example: $E(w, (x_i, y_i)) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{2} (y_i - \langle w, x_i \rangle)^2$

Whenever we don't have a single gradient in a point, we can perform a subgradient. Its indicator function will be:

$$\mathbb{I}[y_i < \langle w, x_i \rangle < 1] = \begin{cases} 1 & \text{if } y_i < \langle w, x_i \rangle < 1 \\ 0 & \text{otherwise} \end{cases}$$

The subgradient of a function f in a point x_0 is any vector v such that for any x hold the condition: $f(x) - f(x_0) \geq v^T (x - x_0)$

LARGE SCALE SVM LEARNING

PSEUDOCODE

- (1) Initialize $w_1 = 0$
- (2) For $t=1$ to T :
 - a. Randomly choose (x_{it}, y_{it}) from D
 - b. set $\gamma_t = \frac{1}{\lambda t}$
 - c. Update w : $w_{t+1} = w_t - \gamma_t \nabla_w E(w; (x_{it}; y_{it}))$
- (3) Return w_{T+1}

CH. 15 KERNEL MACHINES

PROBLEM: with SVM we can deal non-linear problem mapping in non-linear features as combination of features. It becomes computationally unfeasible in high dimensional cases.

SOLUTION: kernel trick

Kernel Trick: replacing the dot products that appear in the dual formulation of non-linear SVMs with an equivalent kernel function:

$$K(x, x') = \Phi(x)^\top \Phi(x')$$

The kernel function uses example in input space (not feature). We can build a function that can compute this without to explicitly build the map.

The dual optimization problem is: $\max_{\alpha \in \mathbb{R}^m} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i y_i \alpha_j y_j \Phi(x_i)^\top \Phi(x_j)$

$$\text{subject to: } \alpha_i \geq 0, \sum_{i=1}^m \alpha_i y_i = 0$$

Dual decision function: $f(x) = \sum_{i=1}^m \alpha_i y_i \Phi(x_i)^\top \Phi(x) = \sum_{i=1}^m \alpha_i y_i K(x_i, x)$

Feature map $\Phi: X \rightarrow \mathcal{H}$ Φ is a function mapping each example to an higher dimensional space \mathcal{H}

Examples x are replaced with their feature mapping $\Phi(x)$.

Maps all possible conjunctions of features {of a certain degree d } \rightarrow HOMOGENEOUS
up to a certain degree \rightarrow INHOMOGENEOUS

what are kernels?

HOMOGENEOUS: $K(x, x') = (x^\top x')^d \rightarrow$ result is a scalar.

A kernel it always correspond to a dot product and some feat. space.

INHOMOGENEOUS: $K(x, x') = (1 + x^\top x')^d$

KERNELS A valid kernel what works properly, is a function defined over the cartesian product of the input space: $K: X \times X \rightarrow \mathbb{R}$ that corresponds to a dot product in SOME feature space $K(x, x') = \Phi(x)^\top \Phi(x')$, kernels can be seen as similarities between objects.

VALIDITY OF A KERNEL

If we built learning systems on non-vector objects we have to invent the kernel ourselves, therefore we need to verify that is a valid one.

GRAM MATRIX: given examples $\{x_1, \dots, x_m\}$ and a kernel function K , the Gram matrix K is the symmetric matrix of pairwise kernels between examples: $K_{ij} = K(x_i, x_j) \quad \forall i, j \rightarrow$ symmetric.

Is a positive definite if $\sum_{i,j=1}^m c_i c_j K_{ij} \geq 0 \quad \forall c \in \mathbb{R}^m$

All eigenvectors are non-negative.

To show validity check the satisfaction of at least one of these conditions:
(1) prove its positive definiteness (difficult)

(2) find corresponding feature map: explication of Φ dot-product combination by making the feature map explicit

(3) use kernel combination properties in order to build a new kernel, the operation that we can perform on kernels will preserve their properties.

Kernelizing different SVMs

Example support vector regression:

$$\text{Dual Problem: } \max_{\alpha \in \mathbb{R}^m} -\frac{1}{2} \sum_{i,j=1}^m (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) \Phi(x_i)^T \Phi(x_j) - \sum_{i=1}^m (\alpha_i^* + \alpha_i) + \sum_{i=1}^m y_i(\alpha_i^* - \alpha_i)$$

$$\text{subject to } \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0$$

$$\text{Regression function: } f(x) = w^T \Phi(x) + w_0 = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \Phi(x_i)^T \Phi(x) + w_0$$

