CZ4041 - Machine Learning

1 Bavesian Classifiers

Probabilities

Sum Rule
$$P(A) = \sum_{B} P(A, B)$$

$$P(A) = \sum_{B} \sum_{C} P(A, B, C)$$

Product Rule

$$P(A, B) = P(B|A) \times P(A) = P(A|B) \times P(B)$$

Bayes Theorem

$$P(A|B) = \frac{P(A,B)}{P(B)} = \frac{P(B|A)P(A))}{P(B)}$$

(Generalised case)

$$P(A_1...A_k|B_1...B_p) = \frac{P(B_1...B_p, A_1...A_k)}{P(B_1...B_p)}$$

Bavesian Classifiers

Bayesian classifiers aim to find the mapping $f: \mathbf{x} \Rightarrow y$ for supervised learning in the form of conditional probability $P(\bar{y}|\mathbf{X})$ Is equivalent to: via Bayes rule.

$$P(y|\mathbf{X}) = \frac{P(y,\mathbf{X})}{P(\mathbf{X})} = \frac{P(\mathbf{X}|y)P(y)}{P(\mathbf{X})}$$

For a classification with C classes, given a **3** Naïve Bayes Classifiers data instance x*:

$$y^* = c^* i f c^* = \underset{c}{\operatorname{argmax}} P(y = c | \mathbf{x}^*)$$

Applying Bayes rule,

$$P(y = c|\mathbf{x}^*) = \frac{P(\mathbf{x}^*|y = c)P(y = c)}{P(\mathbf{x}^*)}$$

Therefore,

$$y^* = \underset{c}{\operatorname{argmax}} \frac{P(\mathbf{x}^* | y = c)P(y = c)}{P(\mathbf{x}^*)}$$

$$= \operatorname*{argmax}_{c} P(\mathbf{x}^{*}|y=c)P(y=c)$$

2 Bayesian Decision Theory

Incorporating cost of misclassification on top of simple Bayesian Classifiers.

Loss/Cost

Actions: a_c , i.e., predict y = cDefine λ_{ij} as the cost of a_i when optimal action is a_i . E.g.:

$$\lambda_{00} = 0$$
 (predict correctly)
 $\lambda_{11} = 0$ (predict correctly)
 $\lambda_{01} = 10$ misclassify 1 as 0
 $\lambda_{00} = 1$ misclassify 0 as 1

Expected Risk

Expected risk for taking action a_i :

$$R(a_i|\mathbf{x}) = \sum_{c=0}^{C-1} \lambda_{ic} P(y=c|\mathbf{x})$$

To classify, for all actions, calculate expected risk, then choose the action with the minimum risk.

Special Case: 0/1 loss

$$\lambda_{ij} = \begin{cases} 0 \text{ if } i = j \\ 1 \text{ if } i \neq j \end{cases}$$

$$\therefore R(a_i|\mathbf{x}) = 1 - P(y = i|\mathbf{x})$$

In this case,

Choose
$$a_i$$
 if $R(a_i|\mathbf{x}) = \min_{a} R(a_c|\mathbf{x})$

Predict
$$y = c^*$$

if $P(y = c^*|\mathbf{x}) = \max P(y = c|\mathbf{x})$

Independence

A is **independent** of B, if:

$$P(A,B) = P(A|B) \times P(B) = P(A) \times P(B)$$

$$P(A,B) = P(B|A) \times P(A) = P(A) \times P(B)$$

Or

$$P(A|B) = P(A)$$

 $P(B|A) = P(B)$

given C if: P(A|B,C) = P(A|C)

Naïve Bayes Classifier

1. Assumption: conditional independence of features given label

$$p(\mathbf{x}|y=c) = P(x_1,...,x_d|y=c)$$

$$= P(x_1|y=c)P(x_2|y=c)...P(x_d|y=c)$$

$$= \prod_{i=1}^{d} P(x_i|y=c)$$

To classify a test record x^* , compute the 4 Bayesian Belief Networks posteriors for each class:

$$p(y=c|\mathbf{x}^*) = \frac{(\prod_{i=1}^d P(x_i^*|y=c))P(y=c)}{P(\mathbf{x}^*)}$$

