

Machine Learning

< Block I: Regression >

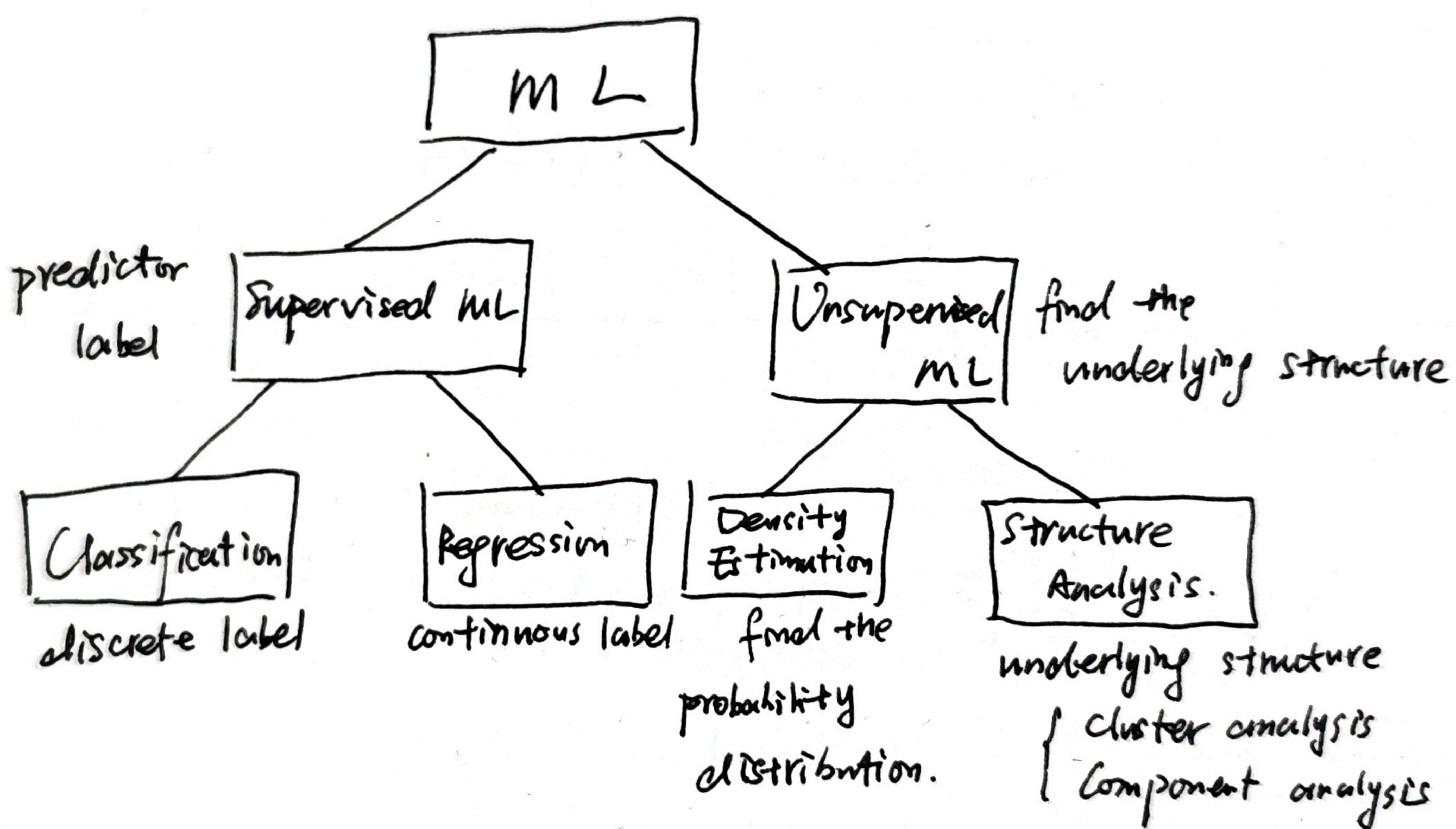
0. Introduction.

ML: The ability to acquire knowledge by extracting patterns from data.

Data: Instance of an observation or a measurement.

Datasets: Data formatted as collections of items described by a set of pre-defined attributes.

ML: A set of tools together with a methodological methodology for solving problems using data.



1. Regression.

Define metric. Let $y = \hat{y} + e$, Sum of squared errors: $SSE = \sum e_i^2$

mean squared error: $MSE = \frac{1}{N} \sum e_i^2$

root mean squared error: $RMSE = \sqrt{\frac{1}{N} \sum e_i^2}$

mean absolute error: $MAE = \frac{1}{N} \sum |e_i|$

R-squared: $R^2 = 1 - \frac{\sum e_i^2}{\sum (y_i - \bar{y})^2}$

Model
Multiple linear regression

$$Xw = Y$$

$$\begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N,1} & \dots & x_{N,k} \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

; 每行代表一个样本，使用最小二乘
 $\min \frac{1}{N} \| Xw - Y \|^2$ 为目标函数
 $w = (X^T X)^{-1} X^T Y$.

Remark 多回旧的高以形为基线而得也，拿一个 predictor 算法。 $Xw = Y$.

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^k \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 1 & x_N & x_N^2 & \dots & x_N^k \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_k \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}, \text{求解方法依然 } w = (X^T X)^{-1} X^T Y.$$

Define Generalisation is the ability of our model to successfully translate what we was learnt during learning stage to deployment. (translate knowledge)

Define 1) Underfitting = unable to describe the underlying pattern.

Large train and deployment errors are produced.

2) Overfitting = memorisation of irrelevant details.

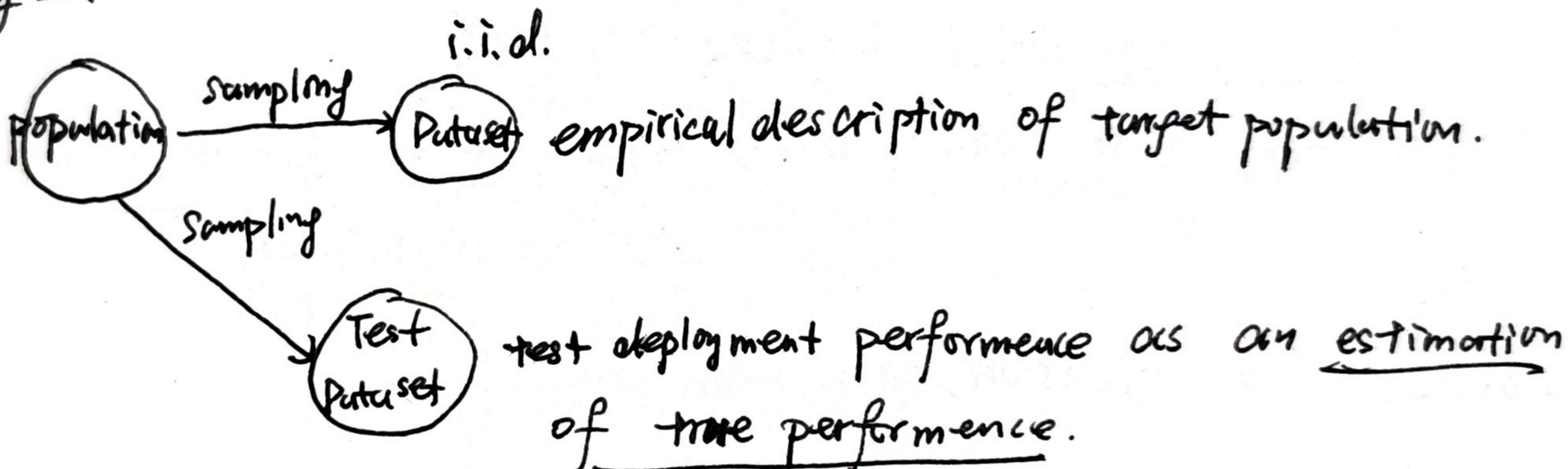
数据少，模型复杂，无法泛化。

Small errors in train, large errors in deployment.

3) Just right: reflects underlying pattern and ignores irrelevant details.

2. Methodology I.

model
target population.



Algorithm
GD

Gradient descent: $w_{t+1} = w_t - \gamma \nabla L(w)$, update iteratively

where γ is called learning rate or step size.

Remark 1) 并不是在 error surface 上找梯度，而是找 Train dataset 描述的 empirical error surface.

2) 实际使用时，将 Training data 分为 Batch, 每次对 N 个样本 (Batch) 进行梯度的计算 \rightarrow Stochastic gradient descent.

原因：small batches produce noisy of gradient of empirical error surface, which can help to escape local minima.

3) other 改进：momentum; RMSProp \mapsto Adam.

Regularisation: modifies the empirical error surface by adding a term that constrains the values that the model parameters can take on. (ex. $\frac{1}{N} \|w\|_2^2 + \lambda \|w\|_1^2$)

Remark 使用 Regularisation 来降低 model's complexity，在训练阶段时候用的是 loss function，但是 deployment 用的是指打，是 target quality metric，不带子啊梯度。Regularisation 只是为了在训练时使用。

Validation Validation allow us to use data for assessing and selecting different families of models, then the one will be trained.

1) Validation set approach Randomly splits the available dataset into training and a validation (holdout) set.

2) Leave-one-out cross-validation (LOOCV), The validation set contains only one sample. Multiple split, final performance is calculated as the average of individual performances.

3) K-fold cross-validation 先对 k 个子集取平均

Remark 这些所有的 cross-validation 都是为了得到 validation error, 可以直接使用, 也可以平均数使用, 但是都是仅得到了一个误差, 我们需要对每一个 model, 或者每部运行一遍 cross-validation, 得到每个 model 对应的 validation error, 再根据进行 "model selection", 把所有模型在全部训练数据上再训练一遍 model 作为最终 output model.

总结: 最最终的 model quality 由 type of model, optimisation strategy and the representativeness of the training data.

<Block 2: Classification>

1. Classification I: 从属类别.

problem classification.
1) no attribute space 在 dataset, label $\in \{0, 1\}$ axis
无特征空间, 以 predictor space 在 dataset, using different symbols for each label in the predictor space.

2) A solution model in classification is a partition of the predictor space into decision regions \Rightarrow separated by decision boundaries.

define
特征向量. In linearly separable datasets, we can find a linear classifier that achieves the maximum ~~accuracy~~ accuracy ($A=1, \epsilon = 1-A=0$). In non linearly-separable datasets, the best accuracy will be $A < 1$, error rate $E > 0$.