Kernelizing a perceptron

Initialize $\alpha^{(0)} = 0$, the coefficients to be learnt

Iterate through the examples, and for each miss-classified example we perform an update: $\alpha_i^{(t+1)} \leftarrow \alpha_i^{(t)} + \gamma y_i$

Classification function becomes $f(x) = \sum_{i=1}^m \alpha_i K(x_i, x)$

STOCHASTIC PERCEPTRON: $f(x) = w^T x = \sum_{i=1}^d w_i x_i + w_0$, update $w^{(t+1)} \leftarrow w^{(t)} + \gamma y_i x_i$

TYPES OF KERNELS

- Linear $K(x, x') = x^T x'$ Inverso della similitudine
- Polynomial $K_{c,d}(x, x') = (x^T x' + c)^d$ IF $x^T x' = 0 \Rightarrow K(x, x') = 1$
- Gaussian: $K_g(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right)$ x' is the mean, x sample depends on width parameter σ

Gaussian kernel can uniformly approximate any arbitrary continuous function

KERNELS ON STRUCTURED DATA

Kernels are generalization of dot products to arbitrary domains.

It's possible to design kernels over structured data like sequences, trees or graphs. Idea is designing a pairwise function measuring the similarity of two objects. This measure has to satisfy the p.d. conditions to be a valid kernel.

MATCH (DELTA) KERNEL: $K_s(x, x') = S(x, x') = \begin{cases} 1 & \text{if } x = x' \\ 0 & \text{otherwise} \end{cases}$ simplest kernel on struct. checks equality.

SPECTRUM KERNEL: feature space is space of all possible k -grams (subsequences)

An efficient procedure based on suffix trees allows to compute kernel without explicitly building feature maps.

Kernel combination. Build a kernel by combining pieces means build a valid kernel by following some rules that ensure validity:

- simpler kernels can be combined using certain operation (+, \times , ...)
- allow to combine simpler kernels to design complex kernel
- correctly using combination operators guarantees that complex kernels are preserved.

SUMMATION: $(K_1 + K_2)(x, x') = [\Phi_1(x) \Phi_2(x)] \begin{pmatrix} \Phi_1(x') \\ \Phi_2(x') \end{pmatrix}$

MULTIPLICATION: $(K_1 \times K_2)(x, x') = \Phi_{12}(x)^T \Phi_{21}(x')$

LINEAR COMBINATION: $K_{\text{sum}}(x, x') = \sum_{k=1}^K \beta_k K_k(x, x')$

The weights of linear combination can be learned simultaneously to the predictor weights (α) → Allows to select which kernel is more useful

NORMALIZATION: $\hat{K}(x, x') = \frac{K(x, x')}{\sqrt{K(x, x) K(x', x')}}$ cosine normalization

COMPOSITION: $(K \circ \phi_K)(x, x') = \exp\left(-\frac{K(x, x) - 2K(x, x') + K(x', x')}{2\sigma^2}\right)$

Is project in higher dimensional space

KERNELS ON GRAPHS

(11)

WEISFELER-LEHMAN GRAPH KERNEL: efficient graph kernel for large graph.
WL ISOMORPHISM TEST. \rightarrow If it says that 2 graphs are different, then they are surely diff.
 Given $G = (V, E)$ and $G' = (V', E')$, with $n = |V| = |V'|$. Let $I(G) = \{I(v) \mid v \in V\}$ be the set of labels in G , and let $L(G) = L(G')$. Let $\text{label}(s)$ be a function assigning a unique label to a string.

