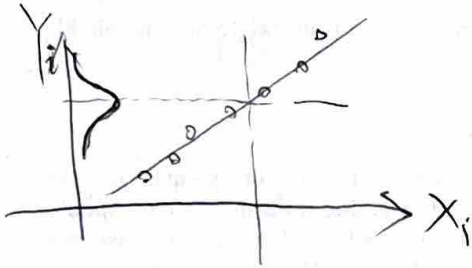


linear regression:

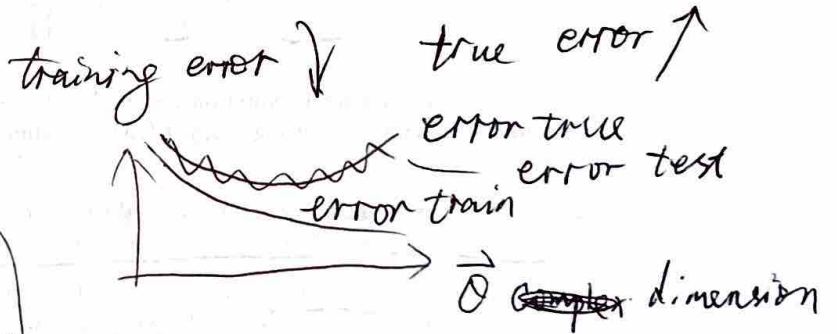
$$\theta = \begin{bmatrix} w \\ b \end{bmatrix}$$

$$P(\vec{Y} | \vec{X}, \theta) = \prod_j \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma^2} \left[Y_j - \sum_i w_i x_{ij} \right]^2\right\}$$

MLE: $\frac{\partial \ln P(\vec{Y} | \vec{X}, \theta)}{\partial \theta} = 0 \Rightarrow$ closed form for w, b

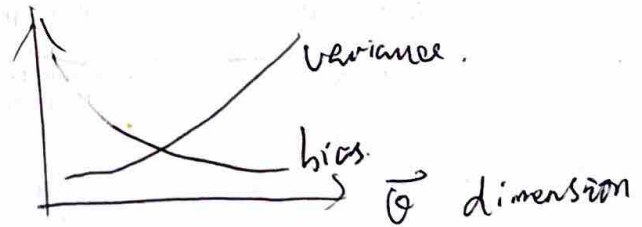


Study on θ : θ dimension



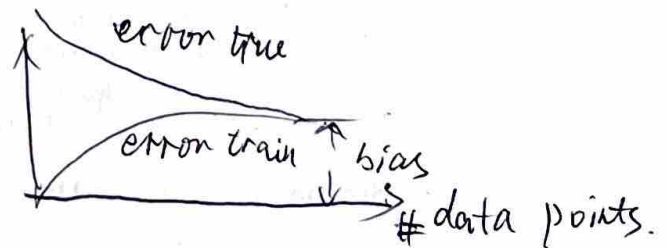
regularization

- $L_1 \Rightarrow$ sparse, less variance.
- $L_2 \Rightarrow$ same variance, more bias.



error true = bias + variance.

Cross-validation for λ using validation data or using training data with LOOCV or k-fold



consistent: $\lim_{\# \text{ data} \rightarrow \infty} (\text{bias}) = 0$

variance can be high

	ridge		LASSO	
	large λ	small λ	large λ	small λ
e_{train}	↑	↓	↑	↓
e_{test}	↑	↓	↑	↓
value of \hat{w}	↓	↑	↓	↑
# of nonzeros in \hat{w}	↓	↑	↓	↑

stay put

logistic regression \Rightarrow classification.

$\theta = w$

$P(Y=0 | \vec{X}, \theta)$

$= \frac{1}{1 + \exp(w \cdot X + w_0)}$

$P(Y=1 | \vec{X}, \theta) = \frac{\exp(w \cdot X + w_0)}{1 + \exp(w \cdot X + w_0)}$

$\ln P(\vec{Y} | \vec{X}; \theta) = \ln \prod P(y_j | x_j, \theta) = \sum_j y_j (w_0 + X_j w) - \ln [1 + \exp(w_0 + X_j w)]$

MLE: $\frac{\partial \ln P(\vec{Y} | \vec{X}; \theta)}{\partial \theta} = 0 \Rightarrow$ no closed form \Rightarrow gradient descent.

convex concave.

L2 regularization

online \Rightarrow stochastic gradient descent.

ada boost.

$w_0^j = \frac{1}{N}$ $\sum w_t^j = 1$

$w_{t+1}^j = w_t^j \exp(-\alpha_t y^j h_t(x^j))$

$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$

$\epsilon_t = \sum_{j=1}^N w_t^j \mathbb{1}[h_t(x^j) \neq y^j]$

$Z_t = \prod_{i=1}^t \sum_{j=1}^N w_{i-1}^j \exp(\alpha_{i-1} y^j h_{i-1}(x^j))$

$w_j^{(t)} = \frac{w_j^{(0)} \exp(\sum_{i=1}^t \alpha_{i-1} y^j h_{i-1}(x^j))}{Z_t}$

$y^j = h(x^j)$

closed form!

