

Probability - L #4

DATE

Uncertainty

→ risks falsehood

→ leads to conclusions that are too weak for decision making.

→ leads to non-optimal decision making.

Probability

$\Omega \rightarrow$ sample space

$w \in \Omega$ is a sample point / possible world / atomic event

Prob. space

func $\Rightarrow P: \Omega \rightarrow \mathbb{R}$ such that $0 \leq P(w) \leq 1$

Random variable

$$\sum_{w \in \Omega} P(w) = 1$$

- A random variable (outcome of a random phenomenon) is a function from the sample space Ω to some range

$X: \Omega \rightarrow \mathcal{B} / \mathbb{R}$

odd: $\Omega \rightarrow \text{Boolean}$

X is a variable of a function

a) $X = x_i$: the random variable X has value $x_i \in \mathcal{B}$

b) $X = x_i = \{w \in \Omega | X(w) = x_i\}$

proposition

A prop. is the event (subset of Ω) where an assignment to a random variable holds

event $a = A = \text{true} = \{w \in \Omega | A(w) = \text{true}\}$

prop. can be combined using standard logic operators e.g.

$\neg A = A = \text{false} = \{w \in \Omega | A(w) = \text{false}\}$, etc.

Prior Prob.

↳ belief prior to arrival of any (new) evidence.

Joint Prob. Dist

DATE					
------	--	--	--	--	--

→ for a set of random variables gives the prob. of every atomic joint event on those random vars.

Conditional / Posterior Prob

$$P(\text{Cavity} = \text{true} \mid \text{weather} = \text{sunny}) \neq P(\text{Cavity} = \text{true}, \text{weather} = \text{sunny}) \\ \neq P(\text{Cavity} = \text{true})$$

Conditional Prob. Distribution

Rep. of all values of conditional prob. of random variables.

Total Probabilities

$$P(a) = P(a|b)P(b) + P(a|\neg b)P(\neg b)$$

in general for any random var Y,

$$P(X) = \sum_{y_i \in D(Y)} P(X \mid Y=y_i) (P(Y=y_i))$$

Int. by Enumeration ↗ set of values for variable Y.

- For any prop. ϕ , sum the atomic events where it is true.

$$P(\phi) = \sum_{w: w=\phi} P(w)$$

$$P(\text{cavity} \vee \text{toothache}) = P(\text{cavity...})$$

Normalization

$$P(A \mid B) = \alpha P(A, B)$$

Conditional Independence

X conditionally independent from Y given Z iff

$$P(X \mid Y, Z) = P(X \mid Z)$$

$$P(X, Y \mid Z) = P(X \mid Y, Z)P(Y \mid Z) = P(X \mid Z)P(Y \mid Z)$$

Y_i conditionally independent from Y_j given Z .

$$P(Y_1, \dots, Y_n | Z) = P(Y_1 | Z) P(Y_2 | Z) \dots P(Y_n | Z).$$

Bayes Rule

DATE

--	--	--	--	--	--	--	--

• Product rule $\Rightarrow P(a \cap b) = P(a|b) P(b) = P(b|a) P(a)$

$$P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

in dist. form,

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)} = \propto P(X|Y)P(Y)$$

$$\propto \frac{1}{P(X)}$$

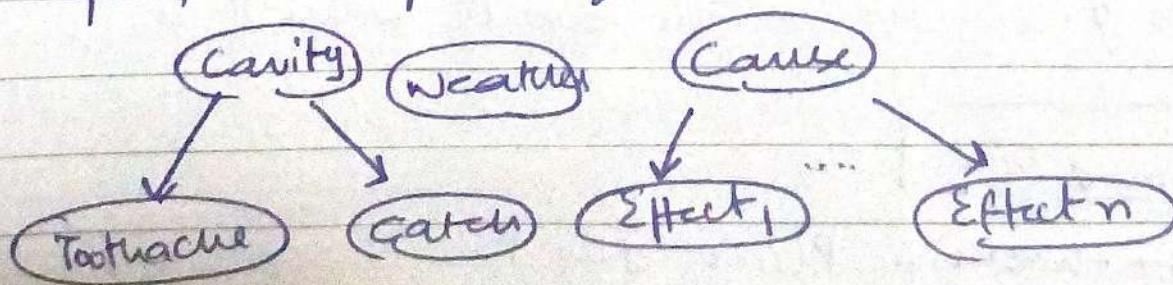
can be useful for assessing diagnostic prob. from causal prob.

$$P(\text{Cause}|\text{Effect}) = \frac{P(\text{Effect}|\text{Cause})P(\text{Cause})}{P(\text{Effect})}$$

Bayesian Networks

P(Effect)

↳ a simple graphical notation for conditional independence assertions and hence for compact specification of full joint dist.