Since $P(\mathbf{x}^*)$ is constant for each class c, it 1. A directed acyclic graph (DAG) enco-Binary split: Divides possible values is sufficient to choose the class that maxi mises the numerator term.

$$y^* = \underset{c}{\operatorname{argmax}} (\prod_{i=1}^{d} P(x_i^* | y = c)) P(y = c)$$

Estimating Cond Prob (Discrete)

$$P(x_i = k | y = c) = \frac{|(x_i - k) \land (y = c)|}{|y = c|}$$

Estimating Cond Prob (Continuous)

$$P(x_i|y=c) = \frac{1}{\sqrt{2\pi\sigma_{ic}^2}} e^{-\frac{(x_i - \mu_{ic})^2}{2\sigma_{ic}^2}}$$

Sample mean:

$$\mu_{ic} = \frac{1}{N_c} \sum_{j=1}^{N_c} x_{ij}$$

Sample variance:

$$\sigma_{ic}^2 = \frac{1}{N_c - 1} \sum_{j=1}^{N_c} (x_{ij} - \mu_{ic})^2$$

Laplace Estimate

Alternative prob estimation for discrete

$$P(x_i = k | y = c) = \frac{|(x_i - k) \land (y = c)| + 1}{|y = c| + n_i}$$

$$P(x_i = k | y = c) = \frac{1}{n_i}$$

M-estimate

A more general estimation:

$$P(x_i = k | y = c) = \frac{|(x_i - k) \land (y = c)| + m\tilde{P}}{|y = c| + m}$$
1) How to split the records?
2. Specifying feature test conduction.

Where *m* is a hyperparameter and \tilde{P} is pri- 2) When to stop splitting? or information of $P(x_i = k | v = c)$. (e.g., domain knowledge) Extreme case with no training data: $P(x_i = Splitting based on binary features$ $k|y=c\rangle = \tilde{P}(x_i=k|y=c)$

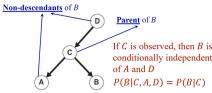
Suppose all features are **discrete** (if there - Multi-way split: Use as many partitions are continuous and discrete, estimation is as distinct values much more difficult)

Two key elements:

- ding dependence relationships between a as 2 subsets, need to find optimal partitioset of variables
- 2. A probability table associating each [Married] node to immediate parent nodes

DAG: Conditional Independence

A node in a Bayesian network is con-Binary split: $(x_j < v)$ or $(x_j \ge v)$ ditionally independent of its nondescendants, if its parents are known. Multi-way split (Discretization)



IMPORTANT! If A and B are condi-Supposing there are N_c instances in class tionally independent given C, we have:

1.
$$P(A|B,C) = P(A|C)$$

2.
$$P(A, B|C) = P(A|C)P(B|C)$$

Important! Using BBN for Inference

Given a BBN, and an inference(prediction) Information Gain

- 1. Translate problem into probabilisite $\Delta_{info} = E(parent) E(children)$ language
- 2. If the probabilities to be estimated of all children nodes and normalize by # cannot be obtained from the probability of training examples in each child node tables of the BBN, then
- A. Identify a subgraph which captures partitions (children), the dependence between input variables $\Delta_{info} = E(t) - \sum_{j=1}^{p} \frac{n_j}{n} E(j)$
- product rule, sum rule and the properties that result in large number of partitions. of conditional independence and indepen-Penalizing large number of partitions dence to induce equivalent forms of the (Gain Ratio) A is **conditionally independent** of B, extreme cases with no training data, found from the probabilities until all probabilities can be given C if:

5 Decision Trees

- Greedy strategy, split records based on feature test that optimises certain criterion Kev issues:

- Specifying feature test condition - Determining best split

Determining Test Conditions

2 Possible outcomes (e.g. Yes/No)

Splitting based on discrete features

e.g.: Marital Status \Rightarrow [Single], [Divorced] [Married]

e.g.: Marital Status ⇒ [Single, Divorced],

Splitting based on continuous features

Consider all possible splits and find the best cut Can be very computationally intensive