(2)

$$\sum_{j=1}^N \mathbb{1}[\cancel{h(x_j)} \neq y_j]$$

$$\sum_{j=1}^N \mathbb{1}\left[\text{sign}\left(\sum_t \alpha_t h_t(x_j)\right) \neq y_j\right]$$

$$\frac{1}{N} \sum_{j=1}^N \exp\left[-y_j \left(\sum_{t=1}^T \alpha_t h_t(x_j)\right)\right]$$

(~~Not~~ MLE) could be.

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

ensemble
vote

include { bagging, boosting, random forest }

① train h_t over $D_t \cdot X$

② $Y = h(x)$

$$D_{t+1}(j) = \frac{D_t(j) \exp(-\alpha_t y_j h_t(x_j))}{\sum_{j=1}^N D_t(j) \exp(-\alpha_t y_j h_t(x_j))}$$

$$\left(\sum_j D_t(j) = 1\right)$$

$$\left(\alpha_t = \frac{1}{2} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right)\right)$$

④ go to ①.

remark: $D_0(j) = \frac{1}{N} \quad | \quad j=1 \dots N$

$$\left(\epsilon_t = \sum_{j=1}^N D_t(j) \mathbb{1}(h_t(x_j) \neq y_j)\right)$$

$$Z(t) = \prod_{j=1}^N \prod_{t=1}^T \exp(-\alpha_t y_j h_t(x_j))$$

$$D_j(t) = \frac{\prod_{t=1}^{T-1} g_j^{(t)}}{\sum_{j=1}^N \prod_{t=1}^{T-1} g_j^{(t)}}$$

mathematical induction

training error

$$\frac{1}{N} \sum_{j=1}^N \mathbb{1}[H(x_j) \neq y_j] \leq \frac{1}{N} \sum_{j=1}^N \exp(-y_j \sum_{t=1}^T \alpha_t h_t(x_j)) = \prod_{t=1}^T Z(t)$$

$$\leq \prod_{t=1}^T 2\epsilon_t(1-\epsilon_t) \leq \exp(-2 \sum_{t=1}^T \epsilon_t^*) \quad \left(\epsilon_t = \frac{1}{2} - \epsilon_t^*\right) \quad \textcircled{3}$$

There are one ways to explain Boosting using χ^2 divergence
 statistical

Decision trees

IG(x) = $H(Y) - H(Y|X)$

↑ before branching ↑ after branching

overfit

↓ variance! (no bias)
 ↓ pruning

feature selected

$$\sum_i -P(Y_i) \log(P(Y_i))$$

$$\sum_i \frac{|S(X_i)|}{|S|} \sum_j -P(Y_j | X_i) \log(P(Y_j | X_i))$$

(Not MLE)

Non-Parametric models
 instance-based learning

(nearest neighbors)
 ⇒ curse of dimensionality (must store and retrieve all data)

1-NN ~ distance metrics $D(x, x^*) = \sqrt{(x-x^*)^T \Sigma (x-x^*)}$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_n^2 \end{bmatrix}$$

k-NN ~ classification ⇒ majority vote

Weighted k-NN : $\hat{y} = \frac{\sum_i \pi_i y_i}{\sum_i \pi_i}$

kernel: $\pi(x, x^*) = \exp\left\{-\frac{\sum_i (x_i - x_i^*)^2}{\rho^2}\right\}$

get ρ with gradient descent

locally weighted regression $\hat{w} = ((\pi X)^T \pi X)^{-1} (\pi X)^T \pi Y$

~~Application Controller~~ ← ~~Action Controller~~ is Base.

- returning User.

- ~~Compute Cn~~

$$\frac{1250}{10} = 125$$

online learning ⇒ perceptron (Not MLE) (report averaged weights at the end)
 (classification) $\hat{y} = \text{sign}(w \cdot x)$ $\frac{1250 \times 3}{10} = 375$ → ~~375~~ → 37.5

if $\hat{y} = y$ $w^{t+1} = w^t$

otherwise $w^{t+1} = w^t + \Delta$ ⇒ ? Δ in perceptron: $\Delta = y^i x^i$
 $= \frac{\partial l(w)}{\partial w}$

linearly separable case:

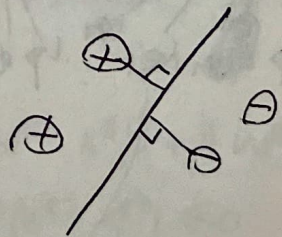
$$y = +1 \text{ or } -1$$

$$[(x^1, y^1), \dots, (x^T, y^T)]$$

$$\begin{cases} y^i = +1 & w^* x^i \geq \gamma \\ y^i = -1 & w^* x^i \leq -\gamma \end{cases} \Rightarrow \forall i \quad y^i (w^* x^i) \geq \gamma$$

$$l = \begin{cases} 0 & \text{if } y(w \cdot x) > 0 \\ -y(w \cdot x) & \text{if } y(w \cdot x) \leq 0 \end{cases} \text{ hinge loss}$$

$$\frac{\partial l}{\partial w} = \mathbb{1}(y(w \cdot x) \leq 0) (-y \cdot x)$$



Batch hinge minimization:

$$\Delta = \eta \frac{1}{N} \sum_{i=1}^N \{ \mathbb{1}[y^{(i)} (w^{(t)} \cdot x^{(i)}) \leq 0] y^{(i)} x^{(i)} \}$$

perceptron is the hinge loss minimization using SGD and $\eta = 1$

anglil@conquest

13au_final.pdf

Sun, 7 Dec 2014 15:29:51 -0800

ps441 / RICOH Aficio MP C305

make a mistake when $ywx \leq 0$

$$\text{loss} = \begin{cases} 0 & \text{if } ywx > 0 \\ -ywx & \text{if } ywx < 0 \end{cases}$$



$$\frac{\partial \text{loss}}{\partial w} = \mathbb{1}(ywx \leq 0) (-yx)$$

(hinge loss minimization using SGD) $(\eta=1)$

① Perceptron: (not MLE).

if $\hat{y} = y_t$ $w^{t+1} \leftarrow w^t$

else $w^{t+1} \leftarrow w^t + y^t x^t$ (mistake when $y^t w^t x^t \leq 0$).

online learning: report average $\hat{w} = \frac{1}{T+1} \sum_{t=0}^T w^t$

(1) linearly separable ($\|w^*\| \triangleq 1$).

$$\left. \begin{aligned} \forall t \text{ if } y^t = +1 \quad w^* x^t \geq \gamma \\ y^t = -1 \quad w^* x^t \leq -\gamma \end{aligned} \right\} \Rightarrow \forall t \quad y^t w^* x^t \geq \gamma$$

mistake at time t: $w^{(t+1)} = w^{(t)} + y^{(t)} x^{(t)}$

△ taking dot product with w^* : $w^* w^{(t+1)} = w^* w^{(t)} + y^{(t)} w^* x^{(t)}$

$$\therefore w^* w^{t+1} - w^* w^t \geq \gamma$$

$$\therefore w^* w^{t+1} - w^* w^0 \geq m\gamma \quad (\text{after } m \text{ mistakes}) \quad (m \leq t+1)$$

$$\therefore w^* w^{t+1} \geq m\gamma$$

△ taking square: $\|w^{(t+1)}\|^2 = \|w^{(t)}\|^2 + 2y^{(t)}(w^{(t)} \cdot x^{(t)}) + \|w^{(t)}\|^2$

$\therefore y^{(t)}(w^{(t)} \cdot x^{(t)})$ is mistaken, $\therefore 2y^{(t)}(x^{(t)} \cdot x^{(t)}) \leq 0$.

$$\therefore \|w^{(t)}\|^2 \leq R^2 \quad \therefore \|w^{t+1}\|^2 - \|w^t\|^2 \leq R^2 \quad \therefore \|w^{t+1}\|^2 - \|w^0\|^2 \leq mR^2$$

△ putting things together

$$m\gamma \leq w^* w^{t+1} \leq \|w^*\| \|w^{t+1}\| \leq \sqrt{m} R$$

$$\therefore m \leq \left(\frac{R}{\gamma}\right)^2$$

(2). Not linearly separable.

Let $M^{(t)}$ be time steps up to t when mistakes are made:

$$W^{(t)} = \sum_{j \in M^{(t)}} \underline{y^j x^j} \quad \left(\sum \frac{\partial \text{loss}}{\partial w} \right)$$

\uparrow $(1 \times k)$ $(k \times 1)$
 \uparrow \uparrow
 $y^j \cdot x \cdot x^j$

$$\therefore \hat{y} = \text{sign}(W^t \cdot x) = \text{sign}\left(x \cdot \sum_{j \in M^t} y^j x^j\right) = \text{sign}\left(\sum_{j \in M^t} y^j \cdot x \cdot x^j\right)$$

$$\therefore \hat{y} = \text{sign}\left(\sum_{j \in M^t} y^j \phi(x) \cdot \phi(x^j)\right) \text{ when using high dimensional features}$$

polynomials of degree exactly d : $\phi(u) \cdot \phi(v) = (\vec{u} \cdot \vec{v})^d$

polynomials of degree up to d : $\phi(u) \cdot \phi(v) = (\vec{u} \cdot \vec{v} + 1)^d$

Gaussian kernel: $\phi(u) \cdot \phi(v) = \exp\left(-\frac{\|\vec{u} - \vec{v}\|^2}{2\sigma^2}\right)$