• weather is independent and toothache and catch are conditionally dep. given cavity.

So the BN model would be:

$\Rightarrow P(X_i | \text{Parents}(X_i))$ for each variable X_i

$P(\text{weather})$,

$P(\text{Cavity})$, $P(\text{Toothache} | \text{Cavity})$

$P(\text{Catch} | \text{Cavity})$

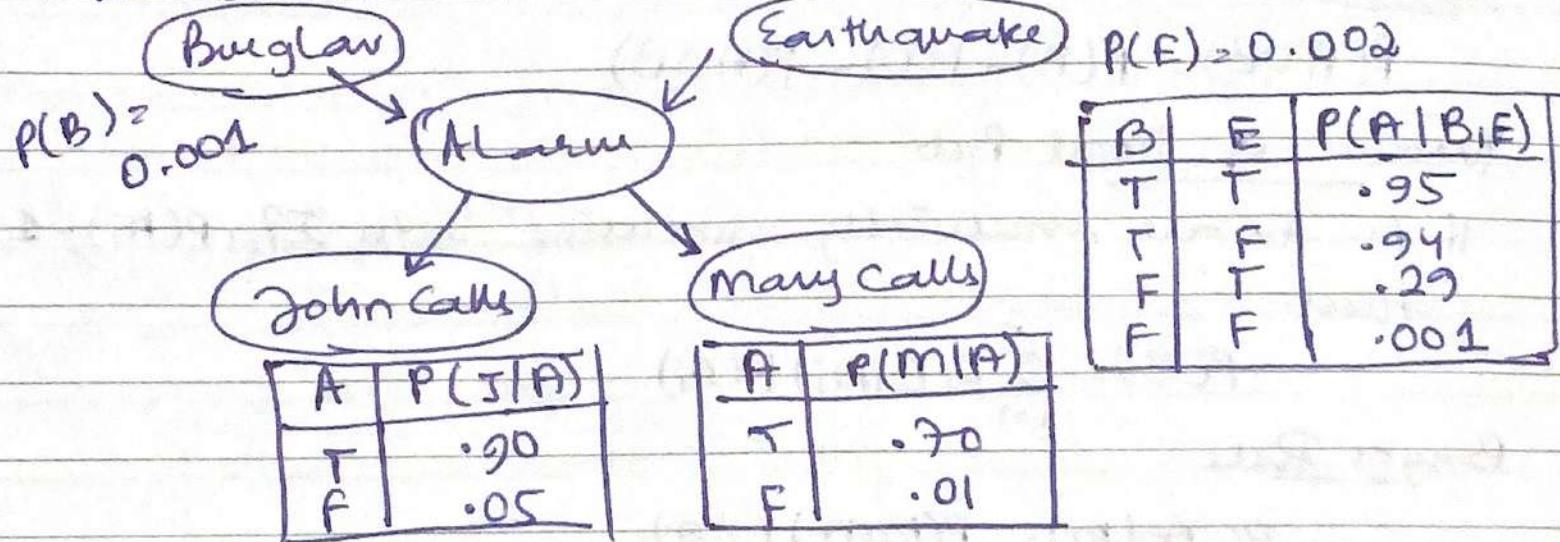
Burglar BN

Var: Burglar, Earthquake, John Calls, Mary Calls

- A burglar can set the alarm.
- An earthquake can set the alarm.
- The alarm can cause Mary to call.
- The alarm can cause John to call.

DATE

--	--	--	--	--



✓ A CPT for Boolean variable X_i with k Boolean parents has 2^k rows for rows of parent values.

- Each row requires one number p for $X_i = \text{true}$ (for $X_i = \text{false}, 1-p$)
- If each variable has no more than k parents, then network requires $O(n \cdot 2^k)$ numbers
- Grows linearly with n as compared to 2^n for full joint prob. dist.