Algorithm
logistic regression. Let $P(\text{class } 1) = p(w^T x_i) = \frac{e^{w^T x_i}}{1+e^{w^T x_i}} = \frac{1}{1+e^{-w^T x_i}}$, then use MLE. That.

$$\max_w L = \max_w \prod_{j=1}^n (1-p(w^T x_j)) \cdot \prod_{i=1}^m p(w^T x_i) \quad \text{or.}$$

$$\max_w \log L = \max_w \sum_{i=1}^m \log [1-p(w^T x_i)] + \sum_{i=1}^n \log p(w^T x_i).$$

-introduction
分类策略

Remark

- 1) logistic regression 在做着一种进行 Linear classifier 的唯一法：keep that boundary away, 即最大似然，只有达到离散边界，才可最大化似然。
- 2) $\text{max}_i \phi_i = \mathbf{w}^T \mathbf{x}_i$, ϕ_i 算 point i 对于 Boundary 的欧式距离, 用 $\phi(\mathbf{x})$ 计算整体的 likelihood.

Algorithm

Nearest neighbors

Non-parametric approaches = K-Nearest Neighbors (KNN)

Instance-based method: The whole training dataset need to be memorised.

该点 \mathbf{x} 找最近的 K 点, 遍历寻找 K 点中多数类。

K 增大, boundary become less complex, overfitting $\xleftarrow[\text{large } K]{\text{less}} \xrightarrow{\text{more}} \text{underfitting}$.

总结：1) classifiers are partitions of the predictor space into decision regions separated by boundaries.

2. Classification II: 极端视角

model

Bayes classifier

$$P(\text{Class}_k | \mathbf{x}) = \frac{P(\text{Class}_k) P(\mathbf{x} | \text{Class}_k)}{P(\mathbf{x})} \text{ of } P(\mathbf{x} | \text{Class}_k) P(\text{Class}_k).$$

对 Bayes classifier
 $P(\text{Class}_1 | \mathbf{x}) / P(\text{Class}_2 | \mathbf{x}) \geq \begin{cases} \text{Class 1} & \text{if } P(\text{Class}_1 | \mathbf{x}) > P(\text{Class}_2 | \mathbf{x}) \\ \text{Class 2} & \text{if } P(\text{Class}_2 | \mathbf{x}) > P(\text{Class}_1 | \mathbf{x}) \end{cases}$
 增加类别的
 likelihood class densities Prior.

Use $\frac{P(\text{Class}_1 | \mathbf{x})}{P(\text{Class}_2 | \mathbf{x})} \geq 1$ 为后验概率, 使用 true posterior probabilities
 为 Bayesian Risk (min the expected cost).

假若类别的 Bayes classifier: Cost $\rightarrow C_1 \times P(\text{Class}_1 | \mathbf{x}) \geq 1$.
 $C_2 \times P(\text{Class}_2 | \mathbf{x})$

Discriminant analysis

algorithm

对 likelihood (class densities) 为假设是 Gaussian

$$P(\mathbf{x} | \text{Class}_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left\{-\frac{(\mathbf{x} - \mu_i)^2}{2\sigma_i^2}\right\}$$

对先验的估计是频率诱导的频率。

Note 1) 对于 Discriminant analysis for 2 类, $\Sigma_1 = \Sigma_2$, boundary 是线性的 \rightarrow Linear DA, LDA.

如果 $\Sigma_1 \neq \Sigma_2$, 也是二分类, Quadratic DA, QDA.

2) 使用 Discriminant analysis, logistic regression, KNN 都是对于 posteriors 的 fit, 在 Bayesian Risk 中, 通过调整 C_1, C_2 (Risk), 用 Risk 或者说是意义上指 Bayes Boundary 的选取。

Metrics

Confusion matrix

对 Class-sensitive problem 于 imbalanced dataset, 使用: $\text{Cost}_1 \times P(\text{Class}_1 | \mathbf{x}) = \text{Cost}_2 \times P(\text{Class}_2 | \mathbf{x})$

Accuracy 与 error rate 可能并不差很好的选择, 可以考虑使用 Confusion Matrix.

| | Actual Class | | 错误的总和 |
|-----------|--------------|--------|----------------|
| | predicted | Actual | |
| predicted | P | N | TPR + FNR = 1. |
| | NP | FN | |
| class | TP | FP | FPR + TNR = 1. |
| | FN | TN | |

| | P | N | |
|---|-----|-----|----------------|
| P | TPR | FPR | TPR + FNR = 1. |
| N | FNR | TNR | FPR + TNR = 1. |

Note) 通过 Confuse matrix 可以衍生出一些指标:

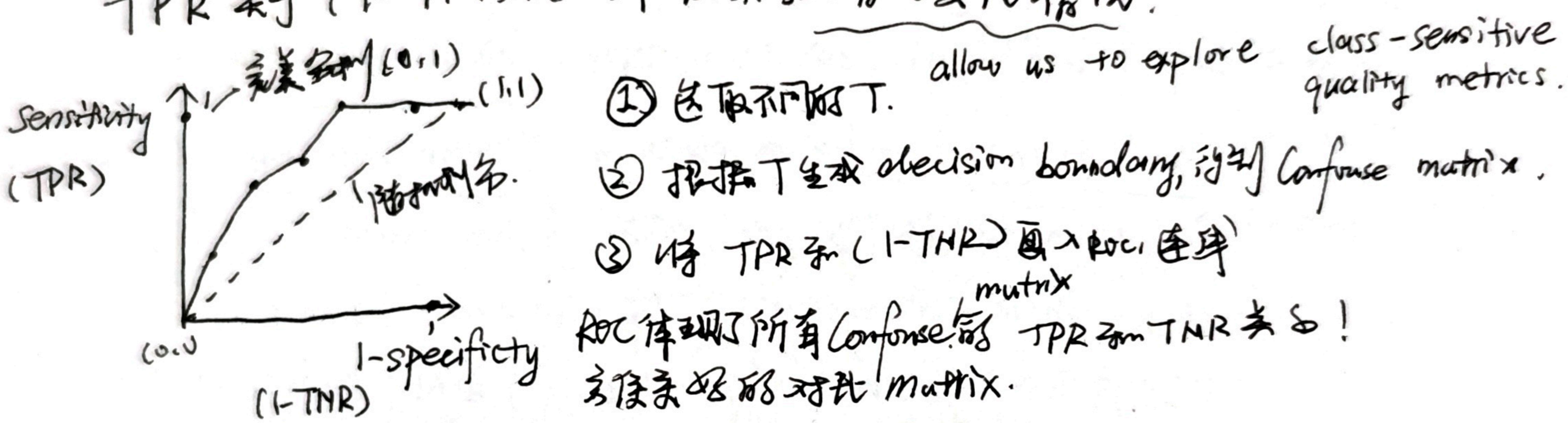
- Sensitivity (recall / true positive rate) = TPR. 灵敏度, 正向率.
- Specificity (true negative rate) = TNR 特异度
- Precision (positive predictive value) = $\frac{TP}{TP+FP}$ 精准, P 中正确的概率.

优化某个指标会带来问题, when improve one metric, others may decrease.
所以有时会考虑 pairs of quality metrics.

$$\text{ex. F1 score: } F1 = \frac{2}{\frac{\text{recall}^{-1} + \text{precision}^{-1}}{2}} = \frac{\text{Precision} \times \text{recall}}{\text{Precision} + \text{recall}}$$

3) 前文中提到 Risk Decision, $\frac{P(\text{class1}|\chi)}{P(\text{class2}|\chi)} \geq T$, 为什么进行了T的选择呢?

可以画出 ROC (Receiver operating characteristic) plane, ROC 是 TPR 关于 $(1-TNR)$ 的曲线, 反映了二者的变关系.



最后选择的 T 距离点 $(0,1)$ 最近那点, 对应的 T.

4) 现在我们如何对模型 model 了, 如何对不同 model 了? $\rightarrow \boxed{\text{AUC!}}$.

AUC - Area under the curve, 即 ROC 曲线下面积

注: AUC 从概率角度进行解释, 反映了 model 捕获到 map 的能力.

该 map 为 score of P > score of N 的映射, 该映射依赖于 latent variable.

好的 AUC $\rightarrow 1$, 差的 AUC $\rightarrow 0.5$, 且 0.5 表示而已, 根据 AUC 进行 model selection.

5) 对 AUC 的评价: AUC is a measure of goodness for a classifier that select T. \rightarrow can be calibrated 校准.
不是 T, 是 model's parameter

< Block 3: Methodology, DL, NNs >.

1. Methodology II.

a. pipeline The term pipeline is often used to describe workflows (sequence of operations).

ML solution are more than a model it include several stages form a processing pipeline: Transformation stages, Several ML models, Aggregation stages.

所有的 pipeline 都是 deployment.

b. Normalisation Data normalisation is a tunable transformation stage that allows us to scale attributes so that their value belongs to similar range.

1) Min-max normalisation.

$$z = \frac{x - \min(x)}{\max(x) - \min(x)} \in [0, 1], \text{ 例如将 } [0, 2] \text{ 转换为 } [0, 1].$$

2) Standardisation.

$$z = \frac{x - \mu}{\sigma}, \text{ if } z=0, \text{ if } z^2=1, \text{ ensures inputs are treated equally.}$$

Note μ, σ 都是 from dataset 中的参数, 在 test dataset 中, 但不包含出-of-range values, 以及 outliers

c. Transformations Transformations are data manipulations that change the way that we represent our samples \Rightarrow move the data to another space.

类型: 同 dim space 变换, 不同 dim space 变换 (dimensionality reduction)

PCA Principal components analysis.

① identifies the directions along which samples are aligned. These directions define a destination space with the same number of dimensions as the original space, i.e. build a linear transformation and additionally assigns eigenvalue to each component.

② We can use the score (eigenvalue) to rank the attributes, define the destination space and remove the least important attributes. That is, use PCA to reduce dimension.

Dimensionality Reduct

① Feature selection. Assume that only a subset of the original attributes are relevant.

② Filtering. Consider each attribute individually, don't consider their interactions.

③ Wrapping. Consider the interaction between features, evaluate them together.

注: 使用 Filtering 和 Wrapping 的方法: 使用不同 subset 进行 validation, validation validation, rank, 通过 validation performance 对比各个 subsets.

④ Feature extraction. New attributes are defined as operations on the original attributes.

Note Transformation not part learn! selecting the right transformation (via validation) - tuning the parameters of given transformation (via training).