Sigmoid: $\phi(u) \cdot \phi(v) = \tanh(\eta \vec{u} \cdot \vec{v} + \theta)$

(3) Plus regularization \Rightarrow SVM

$$\Delta \begin{cases} Wx^+ = \pm 1 \\ Wx^- = \mp 1 \\ x^+ = x^- + 2 \end{cases} \frac{W}{\|W\|} \Rightarrow \gamma = \frac{1}{\|W\|} \Rightarrow \begin{matrix} \text{maximizing margin } \gamma \\ \Downarrow \\ \text{minimizing } \|W\| \end{matrix}$$

$$\Delta \text{loss} = \begin{cases} 0 & \text{if } ywx \geq 1 \\ 1 - ywx & \text{if } ywx < 1. \end{cases}$$

Δ trade-off margin violation with $\|W\| \Rightarrow$

$$\text{loss} = \|W\|^2 + C \sum_{j=1}^N (1 - y_j W x_j)$$

SGD for SVM: $\frac{\partial \text{loss}}{\partial w} = +2W^t + C \mathbb{1}(yW^t x \leq 1) (-y x)$

(4) neural network: SGD, forward propagation (predicting)
backward propagation (learning)

clustering. ③

(1) hard: $\boxed{k\text{-means}}$ \Rightarrow coordinate descent. \Rightarrow MLE !!
 { classify:
 recenter:

(2) soft: \boxed{EM} (density estimation). \Rightarrow MLE !!
 Coordinate descent; local optima

$$\begin{aligned} \theta = \underset{\theta}{\text{argmax}} \quad & \left[\log\text{-likelihood} \right] = \ell(\theta) = \log \prod_{j=1}^m P(x_j | \theta) \\ & = \sum_{j=1}^m \log \sum_z P(x_j, z | \theta) = \sum_{j=1}^m \log \sum_z P(z | x_j, \theta) P(x_j | \theta) \\ & = \sum_{j=1}^m \log \sum_z P(z | x_j, \hat{\theta}) \frac{P(x_j, z | \theta)}{P(z | x_j, \hat{\theta})} \\ & \geq \sum_{j=1}^m \sum_{z=1}^k P(z | x_j, \hat{\theta}) \log \frac{P(x_j, z | \theta)}{P(z | x_j, \hat{\theta})} \Rightarrow \text{KL-divergence} \end{aligned}$$

$\{ M$ step: fix $p(z | x_j, \hat{\theta})$, max θ .

$\{ E$ step: compute $p(z | x_j, \theta)$ using θ . using Bayesian rule.

③ PCA using SVD to scale up.

$$\begin{aligned} \text{error} &= \sum_{i=1}^N (x^i - \hat{x}^i)^2 = \sum_{i=1}^N \left[\left(\bar{x} + \sum_{j=1}^d z_j^i u_j \right) - \left(\bar{x} + \sum_{j=1}^k z_j^i u_j \right) \right]^2 \\ &= \sum_{i=1}^N \left[\sum_{j=k+1}^d z_j^i u_j \right]^2 = \sum_{i=1}^N \sum_{j=k+1}^d (z_j^i)^2 \quad \left(\|u_j\|=1 \right) \\ &\quad \left(u_i \cdot u_j = 0 \text{ if } i \neq j \right) \\ &= \sum_{j=k+1}^d \sum_{i=1}^N [u_j (x^i - \bar{x})]^2 = \sum_{j=k+1}^d u_j^T \left[\sum_{i=1}^N (x^i - \bar{x})(x^i - \bar{x}) \right] u_j \\ &= N \sum_{j=k+1}^d u_j^T \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots \\ \sigma_{21} & \sigma_2^2 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} u_j \end{aligned}$$

$\left\{ \begin{array}{l} \text{where } \sigma_{rs} = \frac{1}{N} \sum_{i=1}^N (x_r^i - \bar{x}_r)(x_s^i - \bar{x}_s). \end{array} \right.$

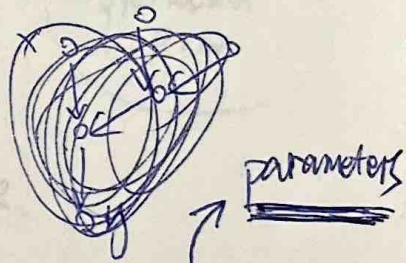
⑧

minimizing reconstruction error equals to picking orthonormal basis with smallest eigen values

1. recenter $X_c \leftarrow X - \bar{X}$
2. covar $\Sigma \leftarrow \frac{1}{N} X_c^T X_c$
3. eigen vectors & values of Σ .
4. principal components: k eigen vectors with highest eigen values

SVD: $X = WS V^T$

\swarrow weight/coordinate \searrow λ \searrow eigenvectors



(MLE!)

④ Bayesian Net

(Conditioned probability table = CPT)

e.g.