Bayesian Learning - L#5

Basic Formulas

Product Rule

DATE						
------	--	--	--	--	--	--

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

Sum Rule

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Theorem of Total Prob

If A_1, A_2, \dots, A_n are mutually exclusive with $\sum_{i=1}^n P(A_i) = 1$

then

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i)$$

Bayes Rule

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

Classification as Probabilistic Estimation

Given a target function $f: X \rightarrow V$, dataset D and a new instance x' , best prediction $\hat{f}(x') = v^*$

$$v^* = \underset{v \in V}{\operatorname{argmax}} P(v|x', D)$$

Given D and x' , compute prob. distribution over V

$$P(v|x'; D)$$

Learning as Prob. Estimation

Given a dataset D and hypothesis space H

compute a prob. dist P_H over D . $P(H|D)$

Bayes rule applying

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

prior prob. of
hypothesis

prior prob. of
training data

MAP Hypotheses

Generally we want the most probable hypothesis h given D .

Maximum a posteriori hypothesis
(h_{MAP})

DATE

--	--	--	--	--	--

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D) = \operatorname{argmax}_{h \in H} \frac{P(D|h)P(h)}{P(D)}$$

$$\boxed{h_{MAP} = \operatorname{argmax}_{h \in H} P(D|h)P(h)} - ①$$

ML Hypothesis

If we assume $P(h_i) = P(h_j)$ we can further simplify eq ① and say

$$\boxed{h_{ML} = \operatorname{argmax}_{h \in D} P(D|h)}$$

Brute Force MAP Hypotheses Learner

- For each $h \in H$, calculate posterior prob.

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

- Output the hypothesis with highest posterior prob.

$$h_{MLP} = \operatorname{argmax}_{h \in H} P(h|D)$$

Bayes Optimal Classifier

Consider target function $f: X \rightarrow V$, $V = \{v_1, \dots, v_k\}$

dataset D and a new instance $x \in D$

$$\boxed{P(v_j|x, D) = \sum_{h \in H} P(v_j|x, h_i)P(h_i|D)}$$

Total prob. over H

$P(v_j|x, h_i)$: prob. that $h_i(x) = v_j$ is indep from D given h .

$$\boxed{P(v_j|x, h_i) = P(v_j|x, h_i, D)}$$

BOC

class of a new instance x

$$\Rightarrow \text{Vor}_B = \underset{v_j \in V}{\operatorname{argmax}} \sum_{h_i \in H} P(v_j | x, h_i) P(h_i | D)$$

DATE

--	--	--	--	--

No other method using same hypothesis space P & same prior knowledge can outperform this method on average.

Very powerful: Labelling new instances x with $\underset{v_j \in V}{\operatorname{argmax}} P(v_j | x, D)$ can correspond to none of hypotheses in H .

Example

$$P(h_i | \{d_i\}) = \alpha P(\{d_i\} | h_i) P(h_i)$$

$$P(h_i | \{d_1, d_2\}) = \alpha P(\{d_1, d_2\} | h_i) P(h_i)$$

$$\geq \alpha \underbrace{P(\{d_2\} | h_i)}_{\text{because independent data samples.}} P(\{d_1\} | h_i) P(h_i)$$

- because independent data samples.

General Approach

Given dataset $D = \{d_i\}$ with $d_i \in \{0, 1\}$, assuming a prob. distribution $P(d_i | \theta)$

$$\text{MLE} \Rightarrow \theta_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} \log P(d_i | \theta)$$

For Bernoulli dist,

$$P(X=k; \theta) = \theta^k (1-\theta)^{1-k}$$

$$\theta_{\text{ML}} = \frac{|\{d_i = 1\}|}{|D|}$$

Bernoulli Dist

→ Prob. dist of binary random variable $X \in \{0, 1\}$

$$P(X=1) = \theta$$

$$P(X=0) = 1-\theta$$

Multi-variate Bernoulli Dist

Joint prob. dist of a set of binary random variables
 X_1, \dots, X_n , each random var.
 following Bernoulli dist.

