

Singular Value Decomposition (SVD)

▶ For an arbitrary matrix $X \in \mathbb{R}^{m \times n}$ there exists a factorization as follows:

$$X = U\Sigma V$$

 $U \in \mathcal{R}^{m \times m}, V \in \mathcal{R}^{n \times n}, UU^T = U^TU = I, VV^T = V^TV = I$ diagonal matrix $\Sigma \in \mathcal{R}^{m \times n}$

• If
$$rank(X) = d$$

 $U \in \mathcal{R}^{m \times d}, V \in \mathcal{R}^{d \times m}, U^T U = I, WV^T = I$
 $\Sigma = \operatorname{diag}(\sigma_1, \sigma_2, \cdots \sigma_d) \ \sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d > 0$

Minimize the mean-squared error:

$$J_n = \frac{1}{n} \sum_{i=1}^{n} (y_i - f(x_i))^2$$

Minimize the residual sum of squares

$$RSS(f) = \sum_{i=1}^{n} (y_i - f(\mathbf{x}_i))^2$$

Optimizing the MSE Criterion

Computing the gradient gives:

$$\nabla J_n = -2X(\mathbf{y} - X^T \mathbf{a})$$

Setting the gradient to zero

$$XX^{T} \mathbf{a} = X\mathbf{y}$$
$$\mathbf{a} = (XX^{T})^{-1} X\mathbf{y}$$

Any problems?

- What is the rank of the matrix XX^T ?
- The solution for a can be obtained uniquely if XX^T is non-singular
- The fitted values at the training inputs are

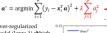
$$\hat{\mathbf{y}} = X^T \mathbf{a} = X^T (XX^T)^{-1} X \mathbf{y}$$
$$\mathbf{a} = (XX^T)^{-1} X \mathbf{y}$$

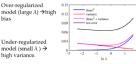
The solution for a can be obtained uniquely if XX^T is non-singular.

LASSO: 平方变绝对值



 $\boldsymbol{a}^* = \left(XX^T + \lambda \boldsymbol{I}\right)^{-1} X \boldsymbol{y}$





Logistic Regression

$$\begin{split} P(y_t = \pm 1 | \mathbf{x}_t, \mathbf{a}) &= \sigma \big(y_t \mathbf{a}^T \mathbf{x}_t \big) = \frac{1}{1 + e^{-y_t \mathbf{a}^T \mathbf{x}_t}} \\ P(D) &= \prod_{l \in I} \sigma \big(y_l \mathbf{a}^T \mathbf{x}_t \big) \\ l(P(D)) &= \sum_{l \in I} \log \Big(\sigma \big(y_l \mathbf{a}^T \mathbf{x}_t \big) \Big) = - \sum_{l \in I} \log \Big(1 + e^{-y_l \mathbf{a}^T \mathbf{x}_t} \Big) \end{split}$$

Objective Function of Linear Regression:

$$E(\boldsymbol{a}) = \sum_{i \in I} (y_i - \boldsymbol{a}^T \boldsymbol{x}_i)^2$$

$$E(\boldsymbol{a}) = \sum_{i \in I} \log \left(1 + e^{-y_i \boldsymbol{a}^T \boldsymbol{x}_i} \right) + \lambda \sum_{i=1}^{p} |\boldsymbol{a}_{ij}^{2}|$$

Two-category Linearly Separable Case

- - w^Tx > 0 for examples from the positive class w^Tx < 0 for examples from the negative class.

$$\min_{\mathbf{w},b} \frac{1}{2} ||\mathbf{w}||^2 + C \sum_{i=1}^n \xi_i$$
$$y(\mathbf{w}^T \mathbf{x}_i + b) \ge 1 - \xi_i$$
$$\xi_i \ge 0$$

Hyper plane of maximum margin is supported by those points (vectors) on the margin. Those are called Support Vectors. Nonsupport vectors can move freely without affecting the position of the hyperplane as long as they don't exceed the margin.