Determining Best Split

Using measure of node impurity – favour split with low degree of impurity

Measure of Impurity: Entropy

Entropy at a given node t:

$$E(t) = -\sum_{c} P(y=c;t) log_2 P(y=c;t)$$

To get entropy for children, get entropy Suppose a parent node t is split into P

$$\Delta_{info} = E(t) - \sum_{j=1}^{P} \frac{n_j}{n} E(j)$$

B. Based on the network topology, apply Disadvantage: Tends to prefer splits

$$\Delta_{InfoR} = \frac{\Delta_{info}}{SplitINFO}$$

SplitINFO =
$$-\sum_{i=1}^{p} \frac{n_i}{n} log_2(\frac{n_i}{n})$$

Stopping Criterias

- 1. All data belong to same class
- 2. Stop expanding when all data have similar feature vals
- 3. Early termination (avoid overfitting)

CZ4041 - Machine Learning

6 Generalisation

Overfitting: Test error rate increase when training error decrease Underfitting: Model too simple, both Euclidean distance: training and test error large

Training errors: error on training set,

Generalisation errors: error on previously unseen testing set, e'(T)

Estimating Generalisation Errors Optimistic Estimate

of overall data e'(T) = e(T)Decision tree induction algo select model with lowest training error rate.

Occam's Razor

Include information of model complexity when evaluating a model.

$$e'(T) = e(T) + N \times k$$

where N is the number of leaf nodes and k Scaling issues - features may need to be is a hyperparameter k > 0

Using Validation Set

Divide training data to 2 subsets, 1 for trai-ferent scales. ning and 1 for estimating generalisation Normalisation

Addressing overfitting

- user-specified threshold
- Stop if expanding current node does not Standardisation (z-score normalisation) The decision boundary of a SVM can be improve generalisation errrors

Post-Pruning

- Grow tree to its entirety
- Trim nodes in bottom-úp fashion
- If generalisation error improves after trimming, replace sub-tree by new leaf

- - "Training" is very efficient
- relatively expensive
- Requires training data to be stored in Activation functions memory

Classification steps:

- 1. Compute distance to other training
- 2. Identify K nearest neighbors
- 3. Use class labels of neighbors to determi-Since function is not differentiable, when ne class of instance

Choosing K

- K too small, sensitive to noise
- K too large, neighborhood may include points from other classes

Distance Metric

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{k=1}^{d} (x_{ik} - x_{jk})^2}$$

Voting Schemes

- Majority voting (sensitive to choice of k)
- Distance-weight voting (weight the influence of neighbor x_i according to distance Assume training set is good representation to test data)

$$w_i = \frac{1}{d(x^*, x_i)^2}$$

$$y^* = \underset{c}{\operatorname{argmax}} \sum_{(x_i, y_i) \in \mathcal{N}_{x^*}} w_i \times I(c = y_i)$$

Other issues with KNN

Solution: Normalisation on features of dif-

- Min-max normalisation

Pre-Pruning
- Stop if number of instances is less than
$$v_{\text{new}} = \frac{v_{\text{old}} - \min_{\text{old}}}{\max_{\text{old}} - \min_{\text{old}}} (\max_{\text{new}} - \min_{\text{new}})$$
Support Vector Machines user-specified threshold

Decision Boundary

(μ : mean, σ :standard deviation)

$$v_{\text{new}} = \frac{v_{\text{old}} - \mu_{\text{old}}}{\sigma_{\text{old}}}$$

8 ANN

$$y = \operatorname{sign}(\mathbf{w} \cdot \mathbf{x})$$

 $w_0 = -\theta, X_0 =$

7 KNN Classifiers $w_0 = -\theta$, $X_0 = 1$ learned. - Instance based, lazy learner - no model Where θ is the threshold term, w is the built weights vector and x is the input vector. For any test example \mathbf{x}^* An additional dimension is added to both - - Classifying unknown test instances are vectors such that the sum of products would minus the threshold term, θ .