DATE						
------	--	--	--	--	--	--

$$P(X_1=k_1, \dots, X_n=k_n; \theta_1, \dots, \theta_n)$$

$$k_i \in \{0, 1\}$$

✓ Under the assumption that random variables X_i are mutually indep., multivariate Bernoulli dist is the product of n Bernoulli distributions.

$$P(X_1=k_1, \dots, X_n=k_n; \theta_1, \dots, \theta_n) = \prod_{i=1}^n P(X_i=k_i; \theta_i) = \prod_{i=1}^n \theta_i^{k_i} (1-\theta_i)^{1-k_i}$$

Binomial Dist

prob. dist of k outcomes from n Bernoulli-trials.

$$P(X=k; n, \theta) = \binom{n}{k} \theta^k (1-\theta)^{1-k}$$

Multinomial Dist

↳ Generalization of binomial dist for discrete valued random variables with d possible outcomes.

Prob. dist of k_1 outcomes for X_1, \dots, k_d outcomes for X_d after n trials with $\sum_{i=1}^d k_i = n$

$$P(X_1=k_1, \dots, X_d=k_d; n, \theta_1, \dots, \theta_d) = \frac{n!}{k_1! \dots k_d!} \theta_1^{k_1} \dots \theta_d^{k_d}$$

Naive Bayes Classifier

↳ uses conditional independence to approximate the solution.

X is conditionally indep. of Y given Z .

$$P(X, Y | Z) = P(X | Y, Z) P(Y | Z) \quad P(X | Z) P(Y | Z)$$

Assume target function

DATE

--	--	--	--	--	--	--

$f: X \rightarrow V$ where each instance x is described by attributes $\langle a_1, a_2, \dots, a_n \rangle$

compute

$$\operatorname{argmax}_{v_j \in V} P(v_j | x, D) = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, \dots, a_n, D)$$

→ No explicit representation of hypotheses.

Given a dataset D and a new instance $x = \langle a_1, \dots, a_n \rangle$ most prob. value of $f(x)$ is:

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} P(v_j | a_1, a_2, \dots, a_n, D)$$

$$= \operatorname{argmax}_{v_j \in V} \frac{P(a_1, a_2, \dots, a_n | D) P(v_j | D)}{P(a_1, a_2, \dots, a_n | D)}$$

$$V_{MAP} = \operatorname{argmax}_{v_j \in V} P(a_1, a_2, \dots, a_n | v_j, D) P(v_j | D)$$

Naive Bayes assumption

$$P(a_1, a_2, \dots, a_n | v_j, D) \propto \prod_i P(a_i | v_j, D)$$

class of new instance x

$$V_{NB} = \operatorname{argmax}_{v_j \in V} P(v_j | D) \prod_i P(a_i | v_j, D)$$

$$\hat{P}(v_j|D) = \frac{|\{<..., v_j>\}|}{|D|}$$

$$\hat{P}(a_i|v_j, D) = \frac{|\{<..., a_i, ..., v_j>\}|}{|\{<..., v_j>\}|}$$

DATE

Typical solution is Bayesian estimate with prior estimates

$$\hat{P}(a_i|v_j, D) = \frac{|\{<..., a_i, ..., v_j>\}| + mp}{|\{..., v_j\}| + m}$$

$p \rightarrow$ prior estimate for $P(a_i|v_j, D)$

$m \rightarrow$ weight given to the prior estimate.

Probabilistic Models for Classification - L#6

Given $f: x \rightarrow C, D = \{(x_i, c_i)\}_{i=1}^n$ and $x \notin D$ estimate

$$P(c_i|x, D)$$

$P(c_i|x) \rightarrow$ posterior

$P(x|c_i) \rightarrow$ class-conditional densities

Two families of models

Generative: estimate $P(x|c_i)$ and then compute $P(c_i|x)$

Discriminative: estimate $P(c_i|x)$ directly

→ Prob. Generative Models

Consider a case of two classes

Find the conditional prob.

$$\begin{aligned} P(c_i|x) &= P(x|c_i)P(c_i) = \frac{P(x|c_i)P(c_i)}{P(x|c_1)P(c_1) + P(x|c_2)P(c_2)} \\ &= \frac{1}{1 + e^{-a}} \quad a = \ln \frac{P(x|c_1)P(c_1)}{P(x|c_2)P(c_2)} \end{aligned}$$