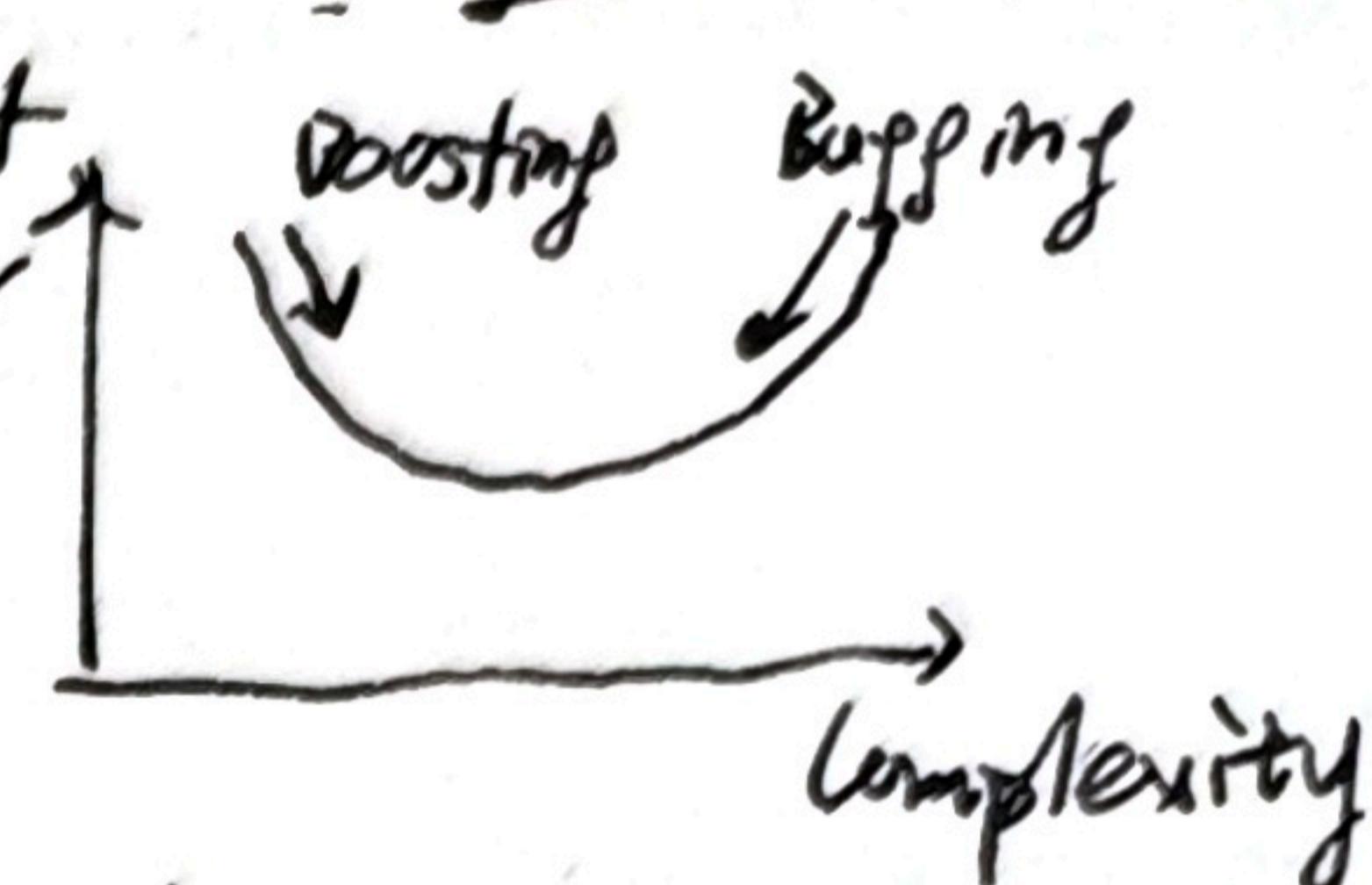
d. Ensembles. Combines the strengths of base models.

- 1) Bagging. Generates K sub-dataset by bootstrapping and trains K simple base models with each-dataset. $f_1(x) = \bar{z} \frac{\sum f_k(x)}{K}$ where Bootstrap is a statistical method that extracts random samples from a dataset. (Randomisation).

- Samples from a dataset.

2) Boosting. Generates a sequence of simple base models, where each successive model focus on the samples that the previous models could not handle properly. Test error \downarrow Boosting Supporting

Note Bagging 是并行算法, Boosting 依赖于基- t model, $t \leq t'$.



Decision trees. Partition the predictor space into multiple decision regions by implementing sequences of splitting rules using only one predictor only. (每次用一个predictor, 也就是Region boundary是axis-parallel的).

Random Forests. Train many individual trees by randomising the training samples and the predictors. Predictions are obtained by averaging the individual predictions. \mapsto Bagging.

2. Neural Networks and Deep Learning.

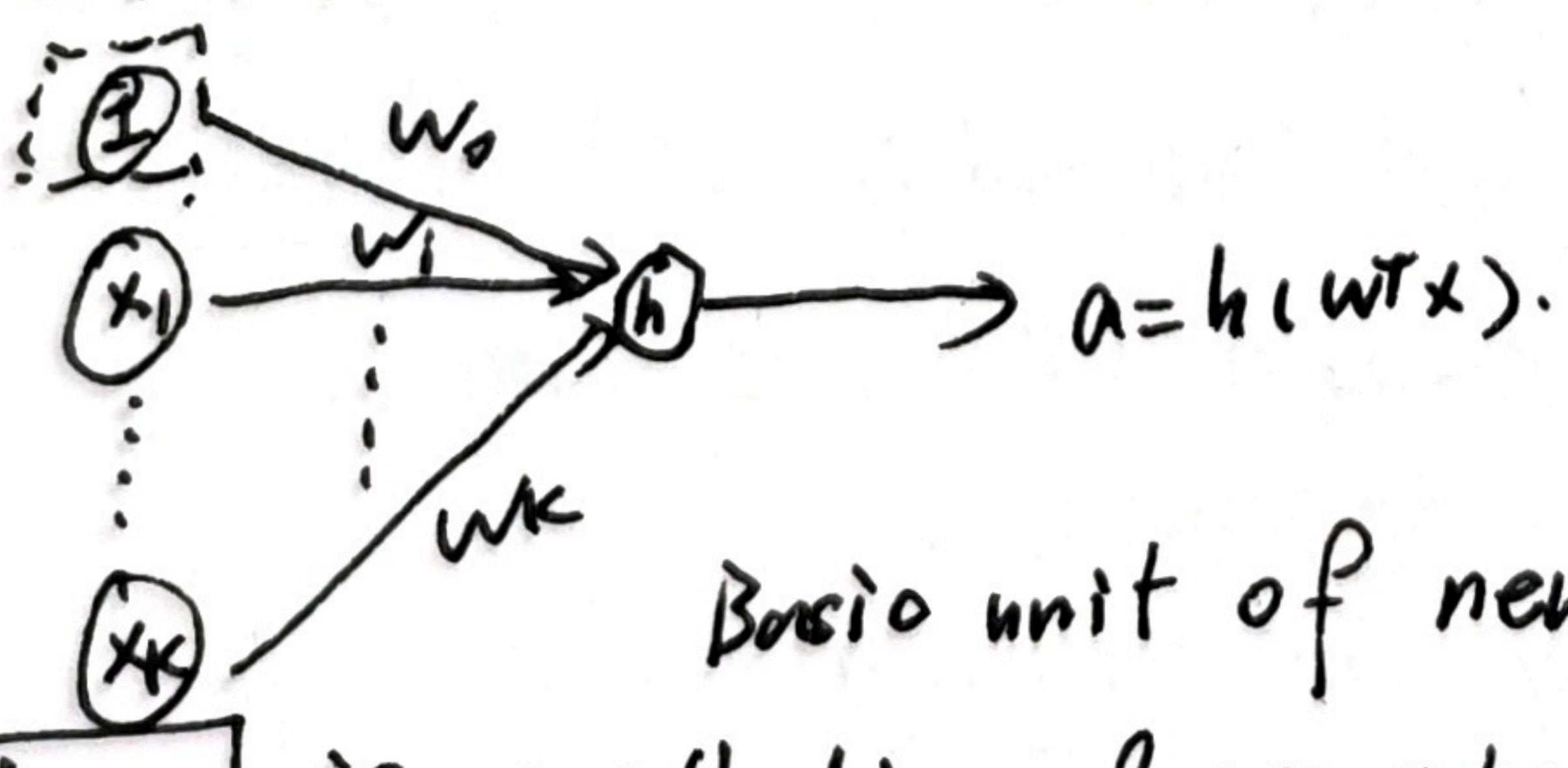
Patterns are laws in our data, structure is a law in our target population, we should discover the underlying structure by identifying the pattern.

perception defined by a weight vector w and an activation function $h(x)$, that maps

The computation on extended x to the output 2. w_0 is bias.

祝扇君，八月廿九

是 perceptions
的颜色！



- | $h \rightarrow$ step activation function \rightarrow linear classifier
- | $h \rightarrow$ logistic function \rightarrow logistic regression.
- | common \mapsto grid pattern detector.

Remark 1) A layer is a collection of perceptrons that use the same input.

ex. L 4 perceptrons, KF input, $(K+1) \times L$ 4 weight.

2) Architecture: describe how layers are connected.

hidden layers produce internal features.

3) A neural network consists of perceptrons arranged in layers that are connected out according to an architecture.

NNs as ML models - NNs can be seen as an entire trainable ML pipeline.

① Cognitive point. large NMs is flexibility to create new and
increasingly complex concepts that might be relevant to predict.

⑤ ~~minimum~~ ML algorithm 对有 model, cost function 与 optimisation method.

对于NNs来说，可以使用 Back-propagation 来有组织地计算 Gradient，始于 output.

③ Transfer learning. 训练好的模型不变, fixed transformation stage $T^{(x)}$

To return back to the stage using new data first ($Z = T_1 x$) by transferring an already learnt transformation.

Types of layers

1) Fully-Connected layer (FC).

2) Convolutional layer.

introduction = Images and time series with equivariance property ,

Same pattern can be expected in different locations of the grid.

Solution: Perceptrons 被排列成 grid, 有 "feature map."

share their parameter \rightarrow convolutional kernel!

卷积层由多个 feature maps 构成，每个由不同 kernel 处理，different
concept，最后 stalk 在一起形成 convolutional layer.

Note. often dimensions of kernel are $H \times W \times D$. 重量及, weight per kernel is $H \times W \times D + 1$ (bias), Train 要教元子对上样例的 kernel to $H \times W \times D + 1$ weight. \rightarrow kernel 中的权值  + bias.

3) Pooling layers.

Reduces the size of feature maps. 由特征图减小尺寸，
插入 Convolutional layers 之间。

① Max pooling ② Average pooling 不仅仅是卷积！

方法：
1) Fractional (map), Cognitive (create new concepts), Computational
2) NNs 构成 m model, ~~通过~~ firm 当作 生产方法, 即 NNs 和 一个
FB architecture, 不同的 architecture 是不同 model. (通过 validation 模型选择).

<Block 4: Unsupervised learning>.

1. Structure analysis.

Unsupervised learning: Understanding how samples are distributed in the attribute space.

Method 1: Structure analysis: identify regions within attribute space (cluster) or directions with high density (component).

Method 2: Density estimation: quantify the probability that find the sample within a region in attribute space.

应用: Discover structure, generate new knowledge, change the way we represent the data.

Clusters: Describe the structure of dataset as group.

Algorithm

k-means.