Sign Activation function

$$sign(z) = \begin{cases} 1, & z \ge 0 \\ -1, otherwise \end{cases}$$

finding derivative of the activation function, we set y = z, and the derivative of y with respect to z would be = 1

Sigmoid Activation function

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

When $\lambda = 1$, it's called the sigmoid functi-

Derivative of sigmoid:

$$\frac{\partial \hat{y}(z)}{\partial z} = y(z) \cdot (1 - y(z))$$

Error/Loss

$$E = \frac{1}{2}(y_i - \hat{y}_i)^2$$

Updating Weights

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}}$$

Applying chain rule:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda \frac{\partial E(\hat{y})}{\partial \hat{y}} \frac{\partial \hat{y}(z)}{\partial z} \frac{\partial z(\mathbf{w})}{\partial \mathbf{w}}$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda(-(y_i - \hat{y_i}))(1)(\mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \lambda(-(y_i - \hat{y}_i))(1)(\mathbf{x}_i)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \lambda (y_i - \hat{y}_i) \mathbf{x}_i$$

Decision Boundary

defined as:

$$w_1x_1 + w_2x_w + b = 0$$

General form:

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

Making predictions

During training, the values of w and b is

$$\begin{cases} f(\mathbf{x}^*) = +1, \text{if } \mathbf{w} \cdot \mathbf{x}^* + b \ge 0 \\ f(\mathbf{x}^*) = -1, \text{if } \mathbf{w} \cdot \mathbf{x}^* + b < 0 \end{cases}$$

Other notes (Linear Algebra): **Inner Product**

$$\mathbf{u} \cdot \mathbf{v} = \sum_{i=1}^{d} (u_i \times v_i)$$

$$\mathbf{u} \cdot \mathbf{v} = \|\mathbf{u}\|_2 \times \|\mathbf{v}\|_2 \times cos(\theta)$$

L2 Norm (Length of vector)

$$\|\mathbf{u}\|_2 = \sqrt{\mathbf{u} \cdot \mathbf{u}} = \sqrt{\sum_{i=1}^d (u_i \times u_i)}$$

Induction

- Direction of **w** is orthogonal (perpendi-minimising the error cular) to the decision boundary.

Parallel hyperplanes:

$$\mathbf{w} \cdot \mathbf{x} + b = k$$

 $\mathbf{w} \cdot \mathbf{x} + b = -k$

(After rescaling $\mathbf{w} = \mathbf{w}/k$, b = b/k)

$$\mathbf{w} \cdot \mathbf{x} + b = 1$$

 $\mathbf{w} \cdot \mathbf{x} + b = -1$

 $\|\mathbf{w}\|_2 \times d = 2$

$$d = \frac{2}{\|\mathbf{w}\|_2}$$

10 Linear Regression

Error for 1-D Linear Regression Model Sum-of-squares (SSE) error:

$$E(w) = \frac{1}{2} \sum_{i=1}^{N} (w \times x_i - y_i)^2$$

Learn linear model in terms of w by

$$w^* = \operatorname*{argmin}_{w} E(w)$$

To solve the unconstrained minimisation problem, set derivative of E(w) w.r.t w to

$$\frac{\partial E(w)}{\partial w} = \frac{\partial (\frac{1}{2} \sum_{i=1}^{N} (w \times x_i - y_i)^2)}{\partial w} = 0$$

Closed form solution:

$$w = \frac{\sum_{i=1}^{N} y_i \times x_i}{\sum_{i=1}^{N} x_i^2}$$

More general case (multi-dimension)

Error for Linear Regression Model

 $E(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} (\mathbf{w} \cdot \mathbf{x}_i - y_i)^2$

Learn linear model in terms of w by mini-

mizing the error (with regularisation term)

 $\mathbf{w}^* = \operatorname{argmin} E(\mathbf{w}) + \frac{\lambda}{2} ||\mathbf{w}||_2^2$

Margin Maximisation

Therefore, decision boundary can be learnt $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$ by maximising the margin, $d = \frac{2}{\|\mathbf{w}\|_2}$. Ho-- By defining $w_0 = b$ and $X_0 = 1$, w wever, this is not easy. Change this into a and x are of d+1 dimensions minimisation problem. $f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x}$