Assume $p(x|C_i) = \mathcal{N}(x; \mu_i, \Sigma)$ - same covariance matrix

$$a = \frac{\ln P(x|c_1)P(c_1)}{P(x|c_2)P(c_2)}$$

$$= \frac{\ln \mathcal{N}(x; \mu_1; \Sigma) P(c_1)}{\mathcal{N}(x; \mu_2; \Sigma) P(c_2)} = w^T x + w_0$$

$$w = \bar{\Sigma}^{-1}(\mu_1 - \mu_2)$$

$$w_0 = -\frac{1}{2}\mu_1^T \bar{\Sigma}^{-1} \mu_1 + \frac{1}{2}\mu_2^T \bar{\Sigma}^{-1} \mu_2 + \ln \frac{P(c_1)}{P(c_2)}$$

Thus,

$$\boxed{P(c_1|x) = \sigma(w^T x + w_0)}$$

Multi-class

$$p(c_k|x) = \frac{p(x|c_k)P(c_k)}{\sum_j p(x|c_j)P(c_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

(Normalized exponential or softmax function)

$$a_k = \ln P(x|c_k)P(c_k)$$

MLE for 2 classes

Assume $P(c_1) = \pi$ ($P(c_2) = 1 - \pi$), $p(x|c_i) = \mathcal{N}(x; \mu_i, \Sigma)$

Given a dataset $D = \{(x_n, t_n)\}_{n=1}^N$, $t_n = 1$ if x_n

belongs to class C_1 , $t_n = 0$ if x_n belongs to class C_2 .

Let N_1 be the # of samples in D belonging to C_1 and N_2 be the # of samples of C_2 ($N_1 + N_2 = N$)

Likelihood

$$p(t|\pi, \mu_1, \mu_2, \Sigma, D) = \prod_{n=1}^N \left[\pi \mathcal{N}(x_n; \mu_1, \Sigma) \right]^{t_n} \left[(1-\pi) \mathcal{N}(x_n; \mu_2, \Sigma) \right]^{(1-t_n)}$$

Maximize log-likelihood

$$N \bar{x} = N_1 / N$$

$$\mu_1 = \frac{1}{N_1} \sum_{n=1}^N t_n x_n$$

$$\mu_2 = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) x_n$$

$$\Sigma = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$$

$$S_i = \frac{1}{N_i} \sum_{n \in C_i} (x_n - \mu_i)(x_n - \mu_i)^T, i=1,2$$

For k -classes

$$P(C_k) = \pi_k; P(x|C_k) = N(x, \mu_k, \Sigma)$$

Dataset $D = \{(x_n, t_n)\}_{n=1}^N$ with t_n 1-to- K encoding

$$\pi_k = \frac{N_k}{N}$$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} x_n$$

$$\Sigma = \sum_{k=1}^K \frac{N_k}{N} S_k, S_k = \frac{1}{N_k} \sum_{n=1}^N t_{nk} (x_n - \mu_k)(x_n - \mu_k)^T$$

Prob. Models \rightarrow represent posterior distributions

with parametric models

For two classes

$$P(C_1|x) = \delta(a)$$

For $k \geq 2$ classes

$$P(C_i|x) = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

$$a_k = w^T x + w_0$$

Compact Notation

$$w^T x + w_0 = (w^0 w) \begin{pmatrix} 1 \\ x \end{pmatrix}$$

$$\tilde{w} = \begin{pmatrix} w_0 \\ w \end{pmatrix}, \tilde{x} = \begin{pmatrix} 1 \\ x \end{pmatrix}$$

DATE

--	--	--	--	--

$$ak = w^T x + w_0 = \tilde{w}^T \tilde{x}$$

MLE for parametric Models

Likelihood $M\theta : P(t|\theta, D) = D = \langle x, t \rangle$

Solution $\theta^* = \underset{\theta}{\operatorname{argmax}} \ln P(t|\theta, x)$

When $M\theta$ belongs to exponential family

$$M\theta : P(t|\theta, D)$$

Likelihood $P(t|\theta, X)$ can be expressed in the form $P(t|\tilde{w}, X)$ with maximum likelihood

$$\tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \ln P(t|\tilde{w}, X)$$

Probabilistic Discriminative Models

Estimate directly

$$P(c_k | \tilde{x}, D) = \frac{\exp(ak)}{\sum_j \exp(a_j)}$$

$$\text{MLE } \tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \ln P(t|\tilde{w}, X)$$

$$\tilde{w}^* = \underset{\tilde{w}}{\operatorname{argmax}} \ln P(t|\tilde{w})$$

without estimating the model parameters.