Minimum Error Rate Classification

$$R(lpha_x \mid oldsymbol{x}) = 1 - P(\omega_i \mid oldsymbol{x})$$

$$R(lpha_i \mid oldsymbol{x}) = \sum_{j=1}^c \lambda(lpha_i \mid \omega_j) P(\omega_j \mid oldsymbol{x})$$

Bayes Risk

$$R = \sum_{m{x}} R(lpha_i \mid m{x})$$

· Bias-variance Decomposition Expected prediction error (expected loss) = (bias)2+variance+noise

> Weakness of the Original Model

SVM缺点

- · When an outlier appear, the optimal hyperplane may be pushed far away from its original/correct place.
- \circ Assign a slack variable ξ to each data point. That means we allow the point to deviate from the correct margin by a distance of ξ

条件风险 condition risk

$$R(\alpha_i|\mathbf{x}) = \sum_{j=1}^{c} \lambda(\alpha_i|\omega_j) P(\omega_j|\mathbf{x})$$

MLE 高斯分布

$$\hat{\mu} = \frac{1}{n} \sum_{k=1}^{n} \mathbf{x}_k$$

$$\widehat{\Sigma} = \frac{1}{n} \sum_{k=1}^{n} (\mathbf{x}_k - \hat{\boldsymbol{\mu}}) (\mathbf{x}_k - \hat{\boldsymbol{\mu}})^t$$

$$SSE = \sum_{i=1}^{K} \sum_{x \in C_i} dist(c_i, x)^2$$

$$\log \prod_{i=1}^{N} p(\mathbf{x}^{(i)}; \mathbf{\Theta}) = \sum_{i=1}^{N} \log \left(\sum_{k=1}^{K} \pi_{k} \mathcal{N}(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{k}, \boldsymbol{\Sigma}_{k}) \right)$$

EM算法

(E-step) for each i:

 $Q^i\big(\boldsymbol{z}^{(i)}\big) = p\big(\boldsymbol{z}^{(i)}|\boldsymbol{x}^{(i)};\boldsymbol{\theta}\big)$

$$\boldsymbol{\theta} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^{m} \sum_{\boldsymbol{x}^{(i)}} Q^{i}(\boldsymbol{z}^{(i)}) \log \frac{p(\boldsymbol{x}^{(i)}, \boldsymbol{z}^{(i)}; \boldsymbol{\theta})}{Q^{i}(\boldsymbol{z}^{(i)})}$$

SVD与PCA不同,PCA是对 数据的协方差矩阵进行矩 阵的分解,而SVD是直接在 原始矩阵上进行的矩阵分 解。并且能对非方阵矩阵 分解,得到左奇异矩阵U、 sigma矩阵Σ、右奇异矩阵 VT。

$$\Sigma = \sum_{i=1}^{2} \sum_{j \in \omega_i} \frac{(x_j - \mu_i)(x_j - \mu_i)^T}{N}$$

Knn 目标函数

$$y\!=\!arg\;max_{c_j}\sum_{x_i\in N_b(x)}\!I(y_i\!=\!c_j)$$

proposition that this is a good thing. · Empirically it works very, very well.

KNN:TAKING MAJORITY VOTE

· There's some theory (using VC dimension) that is related to the

If we've made a small error in the location of the boundary, this gives

Examples

Max margin 优点

Intuitively this feels safest

- Linear kernels $\kappa(x, x') = x^T x'$

us least chance of causing misclassification

- Polynomial kernels $\kappa(x, x') = (x^T x' + 1)^d$
- $\kappa(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} \mathbf{x}'\|^2}{2\sigma^2}\right)$

The information gain of an attribute a is the expected reduction in entropy caused by partitioning on this attribute.

$$Gain(S,a) = \underbrace{Entropy}(S) - \sum_{v \in values(a)} \frac{|S_v|}{|S|} Entropy(S_v)$$

• where S_v is the subset of S for which attribute a has value v

An Illustrative Example (2)

			Humidi	ity Wind	PlayTe	ennis
			High	Weak	No	
Humidity		Wind	High	Strong	No	
	\wedge	\wedge	High	Weak	Yes	
/			High	Weak	Yes	
			Normal	Weak	Yes	
High	Normal	Weak Strong	Normal	Strong	No	9+,5-
3+,4-	6+,1-	6+2- 3+.3-	Normal	Strong	Yes	E=.94
		/-	High	Weak	No	E94
E=.985	E = .592	E=.811 E=1.0	Normal	Weak	Yes	
Calmica	Constitution of	Gain(S,Wind)=	Normal	Weak	Yes	
Gain(S,Humidity)= .94 - 7/14 0.985		.94 - 8/14 0.811	Normal	Strong	Yes	
	0.592=	- 6/14 1.0 =	High	Strong	Yes	
0.151		0.048	Normal	Weak	Yes	
			High	Strong	No	
			∇	$ S_v _{\Gamma_v}$		