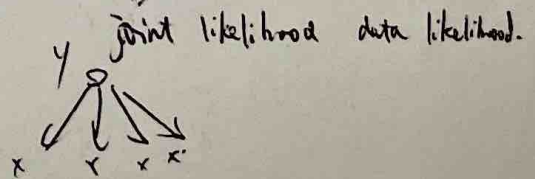
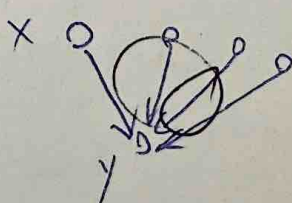
X_1	X_2	...	Y
0	0	---	1
0	1	---	2
1	0	---	0
1	1	---	4

m classes
~~if k binary features~~ ~~and k features~~
 then size of CPT = 2^k (row).
 need $\sum_i m(2^k - 1)$ params.
 (every node counts)

e.g. Naive bayes:

$$\sum_K m(2^k - 1) = km \text{ params.}$$

$$y = \text{arg max}_y P(y) \prod_{i=1}^k P(x_i | y) \quad (\text{MLE})$$



e.g. edgeless graph: least parameters \Rightarrow high bias.

e.g. fully-connected graph: most parameters \Rightarrow high variance. ⑨

Learning the structure of a Bayesian Net. \Rightarrow MLE!

structure + parameters $(P(X | \text{Parents}(X)))$

$$I(X_i, X_j) = \sum_{x_i, x_j} \hat{P}(x_i, x_j) \log \frac{\hat{P}(x_i, x_j)}{\hat{P}(x_i) \hat{P}(x_j)} \quad \text{where } \hat{P}(x_i, x_j) = \frac{\text{count}(x_i, x_j)}{m}$$

1. compute weight $I(X_i, X_j)$ of each possible edge (X_i, X_j)
2. find a maximum weight spanning tree (MST) : greatest total weight.
3. give directions to edges in MST
(pick any node as root, do breadth-first-search)

$$\max_{(X_i, X_j) \in E} \sum I(X_i, X_j)$$

greedy algo:

- {kruskal
- . Prim

NP-hard to learn structures with #parents > 1

(start from Chow-Liu tree,
local search: add/delete/insert edge.
score)

$\ln(N!)$
VC-dimension

decision tree of depth k . $\ln(N!) \approx 2^k \ln 2$

(2) Chernoff bounds

$$P\left(\left| \frac{1}{n} \sum x_i - \mu \right| > \epsilon\right) \leq e^{-n \epsilon^2}$$

$$P(\text{error} > \epsilon) \leq e^{-n \epsilon^2}$$

$$\int_0^\infty P(\text{error} > t) dt$$

$$P(\text{error} > \epsilon) \leq e^{-n \epsilon^2}$$

$$\forall \epsilon, P(\text{error} > \epsilon) \leq e^{-n \epsilon^2}$$

(1) PAC bound.

$$P(\text{error}_{\text{true}}(h) > \epsilon) \leq |H| e^{-N\epsilon} \leq \delta$$

H : hypothesis space.

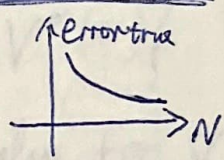
N : i.i.d samples.

h : any learned hypothesis that is consistent on the training data.

$$\epsilon = 0 < \epsilon < 1$$

union bound $\leq P(\exists h: \text{error}_{\text{train}}(h) = 0 \ \& \ \text{error}_{\text{true}}(h) \geq \epsilon) \leq k(1-\epsilon)^N$

(k hypothesis) $\leq |H| \leq |H| e^{-N\epsilon}$



total # of hypothesis

$|H|$ large
bias low
variance large

$|H|$ small
bias high
variance low

$$N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{\epsilon}$$

$$\epsilon \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{N} \Rightarrow O\left(\frac{1}{N}\right)$$

(more generally, decrease as $O\left(\frac{1}{\sqrt{N}}\right)$)

<list $x_1, \dots, x_m \sim y$ table, compute number of poss. of y >
boolean formula with m binary features:

$$\ln |H| = 2^m \ln 2$$

decision tree of depth k : $\ln |H| \leq 2^k \log m$

decision tree with up to k leaves: $\ln |H| \leq k \ln m + 2^k \ln k$

$\ln |H|$:
VC-dimension

(2) Chernoff bound:

$$P\left(\theta - \frac{1}{N} \sum_{j=1}^N x_j > \epsilon\right) \leq e^{-2N\epsilon^2}$$

$\theta = \text{error}_{\text{true}}(h)$ $\hat{\theta} = \text{error}_{\text{train}}(h) = \frac{1}{N} \sum_{j=1}^N \mathbb{1} \{h(x^j) \neq y^j\}$

$$= \int_{\mathcal{X}} p(x) \mathbb{1} \{h(x) \neq t(x)\} dx$$

(3) general: $P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq e^{-2N\epsilon^2}$

$$\forall h: P(\text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h) > \epsilon) \leq |H| e^{-2N\epsilon^2}$$

$$\Rightarrow N \geq \frac{\ln |H| + \ln \frac{1}{\delta}}{2\epsilon^2} \Rightarrow \epsilon \geq \sqrt{\frac{\ln |H| + \ln \frac{1}{\delta}}{2N}} \Rightarrow O\left(\frac{1}{\sqrt{N}}\right)$$

~~Neural Networks~~ MDPs.

x : state
 V : value, R : reward.
 π : policy = a : action

reinforcement learning: training by feedback.