是基于 proximity-based quality metric:

$$\text{① intra-cluster sample scatter: } I(C_i) = \frac{1}{2} \sum_{\substack{x_i, x_j \\ \in C_i}} \|x_i - x_j\|_2^2$$

$$\text{② inter-cluster sample scatter: } OCC(C_0, C_1) = \sum_{\substack{x_i \in C_0 \\ x_j \in C_1}} \|x_i - x_j\|_2^2.$$

Best clustering = lowest intra-cluster and highest inter-cluster.

輔註

即上式即 intra distance + inter

k-means 以 center (mean) of cluster 为 prototypes: $\mu = \frac{1}{N} \sum x_i$.

上述 intra-cluster sample scatter 有另一种写法: $I = N_0 \sum \|x_i - \mu\|_2^2$.

So that good clustering = samples are close to their prototype.

方法: ② Prototypes are obtained as the center of each cluster.

② Samples are re-assigned to the cluster with the ~~closest~~ closest prototype.

Note 1. Hyperparameter K 的选择:

② 美林问题

② validation, 但不是一般情况误差差, 因为当 $K=N$, Err=0.

使用 elbow method, when $K > K_T$, The improvement of quality will be slower.

Note 2. k-means 只能 produces spherical clusters. (convex-set), 不可以 Non-Convex 集合.
使用 density-based clustering 来解决。

Algorithm

Density-based spatial clustering of applications with noise.

参数: radius r , threshold δ , 则所有样本分为 3 类.

Core: - \forall x 内, 样本数大于等于 t .

Border: - \forall x 有 core 点, 但是 core 点数少于 t .

Outlier: others.

core 点, 邻域中非 core 点, db-clusters (非核心点的邻域 cluster, 核心点的邻域).

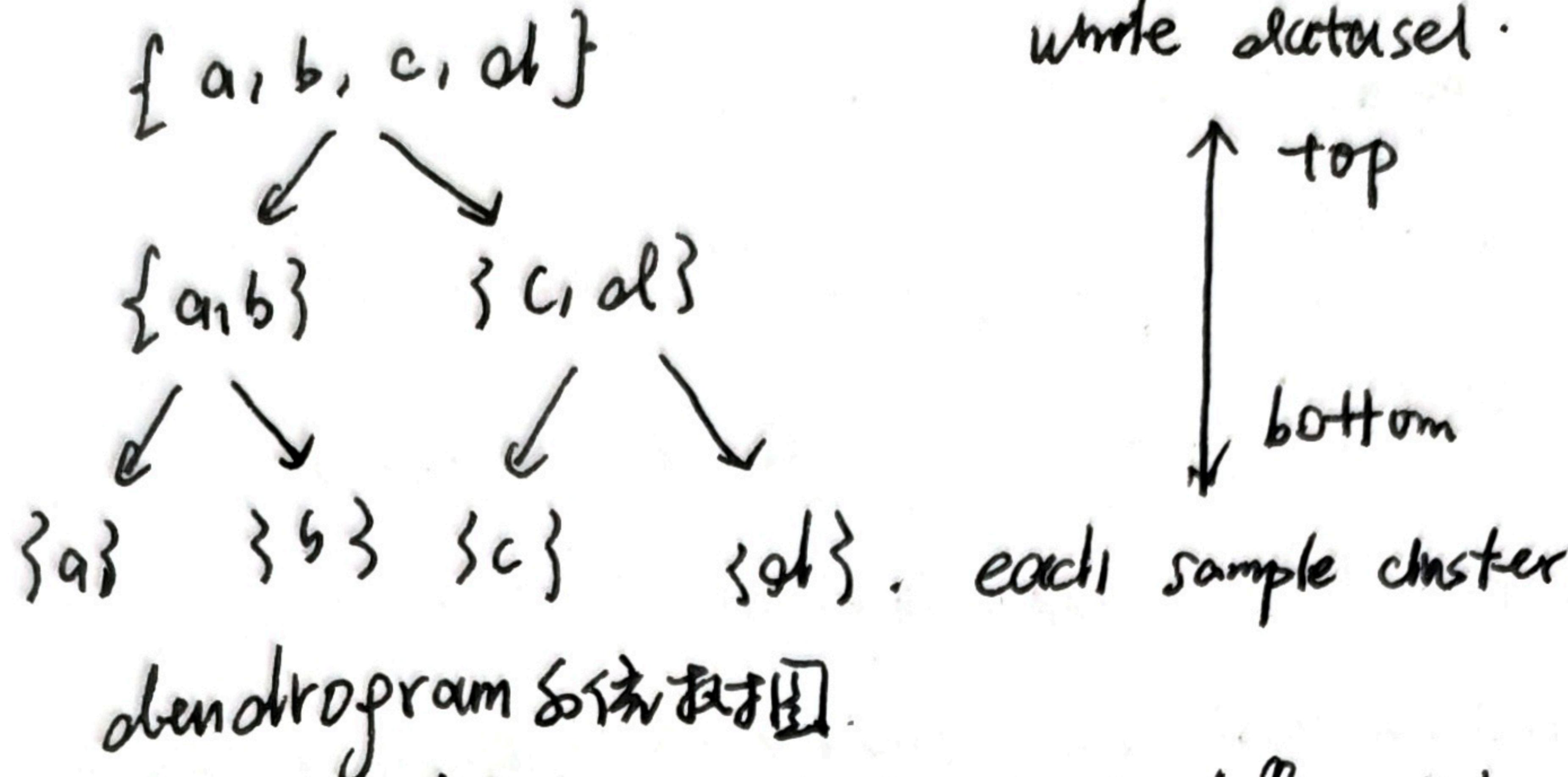
Border 邻域子集域内最多 core 点, 邻域非 cluster. Outlier 不属任何.

Note DBSCAN 不需设置 Number of class, 容错率, 距离和 sample 属于 cluster

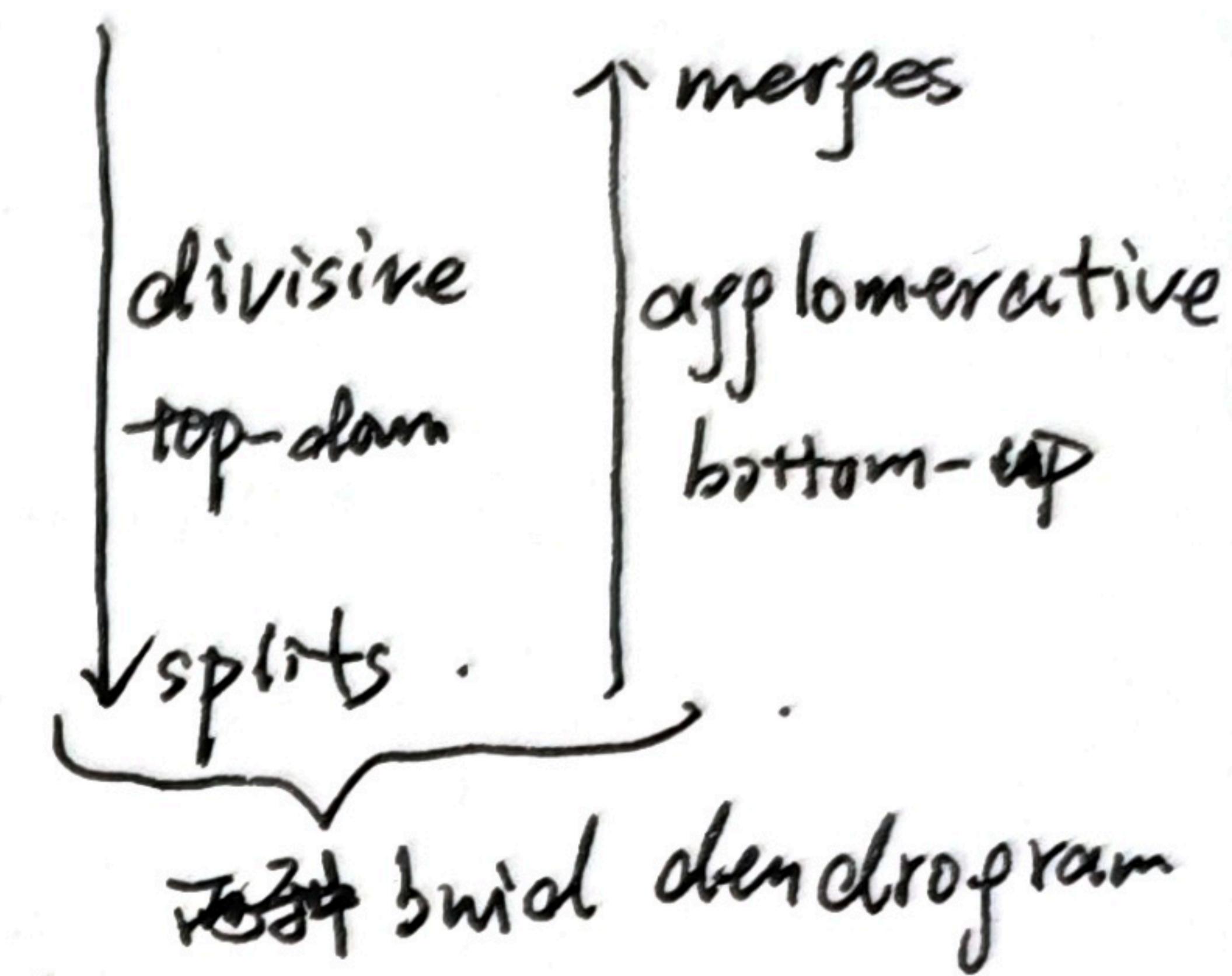
Algorithm

HDBSCAN

hierarchical clustering. The structure of a dataset can be explored at different levels that expose different properties.



whole dataset



Relationship between clusters at different levels.

Note merge or split is option:

① Single linkage $D(A, B) = \min \{d(x_i, y) : x \in A, y \in B\} \rightarrow$ arbitrary shape

② Complete linkage $D(A, B) = \sum_{x \in A, y \in B} d(x, y) \rightarrow$ spherical shape.

③ Group average: $D(A, B) = \frac{\sum_{x \in A, y \in B} d(x, y)}{|A| + |B|} \rightarrow$ shape.

Component analysis allows us to identify the directions in the space

that our data are aligned with (PCA).

应用: transform data, clean it, reduce dimensionality.

2. Density Estimation.

Assumption Samples distributed in the attribute space, and this distribution can be represented or described by probability densities. (不是 model).

Note 该模型 (Probability density) 是 true distribution of data, 而不是一个 model for data in Density Estimation.

Algorithm

Non-parametric

Do not specify the shape of probability.

1) Histogram. $\hat{f}_{1,2} = \frac{1}{Nh} (\text{Number of } x_i \text{ in same bin as } x)$

$b \rightarrow$ small \Rightarrow spiky, $b \rightarrow$ big \Rightarrow flat. (Not flexible to find the underlying structure).