- Set $I_0(v) = I(v)$ for all v .
- For $i \in [1, n-1]$

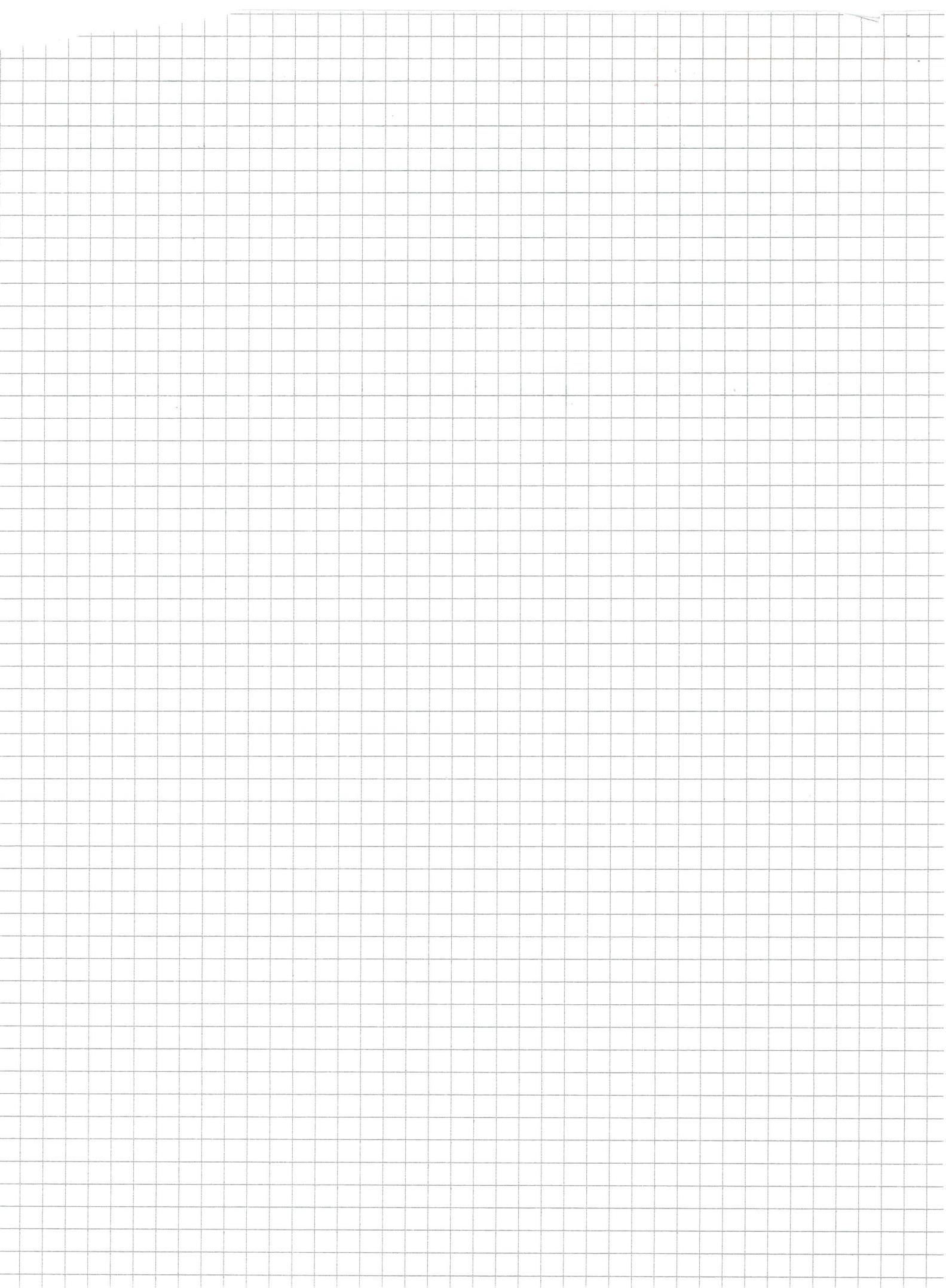
- (1) For each node v in G and G'
- (2) let $M_i(v) = \{I_{i-1}(u) \mid u \in \text{neigh}(v)\}$
- (3) concatenate the sorted labels of $M_i(v)$ into $S_i(v)$
- (4) let $I_i(v) = \text{label}(I_{i-1}(v) \circ S_i(v))$
- (5) IF $L_i(G) \neq L_i(G')$

RETURN FAIL

- RETURN PASS

WEISFELER-LEHMAN graph KERNEL: $K_{WL}^{(n)}(G, G') = \sum_{i=0}^n \kappa(G^{(i)}, G'^{(i)})$

VEDI SLIDE 29-30 x disegni



CH. 17 ENSEMBLE METHODS

Combine multiple ML models can produce better predictions than those of any single model. Training ensembles can be parallelized, with substantial computational savings.

- ENSEMBLING STRATEGIES:
- bagging → diversify the training sets
 - stacking → diversify the models being trained
 - boosting → diversify the importance of examples

Bagging: bootstrap aggregation

Take a learning algorithm A . Extract m different datasets $D^{(i)}$ from the original training set D . Train one base model per dataset $f^{(i)} = A(D^{(i)})$. Combine predictions of different base models $\hat{y} = \text{COMBINE}(f^{(1)}(x), \dots, f^{(m)}(x))$.

Simply partitioning D into m subsets produces very small training sets.

Bootstrap resampling extracts $N = |D|$ samples from D with replacement. Repeating the procedure m times to get diverse datasets of the same size as D .



Each example has probability $(1 - 1/N)$ of not being selected at each draw.

Each example has probability $(1 - 1/N)^N$ of not being selected after N draws.