$$V_{\pi}(x_0) = E_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R(x_t) \right] = E_{\pi} \left[R(x_0) + \gamma \sum_{t=1}^{\infty} \gamma^{t-1} R(x_t) \right]$$
$$= R(x_0) + \gamma \sum_{x_1} P(x_1 | x_0, a = \pi(x_0)) V_{\pi}(x_1).$$

stop when $\|V_{t+1} - V_t\|_{\infty} < \epsilon$. for computing the value of a policy.

compute the optimal value V : Bellman equation

$$V^*(x) = \max_a \left\{ \underline{R(x, a)} + \gamma \sum_{x'} P(x' | x, a) \underline{V^*(x')} \right\}$$

Value iteration on Bellman equation to solve for $V^*(x) \Rightarrow$ then we get $\pi^*(x)$ (dynamic programming).

$$\pi^*(x) = \arg \max_a \left\{ R(x, a) + \gamma \sum_{x'} P(x' | x, a) V^*(x') \right\}$$

start with some guess V^0 (for all states).

iteratively: $V^{t+1}(x) \leftarrow \max_a \left\{ R(x, a) + \gamma \sum_{x'} P(x' | x, a) V^t(x') \right\}$.

stop when $\|V^{t+1} - V^t\| < \epsilon$.

given or guessed.

reinforcement learning: rewards & trans. Prob. are unknown

goal: learn policy that (approx.) maximizes $E_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R(x_t) \right]$

interact with the world at each time step.

exploitation: stick with the known world.

exploration: at the risk of getting no big reward.

max algorithm: $\begin{cases} \text{initialize.} \\ R(x_0, a) = R_{\max} \\ P(x_0 | x_0, a) = 1. \end{cases}$ at each T , either exploit or explore

connectivity models:

hierarchical clustering
builds models based on
distance connectivity

centroid models:

K-means represents each
cluster by a single mean
vector

distribution models:

clusters are modeled using
statistical distributions
such as Gaussian E-M

no "correct" clustering.

~~K-means~~ centroid:

find the k cluster centers and
assign the objects to the nearest
cluster center, such that the
squared distances from the cluster
are minimized.

$$X_i \quad \vec{X} \quad k.$$

$$\min \sum_{i=1}^k \sum_{j=1}^{M_i} \|X_{ij} - C_i\|^2$$

NP-hard.

approximate \Rightarrow local optimum.
heuristics

problem: partition n observations
into k clusters

$$\arg \min_S \sum_{i=1}^k \sum_{X_i \in S_i} \|X_i - C_i\|^2$$

$$C_i = \frac{1}{|S_i|} \sum_{X_i \in S_i} X_i$$

each X_i is assigned to
exactly one S_i .

expectation

assignment step:

$$S_i^{(t)} = \{X_i : \|X_i - C_i^{(t)}\|^2 \leq \|X_i - C_l^{(t)}\|^2\}$$

Euclidean
distance.

maximization

update step:

$$C_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{X_i \in S_i^{(t)}} X_i$$

guarantee.

convergence: the assignments
no longer change.

hard part: initialization

the choice of k .

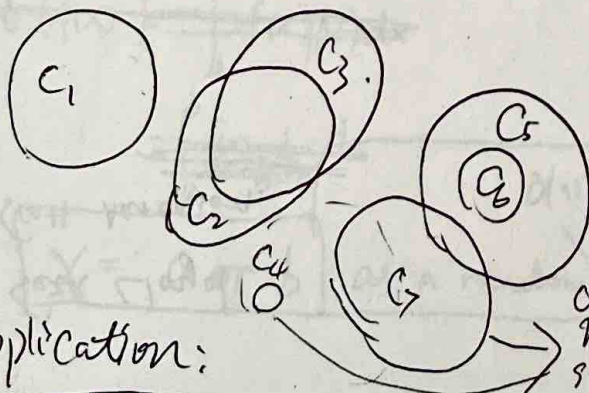
working scenarios
(limitations)

clusters are expected to be of similar size.

spherical clusters separable in a way such that MEAE converges to the cluster center

can't find non-convex clusters

failed example:



application:

Vector quantization \Rightarrow

Color quantization \Rightarrow color palette.

feature learning \Rightarrow

$\text{dist}(\vec{X}, \vec{C})$ as feature
 \Rightarrow named entity recognition

distribution models:

clusters defined as objects belonging most likely to the same distribution

overfitting
 \Downarrow

Choosing the appropriate model \rightarrow difficult

capture correlation and dependence between attributes

Naive Bayes classifier

target $f: X \rightarrow Y, X = X(a_1 \dots a_n)$

$$y_{\text{map}} = \underset{y_j \in Y}{\text{argmax}} P(y_j | a_1 \dots a_n)$$

$$= \underset{y_j \in Y}{\text{argmax}} \frac{P(a_1 \dots a_n | y_j) P(y_j)}{P(a_1 \dots a_n)}$$