2) kernel density estimation. (KDE). Let $h \rightarrow 0$, And $Nh \rightarrow \infty$.

$$\hat{f}_1(x) = \frac{1}{Nh} \sum_{i=1}^N K\left(\frac{x_i - x}{h}\right).$$

Building an individual density around each sample and combining them.

Algorithm
parametric
method
GMM.

Parametric approaches specify the shape of the probability density function by specifying coefficients.
Gaussian Mixture models (GMM).
 $P(x) = \sum_m \pi_m N(x | \mu_m, \Sigma_m)$, where $\sum_m \pi_m = 1$.

Note 1) Why Gaussian? ① CLT ② $\mathbb{E}[X]$ & $\text{Var}[X]$ are well-defined. ③ independent \Leftrightarrow uncorrelated

$$P(x) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}.$$

2) Gaussian fits & parameters estimation (unbiased).

$$\hat{\mu} = \frac{1}{N} \bar{x}_i, \quad \hat{\Sigma} = \frac{1}{N-1} \bar{x}_i (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

最大化似然，(ML) 条件下， $\hat{\Sigma} = \frac{1}{N} \bar{x}_i (x_i - \hat{\mu})(x_i - \hat{\mu})^T$.

3) Algorithm: Expectation-Maximization (EM).

E-step: it's latent variable posterior, sum of expected likelihood over fix α .

M-step: arg max likelihood (latent variable, α).

Outliers samples just depend on noise \Rightarrow outliers or anomalies.

1) Anomaly detection algorithm quantify the probability, if $P(x_i) < T$, x_i anomaly.

2) 在设计系统时，为了 mitigate outliers 造成的影响，提高系统的鲁棒性，
从而达到鲁棒性设计。

App for classifiers

1) LDA, QDA, use Bayes to learn.

2) Naive Bayes., Assume that predictor are independent

$$h_{\text{Bayes}} = \arg \max_{y \in \{0, 1\}} P(Y=y) P(X=x | Y=y).$$

$$= \arg \max_{y \in \{0, 1\}} P(Y=y) \prod_{i=1}^d P(X_i=x_i | Y=y).$$

即 $P(Y=y) \propto \text{prior}(y) \prod_{i=1}^d P(X_i=x_i | Y=y)$.

App for cluster

看 k-means as a version of GMM with $\Sigma = \sigma^2 I$.

if $\Sigma = \sigma^2 I$, GMM produce spherical decision region.

Maching learning.

I. 最小范数，最小二乘、mp 逆。

通用的，存在 $A \in \mathbb{R}^{m \times n}$, 存在问题

$$\min_x \|Ax - b\|_2^2$$

该问题可应对行满秩和列满秩的情况。
行满秩

Solution. $A = U \Sigma V^T = I V_R \underbrace{U_R}_{\text{3行}} \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_R^T \\ V_N^T \end{bmatrix}$. 零空间
3行 \rightarrow 左零空间

$$\begin{aligned} \min_x \|Ax - b\|_2^2 &= \min_x \|D^T A x - D^T b\|_2^2 \\ &= \min_x \left\| \begin{bmatrix} \Sigma_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_R^T \\ V_N^T \end{bmatrix} x - \begin{bmatrix} V_R^T \\ V_N^T \end{bmatrix} b \right\|_2^2 \end{aligned}$$

$$\begin{cases} \text{proj}(A) = A A^T = (V_R V_R^T) \\ \text{proj}(A^T) = A^T A = V_R V_R^T \end{cases} = \min_x \left\| \Sigma_r V_R^T x - V_R^T b \right\|_2^2 + \|V_N^T b\|_2^2$$

见<补充>。
 $= \min_x \|A x - V_R V_R^T b\|_2^2 + \|V_N^T b\|_2^2$

相当于把 b 投影到 A 的列空间 $X^* = \arg \min_x \|A x - V_R V_R^T b\|_2^2 = V_R \Sigma_r^{-1} V_R^T b \triangleq A^+ b$.

此时, $\min_x \|A x - b\|_2^2 = \|V_N^T b\|_2^2$.

对应的最小值解为 $x = A^+ b + v$, $v \in N(A)$.

<对 A 行满秩>. $\boxed{\quad} \mapsto$ 最小范数.

此时, $m < n$; 方程组欠定, 无解, 需求正则化.

$$\min \|x\|_2^2$$

$$\text{s.t. } Ax = b.$$

Analysis: 上述求解 $\|x\|_2^2 = \|A^+ b\|_2^2 + \|V\|_2^2 \geq \|A^+ b\|_2^2$

$$\text{即 } \min \|x\|_2^2 = \|A^+ b\|_2^2.$$

证明. 可以证明, 最小范数解 $x = A^T (A A^T)^{-1} b$

$$\text{此时 } A = U_R \begin{bmatrix} \Sigma_r & 0 \end{bmatrix} \begin{bmatrix} V_R^T \\ V_N^T \end{bmatrix}.$$

$$\text{代入计算有 } x = A^T (A A^T)^{-1} b = A^+ b.$$

2. 最可能的参数是对的！最大似然！

最大似然估计 (Maximum Likelihood Estimation, MLE).

Assumption: 强假设, 目标分布已知, 只是参数未知 (θ 未知).

Motivation 使用分布 Q wrt. $\hat{\theta}$ to approximate 分布 P wrt. θ 有差异.

(频率学派视角: 每个样本都有不完全的(单值)参数!).

Solution. 选用 KL divergence 来度量 P 和 Q .

$$\begin{aligned} D_{KL}(P_\theta(x) \parallel P_{\hat{\theta}}(x)) & \quad (\text{同分布, 参数不同, } \hat{\theta} \text{ 为估计参数}) \\ = \mathbb{E}_{P_\theta} \log \frac{P_\theta(x)}{P_{\hat{\theta}}(x)} \\ = \mathbb{E}_{P_\theta} \log P_\theta(x) - \mathbb{E}_{P_\theta} \log P_{\hat{\theta}}(x). \end{aligned}$$

由大数定律. $\mathbb{E}_{P_\theta} \log P_{\hat{\theta}}(x) \Rightarrow \frac{1}{N} \sum_i^N \log P_{\hat{\theta}}(x_i)$

$$D_{KL}(P_\theta(x) \parallel P_{\hat{\theta}}(x)) = \mathbb{E}_{P_\theta} \log P_\theta(x) - \frac{1}{N} \sum_i^N \log P_{\hat{\theta}}(x_i).$$

(频率学派) θ 必然有最优值, $\mathbb{E}_{P_\theta} \log P_\theta(x)$ 为 constant.

$$\begin{aligned} \hat{\theta}^* &= \arg \min_{\hat{\theta}} D_{KL}(P_\theta(x) \parallel P_{\hat{\theta}}(x)) \\ &= \arg \min_{\hat{\theta}} -\frac{1}{N} \sum_i^N \log P_{\hat{\theta}}(x_i) \\ &= \arg \max_{\hat{\theta}} \frac{1}{N} \sum_i^N \log P_{\hat{\theta}}(x_i) \\ &= \arg \max_{\hat{\theta}} \frac{1}{N} \log \prod_{i=1}^N P_{\hat{\theta}}(x_i). \end{aligned}$$

$$= \arg \max_{\hat{\theta}} \prod_{i=1}^N \underbrace{P_{\hat{\theta}}(x_i)}_{\text{即 } p_{\hat{\theta}}(x_i)}, \text{指概率 } p_{\hat{\theta}}(x_i | \hat{\theta}).$$

$$\begin{aligned} \text{最终, 我们得到. } \hat{\theta}^* &= \arg \max_{\hat{\theta}} \prod_{i=1}^N P_{\hat{\theta}}(x_i) \\ &= \arg \max_{\hat{\theta}} \text{likelihood}(\hat{\theta}) \end{aligned}$$

Remark, 最后的 object 是 "最大似然", 直观理解其黑幕:

| 当 Sample 足够多时, 在取到的样本 $D = \{(x_i, y_i), i \in \mathbb{N}\}$ 就是
| 最可能在取到的 dataset.

缺点: 假设强 (但仅仅对 Probability 建模而并非).

小样本 大数定律不适用, MLE 很偏高.

Remark 2 对于然数函数讨论不成立为止,

Define likelihood function $L(\theta|x) = p(x|\theta) = P_\theta(x)$.

$$\arg\max_\theta P_\theta(x) = \arg\max_\theta P(x|\theta) = \arg\max_\theta L(\theta|x),$$

所以, $P_\theta(x)$, $P(x|\theta)$, $L(\theta|x)$ 都叫似然函数, 但卒子版中, 用作 $P(x;\theta)$, $P(x|\theta)$ 更多, 因为不说 θ 是参数的.

根据 Bayes formula: $p(\theta|x) = \frac{P(x|\theta) P(\theta)}{P(x)}$. 先验

3. 二分类中的最优化器.

对任 y $S = ((x_1, y_1), \dots)$ where $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ 之 loss.

$$R(h) = \mathbb{E}_{(x,y) \sim S} [\mathbb{I}_{f_h(x) \neq y}]$$

其中, h 为假设集内的一种 function, 二分类器 $f_h(x)$ 为

$$f_h(x) = \begin{cases} 1, & h(x) \geq 0 \\ -1, & h(x) < 0 \end{cases}$$

$$\text{且有 } R(h) = \mathbb{E}_{(x,y) \sim S} [\mathbb{I}_{f_h(x) \neq y}] = \mathbb{E}_{x \sim D} [y(x) \mathbb{I}_{h(x) < 0} + (1-y(x)) \mathbb{I}_{h(x) \geq 0}]$$

Define Bayes optimal classifier $h^*(x) = y(x) - \frac{1}{2}$.

即, 为某 label 的概率取 $1/2$, 例如,

Lemma $\forall h \in \mathcal{H}$, s.t. $R(h) - R(h^*) = 2 \mathbb{E}_{x \sim D} [\mathbb{I}_{h^*(x)} \mathbb{I}_{h(x) h^*(x) \leq 0}]$

proof. $R(h) = \mathbb{E}_{x \sim D} [y(x) \mathbb{I}_{h(x) < 0} + (1-y(x)) \mathbb{I}_{h(x) \geq 0}]$

$$= \mathbb{E}_{x \sim D} [y(x) \mathbb{I}_{h(x) < 0} + (1-y(x)) (1 - \mathbb{I}_{h(x) \geq 0})]$$

$$= \mathbb{E}_{x \sim D} [(2y-1) \mathbb{I}_{h(x) < 0} + (1-y(x))]$$

$$= \mathbb{E}_{x \sim D} [2h^*(x) \mathbb{I}_{h(x) < 0} + (1-y(x))].$$

$$W \leq R(h) - R(h^*)$$

$$= \mathbb{E}_{x \sim D_x} [2h^*(x)(\mathbb{I}_{h(x) < 0} - \mathbb{I}_{h^*(x) < 0})].$$

对表处理.

$$\begin{array}{ccc|c} h(x) & h^*(x) & \mathbb{I}_{h(x) < 0} - \mathbb{I}_{h^*(x) < 0} \\ \hline 2 & 1 & 0 \\ -1 & -1 & 0 \\ 1 & -1 & -1 \\ -1 & 1 & + \end{array} \Rightarrow \begin{aligned} -h^*(x) &= |h^*(x)|, \\ +h^*(x) &= |h^*(x)|. \end{aligned}$$

$$W \leq R(h) - R(h^*)$$

$$= \mathbb{E}_{x \sim D_x} [2|h^*(x)| \mathbb{I}_{h(x)h^*(x) \leq 0}], \text{ (证).}$$

□.