For large enough N , this is 37% → 37% of the dataset are not part of a given training set. → Out-Of-Bag instances (OOB) → Used for test set

COMBINATION STRATEGIES

- Majority voting: predict the class with most votes $\hat{y} = \text{argmax}_y \sum_{i=1}^m S(y, \hat{y}^{(i)})$
- Soft voting: predict the class with the largest sum of predicted probability. Assumes base classifier outputs a confidence/probability $\hat{y} = \text{argmax}_y \sum_{i=1}^m f_i^{(i)}(x)$
- Mean: predict the mean of the base model outputs $\hat{y} = \sum_{i=1}^m f_i^{(i)}(x)$

Random Forests

use decision trees as the base models.
Even if $D^{(i)} \neq D^{(j)}$, the learned DTs can be too much correlated. Introduce additional stochasticity in the DT training process. At each node, choose the best feature to split from a random selection of the set of features.

STACKING

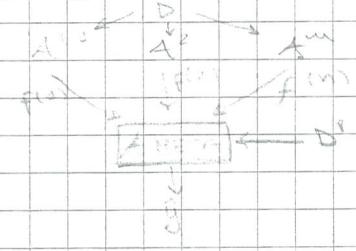
Train m base models $F^{(i)}$ on D with m different learning algorithms $A^{(i)}$
 $f^{(i)} = A^{(i)}(D)$

Use a meta-learner to learn to combine base models

$$g = \text{META}([f^{(1)}, \dots, f^{(m)}], D)$$

Use the meta-model to output ensemble predictions

$$\hat{y} = g([f^{(1)}(x), \dots, f^{(m)}(x)])$$



BOOSTING

Take a learner A and train it on D .

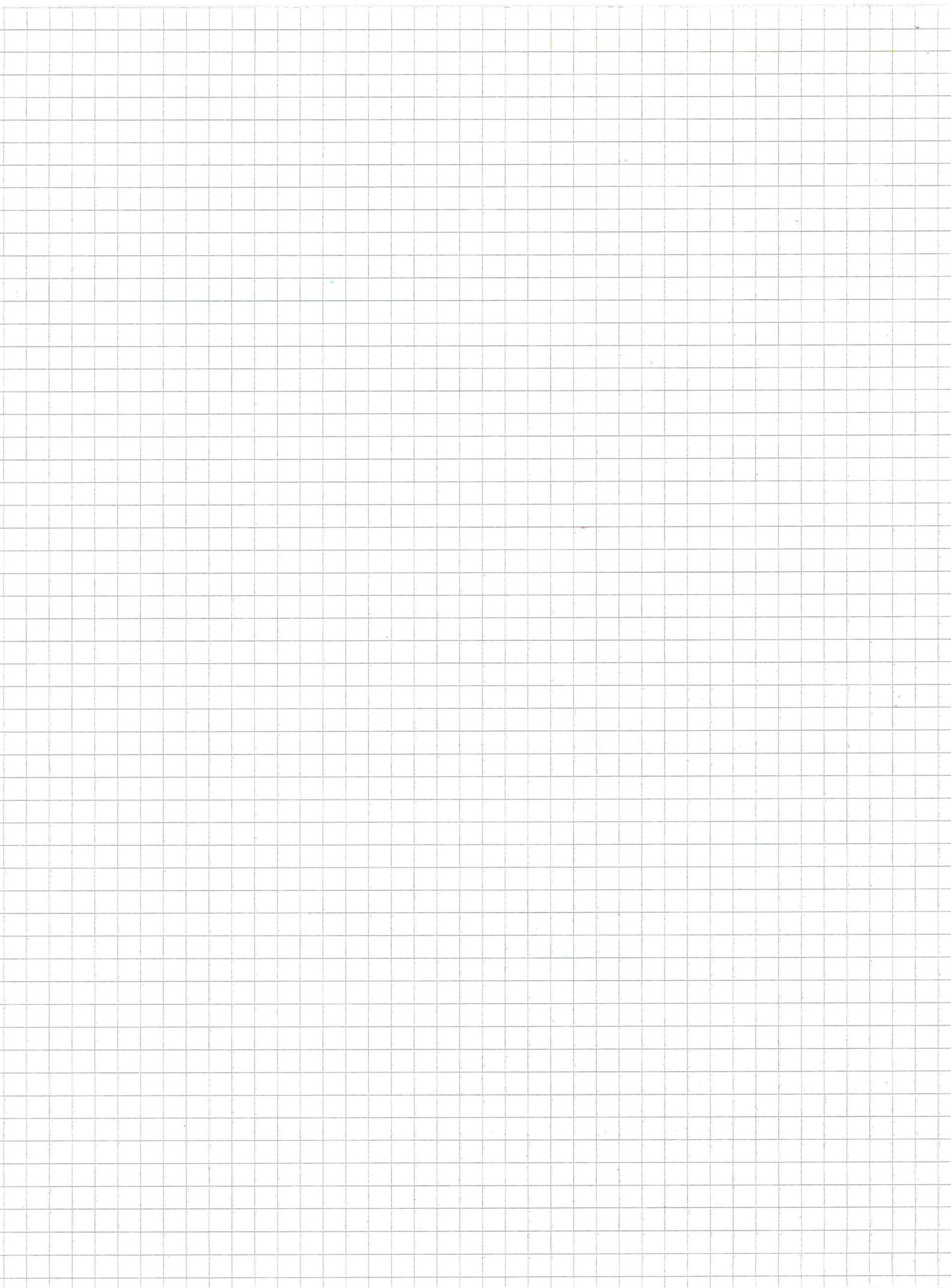
Reweight examples in D based on their accuracy according to the trained model (harder examples get larger weights).