Logistic Regression

↳ prob. discriminative model based on MLE.

DATE						
------	--	--	--	--	--	--

Two classes

Given dataset $D = \{(\tilde{x}_n | t_n)\}_{n=1}^N$ with $t_n \in \{0, 1\}$

Likelihood $f(t)$

$$p(t|\tilde{w}) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n}$$

$$\text{with } y_n = p(c_1 | \tilde{x}_n) = \sigma(\tilde{w}^\top \tilde{x}_n)$$

$t_n \rightarrow$ value in data corresponding to x_n

y_n : posterior pred. of current model \tilde{w} for x_n .

Cross-entropy loss function

$$E(\tilde{w}) = -\ln p(t|\tilde{w}) = -\sum_{n=1}^N [t_n \ln y_n + (1-t_n) \ln(1-y_n)]$$

Solve opt. problem

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

can be solved by \tilde{w} N-R / Gradient-Descent, etc.

Newton-Raphson,

$$\nabla E(\tilde{w}) = \sum_{n=1}^N (y_n - t_n) \tilde{x}_n$$

Gradient descent step,

$$\tilde{w} \leftarrow \tilde{w} - H(\tilde{w})^{-1} \nabla E(\tilde{w})$$

$H(\tilde{w}) = \nabla \nabla E(\tilde{w})$ is the Hessian matrix of $E(\tilde{w})$

Iterative Reweighted Least Squares

Given

$$\tilde{x} = \begin{pmatrix} \tilde{x}_1^\top \\ \vdots \\ \tilde{x}_n^\top \end{pmatrix} \quad t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$$

$y(\tilde{w}) = (y_1, \dots, y_n)^\top$ posterior predictions

$R(\tilde{w}) \Rightarrow$ diagonal matrix $R_{nn} = y_n(1-y_n)$

$$\nabla E(\tilde{w}) = \tilde{x}^T (y(\tilde{w}) - t)$$

$$H(\tilde{w}) = \nabla \nabla E(\tilde{w}) = \sum y_n (1-y_n) \tilde{x}_n \tilde{x}_n^T = \tilde{x}^T R(\tilde{w}) \tilde{x}$$

Iterative method

DATE

--	--	--	--	--	--

- initialize \tilde{w}

- Repeat until termination

$$\tilde{w} \leftarrow \tilde{w} - (\tilde{x}^T R(w) \tilde{x})^{-1} \tilde{x}^T (y(w) - t)$$

For multi-class,

$$P(C_k | \tilde{x}) = \frac{e^{a_k}}{\sum_j e^{a_j}} \quad a_k = \tilde{w}_k^T \tilde{x}$$

Model,

$$p(T | \tilde{w}_1, \dots, \tilde{w}_K) = \prod_{n=1}^N \prod_{k=1}^K P(C_k | x_n)^{t_{nk}} = \prod_{n=1}^N \prod_{k=1}^K y_n^{t_{nk}}$$

$$y_{nk} = Y[n, k]$$

$$t_{nk} = T[n, k]$$

Cross entropy

$$E(\tilde{w}_1, \dots, \tilde{w}_K) = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_{nk}$$

Kernels - L#9

Objects are rep. with fixed-length feature vectors $x \in \mathbb{R}^m$

what about variable length objs or infinite dims? $\phi(x)$

↳ images, time-series, strings, etc

- use a similarity measure b/w instances x, x' .

$$k(x, x') \rightarrow \text{kernel fn}$$

↳ similarity measure

$$\Rightarrow k(x, x') \geq 0$$

If $\phi(x)$, then we can use

$$k(x, x') = \phi(x)^T \phi(x')$$

Kernel function: a real valued function $k(x, x')$ where X is some abstract space for $x, x' \in X$.

Typically k is: symmetric

DATE

--	--	--	--	--	--	--	--

$$k(x, x') = k(x'; x)$$

non-negative

$$k(x, x') \geq 0$$

Note: Not strictly required.