Remark 该 lemma 表示 Bayes classifier 是最优分类器.

其相对于所有 (不B. 二分类), Bayes 是最优的.

$$f_D(x) = \begin{cases} 1, & P(y=1|x=x) > \frac{1}{2} \\ 0, & \text{otherwise.} \end{cases}$$

4. Linear Discriminant Analysis [Intr to ML. Lecture 6].

logistic Regression 模型是 $P(\text{class} | \text{data})$

Discriminant Analysis 模型是 $P(\text{data} | \text{class})$.

A. QDA.

DA 模型是对 likelihood $P(\text{data} | \text{class})$ 做出假设 $f_k(x) \sim N(\mu_k, \Sigma_k)$

它使用算式:

$$\frac{P(x | \text{class} = k)}{P(x | \text{class} = i)} \propto P(\text{class} = k)$$

$$P(\text{class} = k | x) \propto \sum_i P(x | \text{class} = i) P(\text{class} = i)$$

$$\Rightarrow \propto \frac{f_k(x) \pi_{k1}}{P(x)} \propto f_k(x) \pi_{k1}.$$

现在我们假设处理二分类问题, 且 $\pi_{k1} = \pi_{k2}$ ($\pi_1 = \pi_2 = \frac{1}{2}$).

使用 MAP 准则获取 decision boundary

$$P(\text{class } 1 | x) = P(\text{class } 2 | x)$$

$$\Rightarrow P(x | \text{class } 1) \pi_{k1} = P(x | \text{class } 2) \pi_{k2}$$

$$P(x|\text{Class 1}) = P(x|\text{Class 2}) \Leftrightarrow \log f_1(x) = \log f_2(x). \quad ①$$

$$f_k(x) \sim N(\mu_k, \Sigma_k) \Rightarrow f_k(x) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) \right\}.$$

对 ① 式整理，有

$$\log |\Sigma_1| + (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = \log |\Sigma_2| + (x - \mu_2)^T \Sigma_2^{-1} (x - \mu_2)$$

这是由二次型表达的边界，所以称为 QDA。

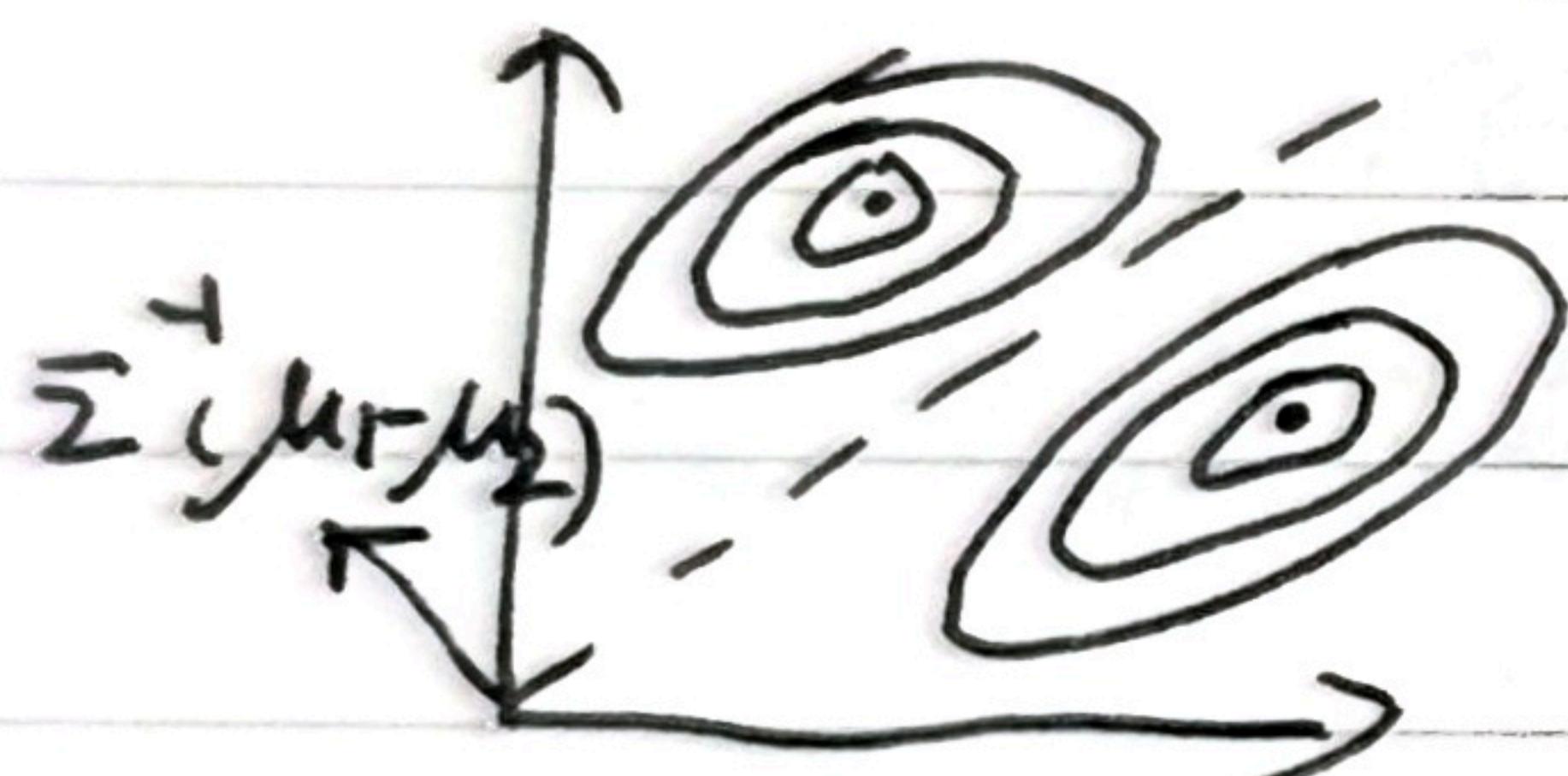
B. LDA.

对 QDA model 增加一条假设： $\Sigma_1 = \Sigma_2 = \Sigma$.

则 ① 式可经整理为：

$$x^T \Sigma^{-1} (\mu_1 - \mu_2) = \frac{1}{2} (\mu_1^T \Sigma^{-1} \mu_1 - \mu_2^T \Sigma^{-1} \mu_2) = \text{Const.}$$

即把 $x^T \Sigma^{-1} (\mu_1 - \mu_2)$ 看成常数。



Σ^{-1} 起到削弱 $(\mu_1 - \mu_2)$ 的作用，如果设有 Σ^{-1} 的作用，或 $\Sigma = \sigma^2 I$ ，那么分离线将退化为中垂线，LDA 退化为 nearest centroid classifier

再者，若 $\Sigma = \sigma^2 I$, $\Sigma_{11} \neq \Sigma_{22}$, 则分离边界以垂直于有质心连结的方式在连线上移动。

C. Estimating the parameters. ($\hat{\mu}_k, \hat{\Sigma}_k$)

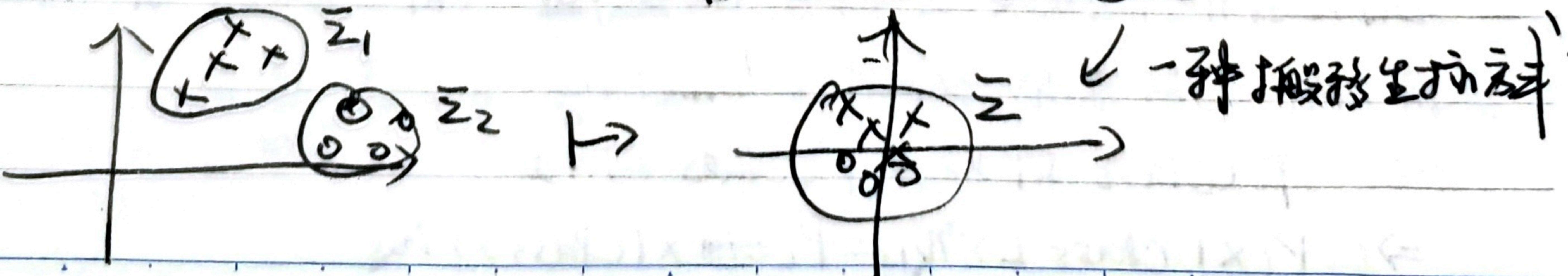
$$\text{QDA: } \hat{\mu}_k = \frac{1}{n_k} \sum_{i \in C_k} x_i$$

$$\hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i \in C_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

$$\hat{\tau}_{nk} = \frac{n_k}{n}.$$

LDA: pooled covariance estimator:

$$\hat{\Sigma} = \frac{1}{n - k} \sum_{k=1}^K \sum_{i \in C_k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$



D. Model Choose.

Σ^{-1} 逆矩阵会有精度问题, $\Sigma^{-1}(\mu_1 - \mu_2)$ 带来

high-variance overfitting 问题, 需要使用

Regularization 方法.

$$\text{Recall: } (\mathbf{X}^T \mathbf{X})^{-1} \mapsto (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1}$$

$$\Rightarrow \Sigma^{-1}(\mu_1 - \mu_2) \mapsto (\Sigma + \lambda \mathbf{I})^{-1}(\mu_1 - \mu_2)$$

$$\Rightarrow ((1-\lambda)\Sigma + \lambda \mathbf{I})^{-1}, \lambda \in [0, 1].$$

可以 mix QDA & LDA.

$$\Sigma^{-1} \mapsto ((1-\lambda)\bar{\Sigma}_k + \lambda \bar{\Sigma})^{-1}.$$

| G covariance matrices | separate | shared |
|-----------------------|-----------------|--------------------|
| Full | QDA | LDA |
| Diagonal | Naive Bayes | Diagonal LDA |
| Spherical | 'Spherical QDA' | 'Nearest Centroid' |

注: Diagonal QDA aka Gaussian Naive Bayes.