Train A again on the reweighted dataset.

Repeat the procedure m times.

Combine the learned models into the final model.

A weak learner learns models with accuracy slightly better than random. A weak learner is easy to implement and train. Applying boosting with weak learners as the base learning algorithm allows to turn them into strong learners.



CH. 18 UNSUPERVISED LEARNING

(13)

We don't have any labels, only data.

CLUSTERING

Decide a prior the number of clusters. Each cluster is represented by its mean μ_i . The idea is to assign example to the clusters with the closer mean. CLUSTER SAMPLE MEAN: $\mu_i = \frac{1}{n_i} \sum_{x \in C_i} x$

(1) Initialize cluster mean μ_1, \dots, μ_n (random)

(2) Iterate until no mean changes:

a. Assign each example to cluster with nearest mean

b. Update cluster means according to assigned examples

Distance Metrics → To compute the nearest mean

(1) Standard euclidean distance in \mathbb{R}^d : $d(x, x') = \sqrt{\sum_{i=1}^d (x_i - x'_i)^2}$

(2) Generic Minkowski metric for $p \geq 1$: $d(x, x') = (\sum_{i=1}^d |x_i - x'_i|^p)^{1/p}$

(3) Cosine similarity: $s(x, x') = \frac{x \cdot x'}{\|x\| \|x'\|}$, cosine of the angle between 2 vectors

Quality of clustering: a criteria for defining a quality of clusters is Sum-of-squared error criterion. This method tells how bad the approximation of the cluster is using the means.

For each cluster, we take the mean μ_i within each cluster. I compare all examples in the cluster with the mean of the cluster they are in. The error will be computed as the squared error: $\mu_i = \frac{1}{n_i} \sum_{x \in C_i} x$

The sum-of-squared errors is defined as the sum over the totality of clusters:

$$E = \sum_{i=1}^k \sum_{x \in C_i} \|x - \mu_i\|^2 \quad k \text{ number of clusters (?)}$$

GAUSSIAN MIXTURE MODEL (GMM)

Each cluster is represented as a Gaussian distribution. Need to estimate the mean and possibly the variance of the Gaussian distribution. PAGE 233 x FIGURE

Parameter estimation: ML estimation cannot be applied as cluster assignment of examples is unknown. Expectation-Maximization approach:

(1) Compute expected cluster assignment given current parameter setting.

(2) Estimate parameters given cluster assignment

(3) Iterate

Example: estimating means of k univariate Gaussians

SETTING

- A dataset of x_1, \dots, x_n examples is observed
- For each example x_i , cluster assignment is modelled as z_{ij}, \dots, z_{ik} binary latent variables. $z_{ij} = 1$ if Gaussian j generated x_i , 0 otherwise.
- $z_{ij} = 1$ if Gaussian j generated x_i , 0 otherwise.
- Parameters to be estimate are the μ_1, \dots, μ_k Gaussians means.
- All Gaussians are assumed to have the same variance σ^2

ALGORITHM

(1) Initialize $h = \langle \mu_1, \dots, \mu_k \rangle$ h = hypothesis \rightarrow unknown must be estimate

(2) Iterate until difference in ML is below a certain threshold:

E-STEP: calculate expected value $E[z_{ij}]$ of each latent variable assuming current hypothesis $h = \langle \mu_1, \dots, \mu_k \rangle$ holds.