Minimise: $\frac{\|\mathbf{w}\|_2^2}{2}$ Constraints: $\mathbf{w} \cdot \mathbf{x}_i + b \ge 1$, if $v_i = 1$ $\mathbf{w} \cdot \mathbf{x}_i + b \le -1$, if $y_i = -1$ OR, $v_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$

Multi-Class SVM

 C_1 , Negative C_2 & C_3

class

mine

 $f_1(\mathbf{x}^*) = -1$

 $f_2(\boldsymbol{x}^*) = 1$

 $f_2(x^*) = -1$

Total Votes:

Optimisation Problem for SVM

$$min_{w,b} \frac{\|\mathbf{w}\|_2^2}{2}$$

Given 3-class problem C_1 , C_2 and C_3

2. Positive C_2 , Negative $C_1 \& C_3$

3. Positive C_3 , Negative C_1 & C_2

$$s.t.y_i \times (\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1$$

Create 3 SVM binary classifiers: 1. Positive

majority voting to

for

1

test

 $\boldsymbol{c_2}$

1

1

1

 $\boldsymbol{c_3}$

1

Closed-Form Solution

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T y$$

11 Ensemble Learning Necessary Conditions

- 1) Base classifiers are independent of each other
- 2) Base classifiers should do better than deter- classifier that performs random guessing example. (i.e., acc > 0.5)

Error rate of ensemble

Supposing N independent base classifiers with error ϵ :

$$P = \sum_{i=(N//2)+1}^{N} \epsilon^{i} (1 - \epsilon)^{N-i}$$

Ensemble Methods Bagging

- Sample examples with replacement and build model on each bootstrap sample.
- Use majority voting to determine class label of ensemble classifier
- A bootstrap sample contains approximately 63.2% of original training data

- equal weights
- is trained from sample
- from training set
- 4) Update weights of examples after the end of boosting round
 - Wrongly classified increase
 - Correctly classified decrease
 - Examples not drawn unchanged
- 5) Use weighted voting, each classifier 1) Select k data as initial centroids would have different weights

Random Forests

- Specifically designed for decision tree instances to closest centroid to grow
- 2) Choose m', number of features used to calculate best split (Typically 20%)
- 3) For each tree
- Choose training set via bootstrapping
- For each node, randomly choose m' featu-Total Sum of squared error (SSE) res and calculate best split
- Trees fully grown and not pruned
- 4) Use majority vote among all trees

Combination Methods

- Majority voting
- Weighted voting
- Simple average:

$$f_M(\mathbf{x}) = \frac{1}{T} \sum_{i=1}^{T} f_i(\mathbf{X})$$
 - Weighted average:

$$f_M(\mathbf{x}) = \frac{1}{T} \sum_{i=1}^T w_i f_i(\mathbf{X})$$

Where $w_i \ge 0$, and $\sum_{i=1}^T w_i = 1$

Combining by Learning

Combiner: second-level learner, or meta-

Combiner takes output of base classifiers SSE as features, and learn to classify based on output label.

12 Clustering

Cluster Analysis

- Finding groups of data instances such Bisecting K-means that data instances in a group are:
 - · Similar to one another

• Different from data in other groups K:

It is NOT:

- available)
- Simple segmentation (e.g. divide ex-2a) Choose the split with the lowest SSE 1) Start with tree with a single point amples into groups by properties/feature from the T number of results

- Partitional (non-overlapping)
- Hierarchical (nested)
- 1) Initially, all examples are assigned Exclusive/Non-exclusive (instance belon--Sizes ging to multiple cluster)
- 2) Bootstrap sample is drawn and a model Fuzzy/Non-fuzzy (point belongs to every Or when clusters have non-globular shacluster with weight (sums to 1))
 - clustered for partial)

K-means

- Partitional clustering Preprocessing: - Normalize data
- Eliminate outliers

Algorithm:

- Loop while centoids' values are updated,
- 2) Form k clusters by assigning data Types:
- (mean of data instances in cluster)

choie of initial centroids.