Gain(S_{sunny}, Humidity) = .97-(3/5) 0-(2/5) 0 = .97Gain(S sunny, Temp) = Rain Overcast 3,7,12,13 Avoid Overfitting A

utlook	? Temperatu			/ - 1
utlook	Temperatu	77 11		
		re Humidi	ity Wind	PlayTeni
ınny	Hot	High	Weak	No
inny	Hot	High	Strong	No
inny	Mild	High	Weak	No
inny	Cool	Normal	Weak	Yes
inny	Mild	Normal	Strong	Yes
1	nny nny nny	nny Hot nny Mild nny Cool	nny Hot High nny Mild High nny Cool Normal	nny Hot High Strong nny Mild High Weak nny Cool Normal Weak

Why Bagging Works: Bias & Variance

 $(E_D(f(\mathbf{x};D)) - E(y|\mathbf{x}))^2 + (E_D\{[f(\mathbf{x};D) - E_D(f(\mathbf{x};D))]^2\}$

Let f denote the ground-truth function and h(x) denote a learner trained from the bootstrap distribution D_{bs} . The aggregated learner generated by Bagging is $H(x) = E_{D_{bs}}[h(x)]$

- Bias: $f_D[h(x)] f \approx H(x) f$
- * Assume there are T bs samples, if data is i.i.d. and each variance is σ^2 . The ensemble variance is $\frac{1}{T}\sigma^2$ $H(x) = E_{D_{20}}[h(x)] = \frac{1}{T}\sum^T h_i(x)$

$$\begin{split} & \operatorname{Var}(H(x)) = \operatorname{Var}\left(\frac{1}{T}\sum_{i}^{T}h_{i}(x)\right) = \frac{1}{T^{2}}\operatorname{Var}\left(\sum_{i=1}^{T}h_{i}(x)\right) = \frac{1}{T^{2}}\sum_{i=1}^{T}\operatorname{Var}(h_{i}(x)) = \frac{1}{T}\sigma^{2} \end{split}$$
 $& \quad \text{In reality, Bootstrap Sampled data is } i.d. \text{ (not necessarily independent)}$

Construct Random Forest

- Let N be the number of trees and K be the feature subset size.
- For each N iterations: (building a tree in forest)
 - 1. Select a new bootstrap sample from original training
 - 2. Growing a tree...
 - 3. At each internal node, randomly select K features from ALL features and then, determine the best split in ONLY the K features.
 - 4. Do not pruning
- 5. Until Test Error never decrease (here means Validation error in RF) At last, overall prediction as the average(or vote) from N trees.

Clusters

Kmeans: (kndt)

- Partition objects into k nonempty
- Compute seed points as the centroids of the clusters of the current partitioning
- Assign each object to the cluster with the nearest seed point
- Go back to Step 2, stop when he assignment does not change
- 1、随机选取K个质心的值
- 2、计算各个点到质心的距离
- 3、将点的类划分为离他最近的质心,形成K个 cluster
- 4、根据分类好的cluster, 在每个cluster内重新 计算质心(平均每个点的值)
- 5、重复迭代2-4步直到满足迭代次数或误差小 于指定的值

K-Medoids (kndt): Instead of taking the mean value of the object in a cluster as a reference point(不在内部), medoids can be used, which is the most centrally located object in a cluster. 1、随机选取K个质心的值 (质心必须是某些样本点的值, 而不是任意值)

2、计算各个点到质心的距离 3、将点的类划 分为离他最近的质心,形成K个cluster 4、根据分类好的cluster,在每个cluster内重新计算质 心: 4.1 计算cluster内所有样本点到其中一个 样本点的曼哈顿距离和(绝对误差) 4.2 选出使 cluster绝对误差最小的样本点作为质心 5、重 复迭代2-4步直到满足迭代次数或误差小于指 定的

值.PAM(n^2dt):0. Calculate the pair-wise distance matrix W

- 1.Initialize: randomly select of the data poi nts as the medoids
- 2. Associate each data point to the closest medoid
- 3. For each cluster, compute its medoid 4. Repeat 2-3 until there is no change in the medoids(怎么选?离之心最近的或者 distance row sum 最小的)

GMM:

$$egin{aligned} p(x;\Theta) &= \sum_{k=1}^{K} \pi_k p_k(x; heta_k) \ \Theta &= \left\{\pi_1, \dots, \pi_k, heta_1, \dots, theta_k\right\}, \sum \pi_k = 1 \ p_k(x; heta_k) &= \mathcal{N}(x; \mu_k; \Sigma_k) \end{aligned}$$

For each point $\mathbf{x}^{(i)}$, associate with a hidden variable $\mathbf{z}^{(i)}$ denotes which Gaussian $x^{(i)}$ belongs to

Spectral clustering:

Divide vertices into two disjoint groups (A,B). Minimize weight of between-group connections

$$\begin{split} cut(A,B) &= \sum_{i \in A, j \in B} w_{ij} \\ assoc(A,V) &= \sum_{i \in A, j \in V} w_{ij} \\ \min Ncut(A,B) &= \frac{cut(A,B)}{assoc(A,V)} + \frac{cut(A,B)}{assoc(B,V)} \end{split}$$

Max cluster间的, min cluster之间的

Generalized Eigen-problem

$$\min_{\mathbf{y}} \frac{\mathbf{y}^T(D-W)\mathbf{y}}{\mathbf{y}^TD\mathbf{y}}, \mathbf{y} \in \mathbf{\mathcal{R}}^n, \mathbf{y}^TD\mathbf{1} = 0$$

Obj of NCut

 $(D - W)\mathbf{v} = \lambda D\mathbf{v}$

► Eigenvector corresponding to the smallest eigenvalue.

Vector 1 is the eigenvector corresponding to the eigenvalue 0. $(D-W)\mathbf{y} = \lambda D^{\frac{1}{2}}D^{\frac{1}{2}}\mathbf{y}$ $D^{-\frac{1}{2}}(D-W)D^{-\frac{1}{2}}D^{\frac{1}{2}}\mathbf{y} = \lambda D^{\frac{1}{2}}\mathbf{y}$ $D^{-\frac{1}{2}}(D - W)D^{-\frac{1}{2}}z = \lambda z$

▶ The eigenvector corresponding to the 2nd small eigenvalue

Spectral Clustering Algorithm

- Graph construction
 - Heat kernel $w_{ij} = \exp \left\{ -\frac{\|x_i x_j\|}{2\sigma^2} \right\}$
- k-nearest neighbor graph
- Eigen-problem
- Compute eigenvalues and eigenvectors of the matrix L
- Map each point to a lower-dimensional representation based on one or more eigenvectors.
- Conventional clustering schemes, e.g. K-Means
- · Assign points to two or more clusters, based on the new representation.

D = EuDist2(X);

 $W = (D \le \max(\min(D, k+1)));$

 $W = W & (D \le threshold);$

W = W - eye(size(D)); % clear the distance to

W = double(W | W'); % merge the two-direction matrix

D = eye(size(W));

D = D .* sum(W,2);

[Y, ~] = eigs(D - W, D, 2, 'smallestabs');

 $idx = kmeans(Y(\cdot 2) k)$

$$f(x) = w^T x + w_0$$
$$y = sign(f(x))$$

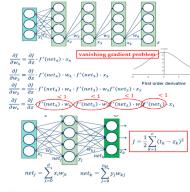
- - o If correct, do nothing
 - \circ If wrong, $w_t = w_{t-1} + xy$
- · Perceptron Criterion Function

$$E(w) = -\sum_{i \in I_M} w^T x_i y_i$$

Gradient Descent (Steepest descent)

$$x_{n+1} = x_n - \gamma \nabla f(x_n), n \geqslant 0$$

Backpropagation Algorithm



$$\frac{\partial J}{\partial w_{kj}} = \frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{kj}} = (z_k - t_k) \cdot f'(net_k) \cdot y_j$$

$$\begin{split} \frac{\partial J}{\partial w_{ji}} &= \frac{\partial J}{\partial y_{j}} \cdot \frac{\partial y_{j}}{\partial net_{j}} \frac{\partial net_{j}}{\partial w_{ji}} \\ &= \left(\sum_{k=1}^{c} (z_{k} - t_{k}) \cdot f'(net_{k}) \cdot w_{kj}\right) \cdot f'(net_{j}) \cdot x_{i} \end{split}$$

Activation Function

$$f(net) = a \, \tanh(b \; net) = a \left[\frac{1 - e^{b \; net}}{1 + e^{b \; net}} \right] = \frac{2a}{1 + e^{-b \; net}} - a$$



Activation function: ReLU(x) = max (0, x)

PCA:

Main steps for computing PCs:

- Form the covariance matrix S.
- Compute its eigenvectors: {a_i}^p_{i=1}
- Use the first d eigenvectors {a_i}_{i=1}^d to form the d PCs.
- The transformation *A* is given by $A = [a_1, \dots a_d]$

几种常见的数据挖掘算法

- Naïve Bayesian classifier
- Linear Regression
- Logistic Regression
- SVM
- Perceptron
- Neural Network
- k Nearest Neighbor
- Decision Tree

正太分布:

Multivariate density: $N(\mu, \Sigma)$ (with dimension d)

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

- $\mathbf{x} = [x_1, \cdots, x_d]^T$
- Σ : $d \times d$ covariance matrix, $|\cdot|$: determinant

MIF:

 $P(x_i \mid \omega_j) =$

Maximum-Likelihood Estimation

- Use the information in training samples to estimate $\theta = (\theta_1, \theta_2, ..., \theta_c); \theta_i \ (i=1,2,...,c)$ is associated with the i-th category
- Suppose sample set D contains n iid samples, $x_1, x_2, ..., x_n$

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

- p(D|θ) is called the likelihood of θ w.r.t. the set of samples.
- ML estimate of θ is, by definition, the value θ that maximizes p(D |
 - "It is the value of 0 that best agrees with the actually observed training samples"

Bayes Estimation:

- We assume that the true values of the a priori probabilities are known or obtainable from a trivial calculation:

 We substitute $P(\omega_1) = P(\omega_1|D)$

Furthermore, we can separate the training samples by class into c subsets D_1, D_2, \dots, D_c , with the samples in D_i belonging to ω_i

$$P(\omega_i|\mathbf{x}, \mathcal{D}) = \frac{p(\mathbf{x}|\omega_i, \mathcal{D}_i) P(\omega_i)}{\sum_{j=1}^{c} p(\mathbf{x}|\omega_j, \mathcal{D}_j) P(\omega_j)}.$$

In essence, we have c separate problems of the following form: use a set D of samples drawn independently according to the fixed but unknown probability distribution p(x) to determine

This is the central problem of Bayesian learning

Bayesian Parameter Estimation: General Theory

 $P(x\mid \textbf{D})$ computation can be applied to any situation in which the unknown density can be parametrized: the basic assumptions are:

- The form of $P(x \mid \theta)$ is assumed known, but the value of θ is not known exactly
- Our knowledge about θ is assumed to be contained in a known prior density $P(\theta)$
- The rest of our knowledge about θ is contained in a set D of n random variables $x_1, x_2, ..., x_n$ that follows P(x)

Discriminant Functions for the Normal Density

The minimum error-rate classification can be achieved by the discriminant function

$$g_i(\mathbf{x}) = \ln P(\mathbf{x}|\omega_i) + \ln P(\omega_i)$$

▶ In case of multivariate normal densities

case of multivariate normal densities
$$P(\mathbf{x}|\omega_l) = \frac{1}{(2\pi)^{\frac{d}{2}}|\Sigma_l|^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_l)^T \Sigma_l^{-1}(\mathbf{x} - \boldsymbol{\mu}_l)\right]$$

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

$$P(error) = \int_{R_2} p(x \mid \omega_1) P(\omega_1) dx + \int_{R_1} p(x \mid \omega_2) P(\omega_2) dx$$

$$\begin{split} P(correct) &= \sum_{i=1}^{c} P(\mathbf{x} \in \mathcal{R}_{i}, \omega_{i}) \\ &= \sum_{i=1}^{c} P(\mathbf{x} \in \mathcal{R}_{i} | \omega_{i}) P(\omega_{i}) \\ &= \sum_{i=1}^{c} \int_{\mathcal{D}_{i}} p(\mathbf{x} | \omega_{i}) P(\omega_{i}) \ d\mathbf{x}. \end{split}$$

Naïve Bayes Classifier

- Given $\mathbf{x} = (x_1, \dots x_p)^T$
 - Goal is to predict class ω
 - lacksquare Specifically, we want to find the value of ω that maximizes $P(\omega|\mathbf{x}) = P(\omega|x_1, \dots x_p)$

$$P\big(\omega|x_1,\cdots x_p\big) \propto P\big(x_1,\cdots x_p|\omega\big)P(\omega)$$

Independence assumption among features

$$P(x_1, \dots x_p | \omega) = P(x_1 | \omega) \dots P(x_p | \omega)$$

How to Estimate Probabilities from Data?

Tid	Refund	Marital Status	Taxable Income	Evade
1	Yes	Single	125K	No
2	No	Married	100K	No
3	No	Single	70K	No
4	Yes	Married	120K	No
5	No	Divorced	95K	Yes
6	No	Married	60K	No
7	Yes	Divorced	220K	No
8	No	Single	85K	Yes
9	No	Married	75K	No
10	No	Single	90K	Yes

- Class: $P(\omega_k) = \frac{N_{\omega_k}}{N}$
 - e.g., P(No) = 7/10, P(Yes) = 3/10
- For discrete attributes $P(x_i|\omega_k) = \frac{|x_{ik}|}{N_{ii}}$
 - where $|x_{ik}|$ is number of instances having attribute x_i and belongs to class ω_k
 - Examples:
 - P(Status=Married | No) = 4/7 P(Refund=Yes | Yes)=0