M-STEP: calculate a new ML hypothesis $h' = \langle \mu'_1, \dots, \mu'_k \rangle$ assuming values of latent variables are their expected values just computed. Replace $h \leftarrow h'$

Probability that x_i is generated by

$$E\text{-STEP: } E[z_{ij}] = \frac{P(x_i | \mu_j)}{\sum_{l=1}^k P(x_i | \mu_l)}$$

Gaussian j assuming hypothesis h

Weighted sample mean

$$M\text{-STEP: } \mu'_j = \frac{\sum_{i=1}^n E[z_{ij}] x_i}{\sum_{i=1}^n E[z_{ij}]}$$

We are given a dataset made of an observed part X and an unobserved part Z . We wish to estimate the hypothesis maximizing the expected log-likelihood for the data, with expectation taken over unobserved data:

$$h^* = \operatorname{argmax}_h E_Z[\ln p(X, Z|h)]$$

The unobserved data Z should be treated as random variables governed by the distribution depending on X and h .

Generic Algorithm

- (1) Initialize hypothesis h
- (2) Iterate until converge:

E-STEP: compute the expected likelihood of an hypothesis h' for the full data, where the unobserved data distribution is modelled according to the current hypothesis h and the obs data.

$$Q(h', h) = E_Z[\ln p(X, Z|h') | h, X]$$

M-STEP: replace the current hypothesis with the one maximizing $Q(h', h)$:

$$h' \leftarrow \operatorname{argmax}_{h'} Q(h', h)$$

HOW TO CHOOSE THE NUMBER OF CLUSTERS?

ELBOW METHOD: Increasing number of clusters allows for better modeling of data. Needs to trade-off quality of clusters with quantity. Stop increasing number of clusters when advantage is limited. → can be ambiguous

AVERAGE SILHOUETTE METHOD: use quality metric that trade-off intra-cluster similarity and inter-cluster dissimilarity.

compute average dissimilarity between i and its cluster

$$a_i = \bar{d}(i, c) = \frac{1}{|c|} \sum_{j \in c} d(i, j)$$

Compute average dissimilarity between i and each cluster $c' \neq c$, take the minimum:

$$b_i = \min_{c' \neq c} d(i, c')$$

The silhouette coefficient is:

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)}$$

Typically only 1 maximum

HIERARCHICAL CLUSTERING

Clustering doesn't need to be flat. Natural grouping of data is often hierarchical. A hierarchy of clusters can be built on examples:

TOP-DOWN → recursively split clusters into subclusters

BOTTOM-UP → recursively aggregate pairs of clusters

ALGORITHM: AGGLOMERATIVE HIERARCHICAL CLUSTERING (Bottom-up)

- (1) INITIALIZE: final number cluster K

Initial cluster number $\hat{K} = n$

Initial clusters $D_i = \{x_j\} \quad j \in [1, n]$

- (2) WHILE $\hat{K} > K$: Find pairwise nearest clusters D_i, D_j

merge D_i and D_j

update $\hat{K} = \hat{K} - 1$

STOPPING CRITERION can be threshold on pairwise similarity.

ALGORITHM STEPWISE OPTIMAL HIERARCHICAL CLUSTERING

- (1) INITIALIZE: final cluster number K

Initial cluster number $\hat{K} = n$

Initial cluster $D_i = \{x_i\} \quad i \in [1, n]$

- (2) WHILE $\hat{K} > K$: find best clusters D_i, D_j to merge according to evaluation criteria

merge D_i and D_j

update $\hat{K} = \hat{K} - 1$

CH. 19 REINFORCEMENT LEARNING

14

The learner is provided a set of possible states S , and for each state, a set of possible actions, a moving it to a next state.

In performing action a from state s , the learner is provided an immediate reward $r(s, a)$. The task is to learn a policy allowing to choose for each state s the action a maximizing the overall reward.

The learner has to deal with problems of delayed reward coming from future moves, and trade-off between exploitation and exploration.

Typical applications include moving policies for robots and sequential scheduling problems in general.

Sequential Decision Making: an agent needs to take a sequence of decisions. The agent should maximize some utility function. There is uncertainty in the result of a decision.

Markov Decision Process (MDP)

A set of states S in which the agent can be at each time instant.

A possibly empty set of terminal states $S_T \subset S$

A set of actions A the agent can make

A transition model providing the probability of going to a state s' with action a from state s : $P(s'|s, a) \quad s, s' \in S, a \in A$ P is the model

A reward $R(s, a, s')$ for making action a in state s and reaching state s' .

→ The probability $P(s'|s, a)$ don't depends on past actions or states, only on s

DEFINING UTILITIES

An environment history is a sequence of states. Utility (U): function defining the total amount of reward of an environment history.