Evaluation of K-means clusters

$$SSE = \sum_{i=1}^{K} \sum_{\mathbf{x} \in C_i} \operatorname{dist}(\mathbf{c}_i, \mathbf{x})^2$$

Solving Initial Centoids Issue

- 1) Multiple runs, choose run with lowest pairwise distance for all points in clusters
- 2) Post-processing:
- Merge clusters with low cluster SSE, undone (errors propagated) which are close to each other
- 3) Bisecting K-means

Empty Clusters Issue

Choose replacement centroid:

- Choose point that contributes most to 2) Difficulty handling different sized
- Choose a point from cluster with highest cluster SSE
- Repeat several times if there are multiple **Divisive Hierarchical Clustering** empty clusters

Start with 1 cluster containing all points 2) While all clusters are not singleton clus-While the number of clusters is less than ters (Or while number of clusters < K):

- Supervised classification (class label info 2) Bisect the cluster with simple K-means Building MST (2 cluster) T times.
 - 2b) Add the split clusters into the list of (random) clusters

K-means limitations

K-means has problems when clusters have other is not. differing:

- Densities
- 3) Model is then used to classify examples Partial/Complete (Only some instances K-means also sensitive when data has out-Goal: Estimate unobservable underlying

Hierarchical Clustering

(Visualised as dendrogram)

Strengths:

- Do not have to assume number of clus-- Parametric density estimation ters, can obtain desired number by cutting - Nonparametric density estimation the dendrogram at the proper level.

- classifiers 1) Choose T, number of trees 3) Recompute the centroid of each cluster instance, then merge until 1 big cluster)
 - points, split at each step until 1 clus-Note: performance affected greatly by ter/instance or when there are K clusters) is sampled from sampled from sampled from sampled from sampled from the sample

Agglomerative

- Merge most similar clusters at each step dependent events. - Update cluster proximity based on type
- of Inter-Cluster Similarity: different clusters
- MAX: Proximity = furthest points in up to parameters, θ .
- different clusters - Group Average: Proximity = average of ling x_i from $p(\mathbf{x}; \theta)$ as likely as possible.

Limitations

- Decompose cluster with high Cluster SSE Once clusters are combined, can't be Maximum Likelihood Estimation
 - No objective function is minimized
 - following:
 - 1) Sensitive to noise/outliers

- Algorithm:
- lect all instances as single cluster

1) Select the cluster with the largest SSE link corresponding to largest distance. product can be converted into sum

Algorithm:

While there are points not in the tree: 2) Look for the closest pair of points such

- that 1 point is in current tree and the
- 3) Add the point into the tree and an edge (with a value of the distance) between the two points

Density Estimation

probability density from observed data.

- Nested clusters organised as hierarchical Note, probability = area under density fn curve. Integrating entire curve gives

Approaches for Density Estimation

General Principle

Observed data points assumed to be sam-- Agglomerative (Start with 1 cluster per ple of N random variables independent and identically distributed (i.i.d.) - Divisive (Start with 1 cluster with all Identically Distributed: For any $x_i \in D$, it

is sampled from sampled from same pro-

Parametric Density Estimation

- MIN: Proximity = closest points in Assume that data are drawn from known probability density family, $p(\mathbf{x}; \theta)$, defined First bin start from 0, therefore determine

Task: Find parameter θ that makes samp-Naïve Estimator

Approach: Maximum Likelihood Estima-

Likelihood of parameter θ given sample Windowing function \mathcal{D} : $l(\theta|\mathcal{D}) \triangleq p(\tilde{\mathcal{D}};\theta)$

- Different schemes have issues with the Since data are i.i.d, the likelihood is product of likelihoods of individual data points:

$$l(\theta|\mathcal{D}) \triangleq p(\mathcal{D};\theta) = \prod_{i=1}^{N} p(x_i;\theta)$$

1) Generate minimum spanning tree to col-Therefore, find θ that makes \mathcal{D} most likely to be drawn:

$$\hat{\theta} = \operatorname*{argmax}_{\theta} l(\theta|\mathcal{D})$$

3) Create a new cluster by breaking the Typically, log-likelihood is used such that

$$ln \ l(\theta|\mathcal{D}) = \sum_{i=1}^{N} ln \ p(x_i;\theta)$$

Univariate Gaussian

$$\mathbf{x}_i \sim p(\mathbf{x}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)}$$