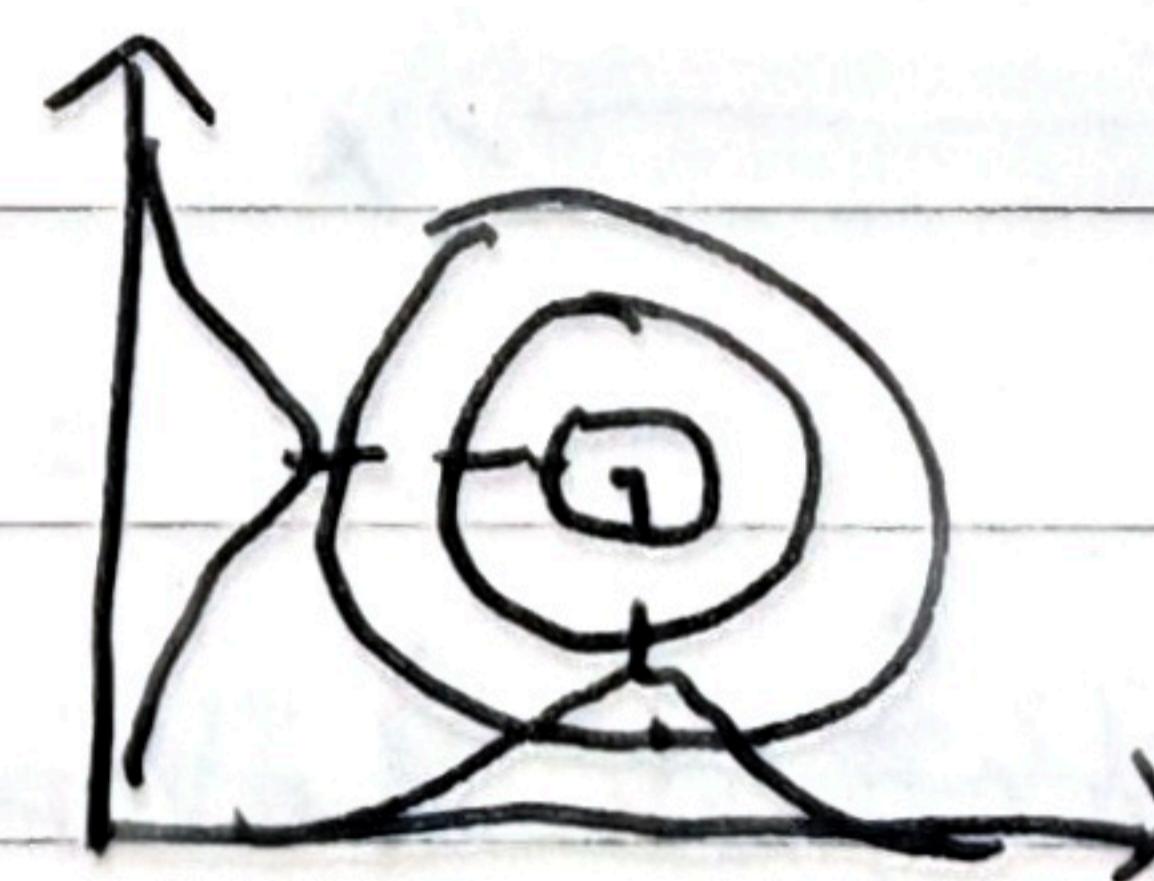
$$\text{Assume } f_k(x) = \prod_{j=1}^p f_{kj}(x_j)$$

去除各个特征相关性! (Naive).

$$f_k \sim N(\mu_k, \Sigma_k), \Sigma_k = \text{diag}\{ \sigma_{k1}^2, \dots, \sigma_{kp}^2 \}, f_{kj} \sim N(\mu_{kj}, \sigma_{kj}^2)$$



→
Naive.



E. Fisher Discriminant Analysis.

LDA 是线性推导.

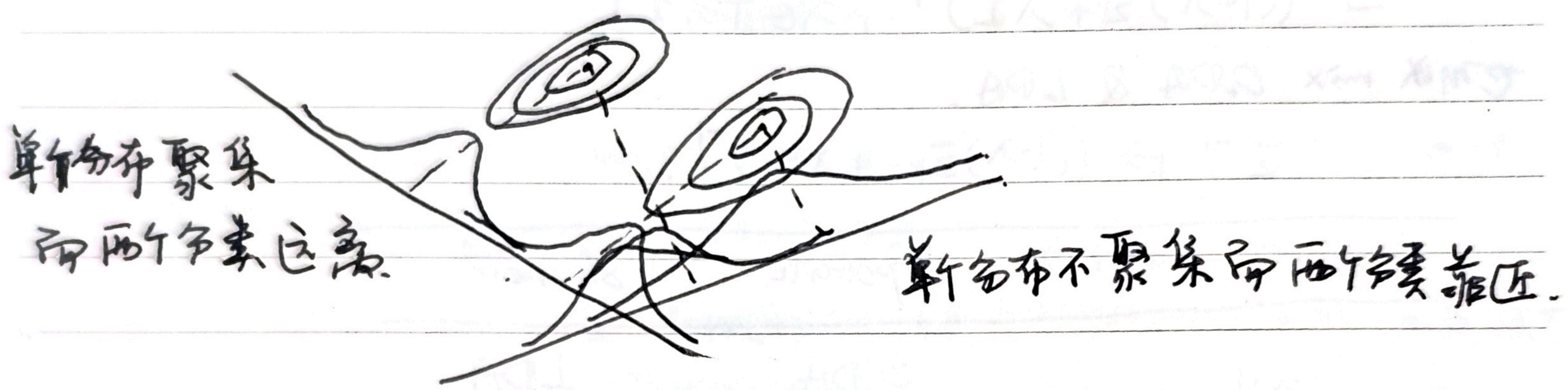
Problem: Find a linear projection that would maximize the ratio of the between-class 'spread' to the within-class 'spread'.

Definition: Fisher's Ratio.

$$\Phi = \frac{(m_1 - m_2)^2}{S_1^2 + S_2^2} \quad \frac{\text{均值之差的平方}}{\text{方差之和}}$$

$$\sim \frac{(w^T(\mu_1 - \mu_2))^2}{w^T \Sigma w}$$

By $\max \Phi \mapsto \hat{w} \sim \Sigma^{-1}(\mu_1 - \mu_2)$.

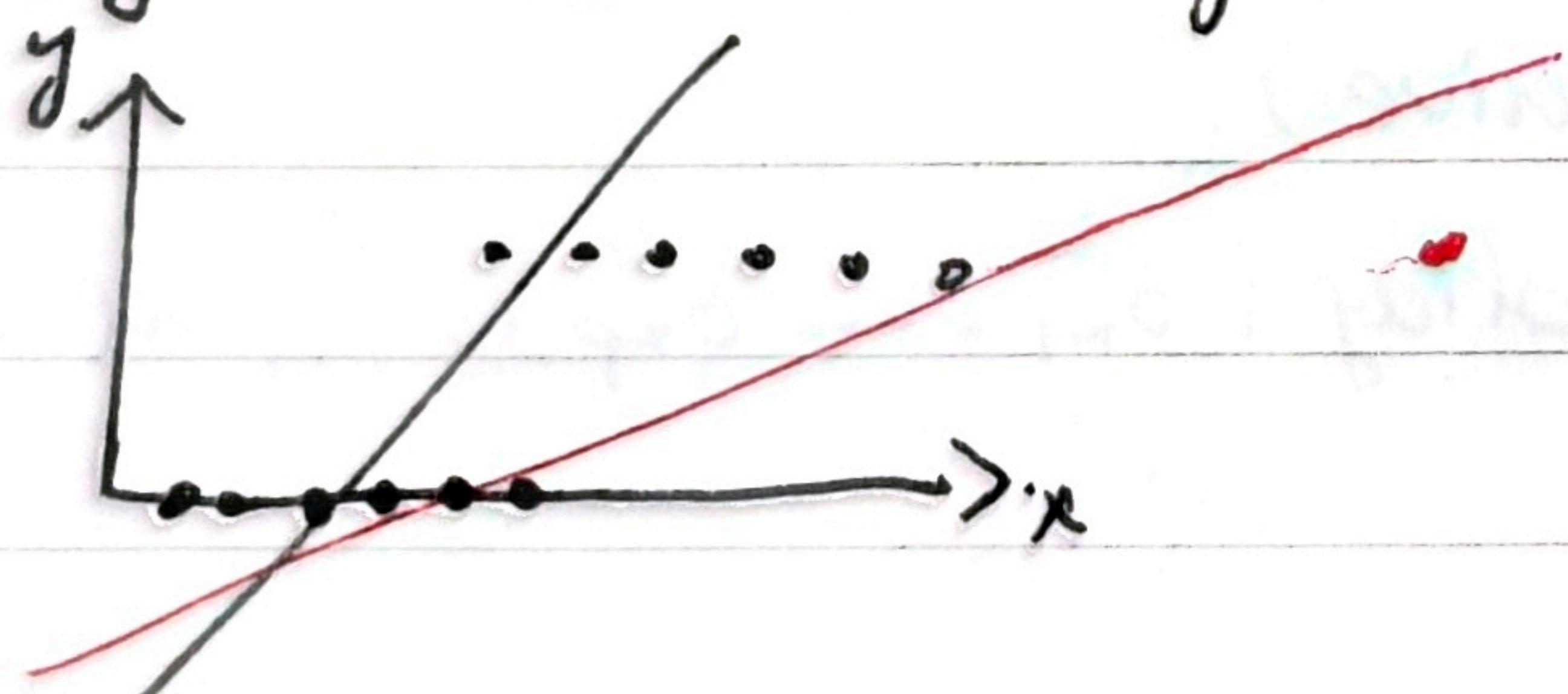


Remark. logistic Regression 与 DA 都能输出无监督结果！

5. Logistic regression [Intro to ML Lecture I] (概率视角).

A. Classification.

Why not use linear regression?