We assume stationary preferences. There are 2 ways to define utilities:

- ADDITIVE REWARDS: $U([s_0, s_1, s_2]) = R(s_0) + R(s_1) + R(s_2) + \dots$
 - DISCOUNTED REWARDS: $U([s_0, s_1, s_2]) = R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots \quad \gamma \in [0, 1]$
- * we consider rewards that only depend on the (destination) state
- * lower discount factor γ , less important the future rewards

MDP: TAKING DECISIONS → Optimal Policy

A policy π is a full specification of what action to take at each state. The expected utility of a policy is the utility of an environment history, taken in expectation over all possible histories generated with that policy. An optimal policy π^* is a policy maximizing expected utility. For infinite horizons, optimal policies are stationary, i.e. they only depend on the current state.

Optimal policies varying π : utility is made with additive rewards.

↳ Determined by the moving reward r

Arrows indicate the best action to take.

Star indicates all actions are equally optimal.

MOVING EXPENSIVE → optimal policy is to reach any terminal state asap.

↳ **CHEAP** → optimal policy is avoiding the bad terminal state at all costs.

↳ **GIVES POSITIVE REWARDS** → optimal policy is to stay away of terminal states

Utility of states. The utility of a state given a policy π is

$$U^\pi(s) = E_\pi \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid S_0 = s \right]$$

where s_t is the state reached after t steps using policy π starting from $S_0 = s$. The true utility of a state is its utility under an optimal policy

$$U(s) = U^{\pi^*}(s) = \max_\pi U^\pi(s)$$

Given the true utility, an optimal policy is as follows:

$$\pi^*(s) = \operatorname{argmax}_{a \in A} \sum_{s' \in S} P(s'|s, a) U(s')$$

- The utility of a state is its immediate reward plus the expected discounted utility of the next state, assuming that the agent chooses an optimal action.

Bellman Equation $U(s) = R(s) + \gamma * \max_{a \in A} \sum_{s' \in S} P(s'|s, a) U(s')$

There is a Bellman equation for each state $s \in S$. Utilities of states are solutions of the set of Bellman equations.

The solutions to the set of Bellman equations are unique.

Directly solving the set of equations is hard.

Computing an optimal policy:

(1) Initialize $U_0(s)$ to zero for all s

(2) Repeat

 a. Do Bellman update for each state s : $U_{i+1}(s) \leftarrow R(s) + \gamma * \max_{a \in A} \sum_{s' \in S} P(s'|s, a) U_i(s')$

 b. $i \leftarrow i + 1$

(3) Until max utility difference below a threshold

(4) return U

The optimal policy can be set as: $\pi^*(s) = \arg\max_{a \in A} \sum_{s' \in S} P(s'|s, a) U(s')$

POLICY ITERATION:

(1) Initialize π to randomly

 Action a in state s is prescribed by policy $\pi^{(t)}$

(2) Repeat

 a. Policy evaluation $U_i(s) = R(s) + \gamma \sum_{s' \in S} P(s'|s, \pi_i(s)) U_i(s')$

 Instead of maximizing w.r.t. all $a \in A$

 b. Policy improvement $\pi_{i+1}(s) \leftarrow \arg\max_{a \in A} \sum_{s' \in S} P(s'|s, a) U_i(s')$

 c. $i \leftarrow i + 1$

(3) Until no policy improvement

(4) Return π

PARTIAL KNOWLEDGE

Value iteration and policy iteration assume perfect knowledge.

In most cases, some of these aspects are not known.

Reinforcement learning aims at learning policies by space exploration.

Policy EVALUATION: Policy is given, environment is learned.

Policy IMPROVEMENT: both policy and environment are learned.

ADAPTIVE DYNAMIC PROGRAMMING (ADP): Policy evaluation

Loop

a. Initialize $P(s'|s, a) = 0$ Also s times a is taken in S

b. Repeat: - receive award r , set $R(s) = r$ Also s times a is taken in S , leading

 - choose next action $a' \leftarrow \pi(s)$ to S'

 - take action a' , reach step s'

 - update counts $N_{sa} \leftarrow N_{sa} + 1$; $N_{s'a} \leftarrow N_{s'a} + 1$

 - update transition model $P(s''|s, a) \leftarrow N_{s''a} / N_{sa}$

 - update utility estimate $U \leftarrow \text{PolicyEvaluation}(r, U, P, R, \gamma)$

c. Until s is terminal

The algorithm performs ML estimation of transition probabilities.

Standard policy evaluation to update utility estimate.

Each step is expensive because runs Policy evaluation.

POLICY EVALUATION IN UNKNOWN ENVIRONMENT

TEMPORAL-DIFFERENCE (TD) POLICY EVALUATION: rationale

Avoid running policy evaluation at each iteration. Locally update utility.