Solutions to parameters μ and σ^2 (Unbiased estimation):

$$\begin{cases} \hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_i \\ \hat{\sigma^2} = \frac{1}{N-1} \sum_{i=1}^{N} (\mathbf{x}_i - \hat{\mu})^2 \end{cases}$$

Multivariate Gaussian

Solutions to parameters μ (d-dimensional mean vector)

and Σ (d x d covariance matrix) (Unbiased estimation):

$$\begin{cases} \hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}_i \\ \hat{\Sigma} = \frac{1}{N-1} \sum_{i=1}^{N} (\mathbf{x}_i - \hat{\mu}) (\mathbf{x}_i - \hat{\mu})^T \end{cases}$$

Independent: all data points $x_i \in D$ are in-Nonparametric density estimation Histogram Estimator

$$\hat{p}(\mathbf{x}) = \frac{\#x_i \text{ in same bin as } x}{N\Lambda}$$

which bin to use by using x

$$\hat{p}(\mathbf{x}) = \frac{1}{N\Delta} \sum_{i=1}^{N} w(\frac{\mathbf{x} - \mathbf{x}_i}{\Delta})$$

$$w(u) = \begin{cases} 1 \text{ if } -\frac{1}{2} \le u < \frac{1}{2} \\ 0 \text{ otherwise} \end{cases}$$

Generalisation to Multivariate

If the data is d-dimensional, the bins become a d-dimensional hypercube with volume: $V = h^d$, where h is the length of each edge and d is the number of dimensions. Windowing function:

$$w(u) = \begin{cases} 1 \text{ if } -\frac{1}{2} \le u_j < \frac{1}{2} \text{ for all } j \in \{1,.,d\} \\ 0 \text{ otherwise} \end{cases}$$

C74041 -	Machine	Lagrning
CZ-10-11 -	Maciline	Learning

$$\hat{p}(\mathbf{x}) = \frac{1}{NV} \sum_{i=1}^{N} w(\frac{\mathbf{x} - \mathbf{x}_i}{h})$$

KNN Estimator

Set the numerator as a constant (K), adapt smoothing to local density of data. Degree of smoothing controlled by K (number of neighbours).

$$p(\mathbf{x}) = \frac{K}{NV_x}$$

 V_x is the volume of the space centered at ${\bf x}$ that exactly contains K nearest neighbors of ${\bf x}$.

With a predefined K, for a data point \mathbf{x} ,
1) Compute distance of all points to x (e.g.
euclidean dist)

²⁾ Sort observed data points based on distances in asc order

distances in asc order
3) Obtain the k-th distance to compute $V_{d_k(\mathbf{x})}$, which is the volume of the d-ball

Dimension	Volume of a ball of radius ${\it R}$
0	1
1	2R
2	$\pi R^2 \approx 3.142 \times R^2$
3	$rac{4\pi}{3}R^3pprox 4.189 imes R^3$
4	$\frac{\pi^2}{2}R^4\approx 4.935\times R^4$
5	$\frac{8\pi^2}{15}R^5\approx 5.264\times R^5$
6	$\frac{\pi^3}{6}R^6\approx 5.168\times R^6$
7	$\frac{16\pi^3}{105}R^7\approx 4.725\times R^7$
8	$\frac{\pi^4}{24}R^8\approx 4.059\times R^8$
9	$\frac{32\pi^4}{945}R^9\approx 3.299\times R^9$
10	$\frac{\pi^5}{120} R^{10} \approx 2.550 \times R^{10}$
11	$\frac{64\pi^5}{10395}R^{11}\approx 1.884\times R^{11}$
12	$\frac{\pi^6}{720} R^{12} \approx 1.335 \times R^{12}$
13	$\frac{128\pi^6}{135135}R^{13}\approx 0.911\times R^{13}$
14	$\frac{\pi^7}{5040}R^{14}\approx 0.599\times R^{14}$
15	$rac{256\pi^7}{2027025}R^{15}pprox 0.381 imes R^{15}$