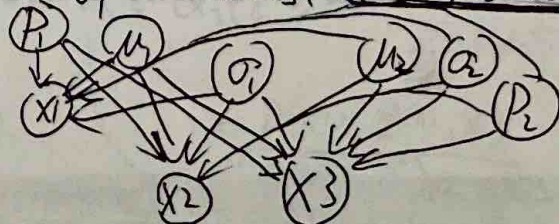
$$= \underset{y_j \in Y}{\text{argmax}} P(a_1 \dots a_n | y_j) P(y_j)$$

$$= \underset{y_j \in Y}{\text{argmax}} P(y_j) \prod_i P(a_i | y_j)$$

(Naive Bayes assumption: ~~a_i~~

a_j conditionally independent on one another)

Mixture of Gaussians: (Bayes nets)



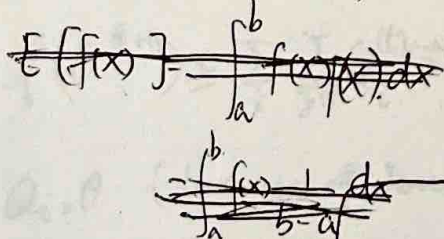
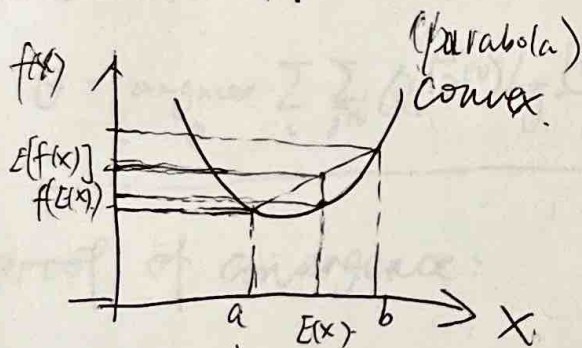
Jensen's inequality

f : convex: $f''(x) \geq 0$ for all $x \in \mathbb{R}$

$f(\vec{x})$: hessian $H \geq 0$ semi-definite

$$\boxed{E[f(x)] \geq f(E(x))}$$

if f is strictly convex
then $E[f(x)] = f(E(x))$ iff $P(X = E(x)) = 1$
($x = \text{constant}$)



$X = a \text{ or } b$ as a random variable

training data: $\vec{X} \in \mathbb{R}^{(m)}$, $i = 1, \dots, m$
 hidden/latent variables: \vec{z} (~~latent~~ cluster assigned for each i)
 known variables: $\vec{\theta}$ (~~cluster assignment~~ distribution param.)

$$l(\vec{\theta}) = \sum_{i=1}^m \log P(x^{(i)}; \vec{\theta})$$

$$= \sum_{i=1}^m \log \sum_{\vec{z}} P(x^{(i)}; \vec{z}; \vec{\theta})$$

goal = finding MLE of $\vec{\theta}$ (~~prob~~ cluster assignment)

$\vec{\theta}$ = discrete (3)

$$\vec{\theta} = \arg \max_{\vec{\theta}} l(\vec{\theta})$$

$$l(\vec{\theta}) = \sum_{i=1}^m \log P(x^{(i)}; \vec{\theta})$$

$$= \sum_{i=1}^m \log \sum_{\vec{z}^{(i)}} P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})$$

$$= \sum_{i=1}^m \log \sum_{\vec{z}^{(i)}} Q_i(\vec{z}^{(i)}) \frac{P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}{Q_i(\vec{z}^{(i)})}$$

where $Q_i(\vec{z}^{(i)})$ is some distribution over the $\vec{z}^{(i)}$'s.

$$\left(\sum_{\vec{z}^{(i)}} Q_i(\vec{z}^{(i)}) = 1 \right)$$

$$= \sum_{i=1}^m f[E(\vec{z}^{(i)})]$$

$$\geq \sum_{i=1}^m E[f(\vec{z}^{(i)})]$$

$$= \sum_{i=1}^m \sum_{\vec{z}^{(i)}} Q_i(\vec{z}^{(i)}) \log \frac{P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}{Q_i(\vec{z}^{(i)})}$$

(lower bound on $l(\vec{\theta})$)

for = to hold, we have

$$\frac{P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}{Q_i(\vec{z}^{(i)})} = \text{constant}$$

$$\therefore Q_i(\vec{z}^{(i)}) \propto P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})$$

$$\therefore \sum_{\vec{z}^{(i)}} Q_i(\vec{z}^{(i)}) = 1$$

$$\therefore Q_i(\vec{z}^{(i)}) = \frac{P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}{\sum_{\vec{z}^{(i)}} P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}$$

$$= \frac{P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}{\sum_{\vec{z}^{(i)}} P(x^{(i)}; \vec{z}^{(i)}; \vec{\theta})}$$

$$= \frac{P(x^{(i)}; \vec{\theta})}{P(\vec{z}^{(i)} | x^{(i)}; \vec{\theta})}$$

rewrite:

$$l(\vec{\theta}) \geq \sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \vec{\theta})}{Q_i(z^{(i)})}$$

lower bound: $Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \vec{\theta})$