If transition from s to s' is observed:

- if s' was always the successor of s , the utility of s should be $U(s) = R(s) + \gamma U(s')$
- The temporal-difference update rule updates the utility to get closer to that situation: $U(s) \leftarrow U(s) + \alpha (R(s) + \gamma U(s') - U(s))$ α is learning rate.

No need for a transition model for utility update. Each step is much faster than ADP. Takes longer to converge.

GREEDY AGENT:

Policy learning requires combining learning the environment and learning the optimal policy for the environment. A simple option consists of replacing policy evaluation in ADP with optimal policy computation

$$U(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s'|s,a) U(s')$$

- Greedy agents always follow the local optimal action \rightarrow lack of exploration.
- The knowledge of environment is incomplete \rightarrow suboptimal policy.

LEARNING OPTIMAL POLICY

EXPLOITATION: consists in following promising directions given current knowledge.

EXPLORATION: consists in trying novel directions looking for better alternatives.

A reasonable trade-off should be used in defining the search scheme:

ϵ -greedy strategy: choose a random move with probability ϵ , be greedy otherwise.

Assign higher utility estimates to unexplored state-action pairs:

EXPLORATION FUNCTION: $U^+(s) = R(s) + \gamma \max_{a \in A} \sum_{s' \in S} P(s'|s,a) U^+(s')$

$f(u; n)$ trades off between greed (high values of u) and curiosity (high value of n). Idea is to give more weight to actions that the agent has not tried very often.

TD LEARNING \rightarrow Learning utilities of actions

- TD Policy evaluation can also be adapted to learn an optimal policy.
- If TD is used to learn a state utility function, it needs to estimate a transition model P to derive a policy.
- TD can instead be applied to learn a state-action utility function $Q(s,a)$
- The optimal policy corresponds to: $\pi^*(s) = \arg \max_{a \in A} Q(s,a)$

ON/OFF POLICY TD LEARNING

SARSA: on-policy TD learning

LOOP

1. Initialize $Q(s,a) = 0 \quad \forall s \in S, a \in A$

2. Repeat

a. Receive award r

b. Choose next action $a \leftarrow \pi^*(s)$

c. Take action a , reach step s'

d. Choose action $a' \leftarrow \pi^*(s')$

e. Update local utility estimate $Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma Q(s',a') - Q(s,a))$

3. Until s is terminal

π^* is an ϵ -greedy policy based on Q .

In Q-Learning off-policy TD learning eliminate action d. and utility estimate is $Q(s,a) \leftarrow Q(s,a) + \alpha (r + \gamma \max_{a' \in A} Q(s',a') - Q(s,a))$

\rightarrow Greedy policy's action

Sarsa vs Q-Learning

- Sarsa is on-policy: it updates Q using the current policy's action. Tends to converge faster, and are easier to use for linear function approximators.
- Q-Learning is off-policy: it updates Q using the greedy policy's action. This method is more flexible, it can even learn from traces generated with an unknown policy.

SCALING TO LARGE STATE SPACES

FUNCTION APPROXIMATION:

All techniques seen so far assume a tabular representation of utility functions. Tabular representations don't scale to large state spaces.

The solution is to rely on function approximation: approximate $U(s)$ or $Q(s, a)$ with a parameterized function.

The function takes a state representation as input.

The function allows to generalize to unseen states.

TD Learning: state utility

$$\text{TD error: } E(s, s') = \frac{1}{2} (R(s) + \gamma U_\theta(s') - U_\theta(s))^2 \quad U_\theta \approx U(s) \text{ utility approximation}$$

$$\text{Error gradient wrt function parameters: } \nabla_\theta E(s, s') = (R(s) + \gamma U_\theta(s') - U_\theta(s)) (-\nabla_\theta U_\theta(s))$$

$$\text{Stochastic gradient update rule: } \theta = \theta - \alpha \nabla_\theta E(s, s')$$

TD learning: action utility (Q-learning)

$$\text{TD error: } E(s, a, s') = \frac{1}{2} (R(s) + \gamma \max_{a'} Q_\theta(s', a') - Q_\theta(s, a))^2 \quad Q_\theta(s, a) \approx Q(s, a)$$

$$\text{Error gradient wrt function parameters: } \nabla_\theta E(s, a, s') = (R(s) + \gamma \max_{a' \in A} Q_\theta(s', a') - Q_\theta(s, a)) (-\nabla_\theta Q_\theta(s, a))$$

$$\text{Stochastic gradient update rule: } \theta = \theta - \alpha \nabla_\theta E(s, a, s')$$