令因高群点产生奇怪的
没有道理的分类情况。

有时，预测概率比真数多来更有意义！

$$f(x) = \frac{1}{1+e^{-x}} \quad \text{logistic} \quad \rightarrow h(x) = f(\beta^T x) = \frac{1}{1+e^{\beta^T x}}$$

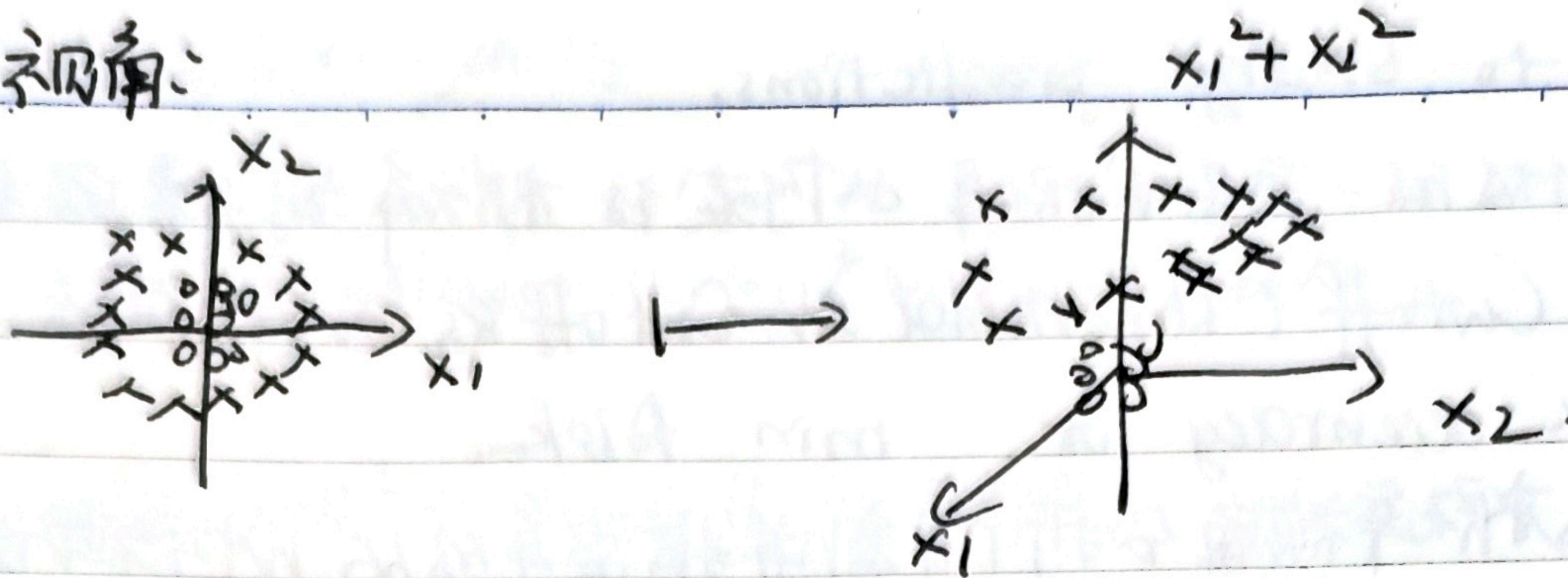
Why regression is classifier?

[ML] 课件PPT中已有n个样本，那就用 $\beta^T x$ 来当决策边界。

该决策的解释是：线性分类后，再由 logistic 函数制一个概率。

即为伯努利分布指定概率，类似于参数估计。

高发视角:



通过添加特征的方差使得线性参数!

B. Loss function and model.

Let $\hat{y} = P(y=1) = \frac{1}{1+e^{-\beta^T x}} = f(\beta^T x) = h(x)$.

为什么不用MSE当loss? $l = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2$.

$\hat{y}_i = 0.99$ 和 $\hat{y}_i = 0.001$ 在 $y_i = 0$ 时有相同的 loss, 无法很好得衡量误差.

使用MLE.

$$\begin{aligned} L &= -\log \text{Likelihood} \\ &= -\sum_{i:y_i=1} \log h(x_i) - \sum_{i:y_i=0} \log (1-h(x_i)) \\ &= -\sum_i [y_i \log h(x_i) + (1-y_i) \log (1-h(x_i))] \end{aligned}$$

Generalized linear model, L 是 convex 的, 但没有闭式解.

C. 优化.

类似于 logistic 的优化是 2 阶方法, 这里分析 GD.

By $f'(x) = g(x)(1-f(x))$

$$f'(x) = \frac{1}{1+e^{-x}} \quad f(x) = \frac{1}{1+e^{-x}}$$

$$\Rightarrow \nabla_{\beta} \log h(x) = \nabla \log f(\beta^T x) = x(1-h(x))$$

$$\nabla_{\beta} \log (1-h(x)) = -h(x)x.$$

$$\nabla L = -\nabla \sum_i [y_i \log h(x_i) + (1-y_i) \log (1-h(x_i))]$$

$$\Rightarrow -\sum_i [x_i(y_i - h(x_i))] \triangleq -X^T(y - \hat{y})$$

D. 补充.

① Overfitting 时, loss $\rightarrow 0$, $\beta \rightarrow \infty$. $\Rightarrow \hat{y} \rightarrow 0$ or $\hat{y} \rightarrow 1$.

需要一些正则化项防止 N 次参数估计.

② Probabilities to binary predictions.

我们在做的一直是二分类，那么在做二分类判决时呢？
需要一个 Cutoff (Threshold)，Cutoff 的选取原则是：

max accuracy or min Risk.

注：这里从概率视角分析问题，并非简单地从 MLE 视角，如此选择 Cutoff 可以融合更多先验信息，使决策更灵活。

③ Multinomial logistic regression.

By Softmax ! $P(y=f) = \frac{e^{\beta_f^T x}}{\sum e^{\beta_k^T x}}$

总结：从概率视角分析 logistic，与 MLE 视角相互照应。

logistic 不止是一种分类标准，其实应该是一种 Regression！

6. 一个关于 k-means 的证明。

Define intra-cluster sample scatter

$$I(C_0) = \frac{1}{2} \sum_{i,j} d_{ij} = \frac{1}{2} \sum_{i,j} \|x_i - x_j\|_2^2$$

Define Center(means) $\mu_0 = \frac{1}{N_0} \sum x_i$

Lemma $I(C_0) = \underbrace{\frac{1}{2} \sum_{i,j} \|x_i - x_j\|_2^2}_{A} = N_0 \sum_i d_i$
 $= N_0 \sum_i \underbrace{\|x_i - \mu_0\|_2^2}_{B}$.

proof.

$$\begin{aligned} A &= \frac{1}{2} \sum_{i,j} \|x_i - x_j\|_2^2 \\ &= \frac{1}{2} \sum_{i,j} [\|x_i\|_2^2 + \|x_j\|_2^2 - 2 \langle x_i, x_j \rangle] \\ &= \frac{1}{2} \sum_i N_0 \|x_i\|_2^2 + \frac{1}{2} \sum_j N_0 \|x_j\|_2^2 - \sum_{i,j} \langle x_i, x_j \rangle \\ &= N_0 \sum_i \|x_i\|_2^2 - \sum_{i,j} \langle x_i, x_j \rangle. \end{aligned}$$

$$B = N_0 \sum_i \|x_i - \mu_0\|_2^2$$

$$= N_0 \sum_i \left[\|x_i\|_2^2 + \|\mu_0\|_2^2 - 2 \langle x_i, \mu_0 \rangle \right]$$

$$= N_0 \sum_i \|x_i\|_2^2 + N_0 \cdot \sum_i \left[\frac{1}{N_0} \sum_j \|x_j\|_2^2 - 2 \frac{N_0}{N_0} \sum_i \langle x_i, \sum_j x_j \rangle \right]$$

$$= N_0 \sum_i \|x_i\|_2^2 + \frac{N_0 \cdot N_0}{N_0^2} \sum_j \langle \sum_k x_k, \sum_k x_k \rangle - 2 \sum_{i,j} \langle x_i, x_j \rangle$$

$$= N_0 \sum_i \|x_i\|_2^2 + \sum_{j,k} \langle x_j, x_k \rangle - 2 \sum_{i,j} \langle x_i, x_j \rangle$$

$$= N_0 \sum_i \|x_i\|_2^2 - \sum_{i,j} \langle x_i, x_j \rangle = A.$$

□

使用该 Lemma 以使用 $I(\mu_0) = N_0 \sum_i \|x_i - \mu_0\|_2^2$ 来简化
计算 intra-cluster sample scatter!

7. Kernel Density Estimation.

Empirical distribution: $\hat{f}(x) = \frac{1}{n} \sum_i \mathbb{I}\{x_i = x\}$.

通过把 $\hat{f}(x)$ 拓展化的累加，得到 $\hat{f}(x)$ 是连续函数。

$$f(x_0) = \lim_{h \rightarrow 0} \frac{F(x_0+h) - F(x_0-h)}{2h}$$

$$= \lim_{h \rightarrow 0} \frac{1}{2h} \left[\frac{1}{N} \sum_i \mathbb{I}\{x_0+h \geq x_i\} - \frac{1}{N} \sum_i \mathbb{I}\{x_0-h \geq x_i\} \right]$$

$$= \lim_{h \rightarrow 0} \frac{1}{2N} \cdot \frac{1}{h} \left[\sum_i \mathbb{I}\{x_0-h \leq x_i \leq x_0+h\} \right]$$

$$= \lim_{h \rightarrow 0} \frac{1}{Nh} \left[\frac{1}{2} \sum_i \mathbb{I}\left\{ \left| \frac{x_i - x_0}{h} \right| \leq 1 \right\} \right].$$

$\triangleq \frac{1}{Nh} \sum_i k\left(\frac{x_i - x_0}{h}\right)$, where $k(\cdot)$ is kernel function.
 h is bandwidth.

$$\begin{aligned}
 \underset{D}{\text{MSE}}\{\hat{f}(x)\} &= \underset{D}{\mathbb{E}}[(\hat{f}(x_0) - f(x_0))^2] \\
 &= [\underset{D}{\mathbb{E}}(\hat{f}(x_0) - f(x_0))]^2 + \underset{D}{\text{Var}}(\hat{f}(x_0) - f(x_0)) \\
 &= \underbrace{[f(x_0) - \underset{D}{\mathbb{E}}\hat{f}(x_0)]^2}_{\text{estimator bias}} + \underbrace{\underset{D}{\mathbb{E}}[\underset{D}{\mathbb{E}}\hat{f}(x) - \hat{f}(x)]^2}_{\text{estimator variance}}
 \end{aligned}$$

通过式子 $\text{MISE} = \underset{D}{\mathbb{E}} \int_{-\infty}^{+\infty} (\hat{f}(x) - f(x))^2 dx$ 理解 MISE.

$$\begin{aligned}
 &= \int \text{bias}^2 dx + \int \text{var}^2 dx.
 \end{aligned}$$

当 h 变大时, $\int \text{bias}^2 dx$ 变大, $\int \text{var}^2 dx$ 变小., 那

- Large h : $\hat{f}(x)$ is smoother (low model flexibility), low variance, high bias.

- Small h : $\hat{f}(x)$ is less smoother (high model flexibility, spiky), high variance, low bias.

Theorem Consistency when $N \rightarrow \infty, h \rightarrow 0, Nh \rightarrow \infty$, exist
 $\text{bias} \rightarrow 0, \text{var} \rightarrow 0, \hat{f}(x) \xrightarrow{P} f(x)$.

Asymptotic Normality when $N \rightarrow \infty, h \rightarrow 0, Nh \rightarrow \infty$, exist
 $\sqrt{Nh}(\hat{f}(x) - f(x)) \xrightarrow{d} N(0, f''(x) \int_k(z) dz)$.

Note 通过 $\min \text{MISE}$ 选择 h , 一般选择 h^* 称为 kernel: Spannweite.
 $h^* = \arg \min_h \text{MISE}(\hat{f})$.