$\therefore l(\vec{\theta})$ lower bounded by

$$\sum_{i=1}^m \sum_{z^{(i)}} P(z^{(i)} | x^{(i)}; \vec{\theta}) \log P(x^{(i)}; \vec{\theta})$$

maximization:

$$\vec{\theta} = \text{argmax}_{\vec{\theta}} \sum_i \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \vec{\theta})}{Q_i(z^{(i)})}$$

proof of convergence:

$$l(\vec{\theta}^{(t+1)}) \geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \vec{\theta}^{(t+1)})}{Q_i^{(t)}(z^{(i)})}$$

$\forall Q_i > 0$ (because of Jensen's inequality)

$$\geq \sum_i \sum_{z^{(i)}} Q_i^{(t)}(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \vec{\theta}^{(t)})}{Q_i^{(t)}(z^{(i)})}$$

(because of the maximization step)

$$= l(\vec{\theta}^{(t)})$$

$\therefore l(\vec{\theta})$ converges monotonically under EM

$z = 0 > 1$ in MoG example.

in Mixture of Gaussians

~~in M-step, we solve for ϕ, μ, Σ~~

~~$\vec{\theta} = (\phi, \mu, \Sigma)$~~

k -clusters

in E-step, we calculate

$$Q_i(z^{(i)}=j) = P(z^{(i)}=j | x^{(i)}; \vec{\theta})$$

$$= \frac{P(z^{(i)}=j; x^{(i)}; \vec{\theta})}{P(x^{(i)}; \vec{\theta})}$$

$$= \frac{P(x^{(i)}; \vec{\theta} | z^{(i)}=j) P(z^{(i)}=j)}{\sum_{j=1}^k P(x^{(i)}; \vec{\theta} | z^{(i)}=j) P(z^{(i)}=j)}$$

$$\triangleq w_j^{(i)}$$

in M-step, we solve for

$$\vec{\theta} = (\phi, \mu, \Sigma)$$

rewrite: lower bound.

$$l(\vec{\theta}) =$$

$$\sum_{i=1}^m \sum_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \vec{\theta})}{Q_i(z^{(i)})}$$

$$= \sum_{i=1}^m \sum_{j=1}^k Q_i(z^{(i)}=j) \log \frac{P(x^{(i)} | z^{(i)}=j; \vec{\theta}) P(z^{(i)}=j)}{Q_i(z^{(i)}=j)}$$

$$= \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \frac{1}{w_j^{(i)}} \frac{1}{P(z^{(i)}=j) \sum_j \frac{1}{P(z^{(i)}=j)}}$$

$$\exp[-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma_j^{-1} (x^{(i)} - \mu_j)] \phi_j$$

$$\begin{cases} \mu_j = \\ \phi_j = \\ \Sigma_j = \end{cases}$$

agglomerative: bottom-up. tree divisive: top-down.

$O(n^3) - O(n^2)$

merge by distance of two clusters:

$\|a-b\|_2 = \sqrt{\sum_i (a_i - b_i)^2}$

$\|a-b\|_1 = \sum_i |a_i - b_i|$

$\|a-b\|_\infty = \max_i |a_i - b_i|$

non-numeric: hamming distance

K-L divergence (e.g. decision tree)

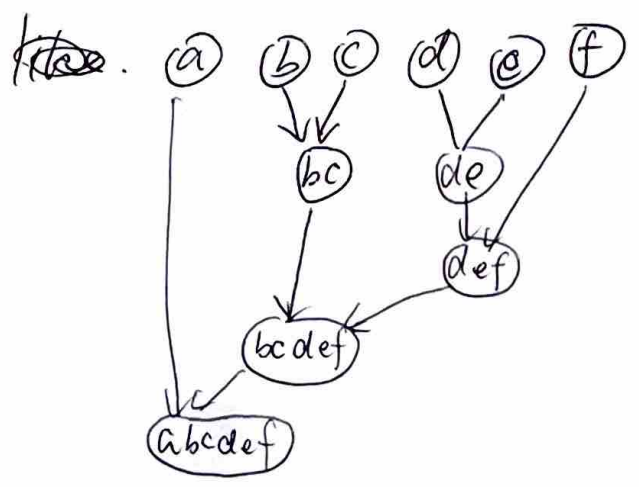
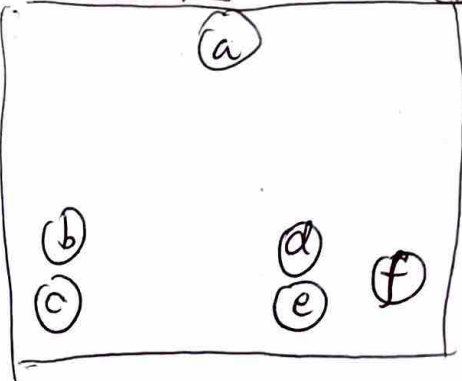
entropy

$O(2^n)$

min-cut.

k-means ..

etc.



evaluating clustering results