- Input data D must be normalized in order for the kernel to be a good similarity measure

i) Min-max normalization: $\bar{x} = \frac{x - \min}{\max - \min}$

ii) Standardization: $\bar{x} = \frac{x - \mu}{\sigma}$

iii) Unit Vector: $\bar{x} = \frac{x}{\|x\|}$

Kernel families

Linear $\Rightarrow k(x, x') = x^T x'$

Poly $\Rightarrow k(x, x') = (\beta x^T x' + \gamma)^d$, $d \in \{2, 3, \dots\}$

RBF $\Rightarrow k(x, x') = \exp(-\beta \|x - x'\|^2)$

Sigmoid $\Rightarrow k(x, x') = \tanh(\beta x^T x' + \gamma)$

Consider a linear model $y(x; w) = w^T x$ with dataset

$$D = \{(x_n, t_n)\}_{n=1}^N$$

Minimize $J(w) = (t - Xw)^T(t - Xw) + \lambda \|w\|^2$

$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix}$ design matrix, $t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}$ output vector

Optimal solution

$$\hat{w} = (X^T X + \lambda I_N)^{-1} X^T t = X^T (X X^T + \lambda I_N)^{-1} t$$

with I_N the $N \times N$ identity matrix

DATE

Kernelized Linear Models

$$\alpha = (X X^T + \lambda I_N)^{-1} t$$

$$\hat{w} = X^T \alpha = \sum_{n=1}^N \alpha_n x_n x_n^T$$

If we consider a linear kernel $k(x, x') = x^T x'$, we can re-write the model as:

$$y(x; \hat{w}) = \sum_{n=1}^N \alpha_n k(x, x')$$

$$\alpha = (K + \lambda I_N)^{-1}$$

$K = X X^T \rightarrow$ Gram matrix

$$k(x_i, x_j) \left[\begin{array}{ccc} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{array} \right]$$

Linear model with kernel $k(x, x') = x^T x'$

$$y(x; \alpha) = \sum_{n=1}^N \alpha_n x_n^T x$$

$$\alpha = (K + \lambda I_N)^{-1} t$$

Gram matrix

$$K = \left[\begin{array}{ccc} x_1^T x_1 & \dots & x_1^T x_N \\ \vdots & \ddots & \vdots \\ x_N^T x_1 & \dots & x_N^T x_N \end{array} \right]$$

$$K = \left[\begin{array}{ccc} k(x_1, x_1) & \dots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \dots & k(x_N, x_N) \end{array} \right]$$

Kernel Trick

(kernel substitution)

- If input vector x appears in an algorithm only in the form of an product $x^T x'$, replace the inner product with some kernel $k(x, x')$
- Can be applied to any x (even infinite size)
- No need to know $\phi(x)$
- Directly extend many well-known algorithms

Kernelized SVM - classification

In SVM, the solution has the form

$$\hat{w} = \sum_{n=1}^N \alpha_n x_n$$

Linear model (linear kernel)

$$y(x, \alpha) = -\left(w_0 + \sum_{n=1}^N \alpha_n x_n^T x\right)$$

with Kernel trick,

$$y(x, \alpha) = -\left(w_0 + \sum_{n=1}^N \alpha_n k(x_n, x)\right)$$

Lagrangian Problem for Kernelized SVM

$$L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^M \alpha_n \alpha_m t_n t_m k(x_n, x_m)$$

$$\alpha_n = \dots$$

$$\rightarrow w_0 = \frac{1}{|SV|} \sum_{x_i \in SV} \left(t_i - \sum_{x_j \in S} \alpha_j t_j k(x_i, x_j) \right)$$

Kernelized Linear Regression

Linear model for regression $y = w^T x$ and data set
 $D = \{(x_n, t_n)\}_{n=1}^N$

DATE

--	--	--	--	--	--	--	--

Minimize the regularized loss function

$$J(w) = \sum_{n=1}^N E(y_n, t_n) + \lambda \|w\|^2$$

where $y_n = w^T x_n$

$$E(y_n, t_n) = (y_n - t_n)^2$$

$$\hat{w} = (X^T X + \lambda I_N)^{-1} X^T t = X^T \alpha$$

$$\alpha = (X^T X + \lambda I_N)^{-1} t$$

To make predictions

$$\Rightarrow g(x; \hat{w}) = \sum_{n=1}^N \alpha_n x_n^T x$$

$$\sum_{n=1}^N \alpha_n K(x, x_n)$$

$$\alpha = (K + \lambda I_N)^{-1} t$$

gram matrix

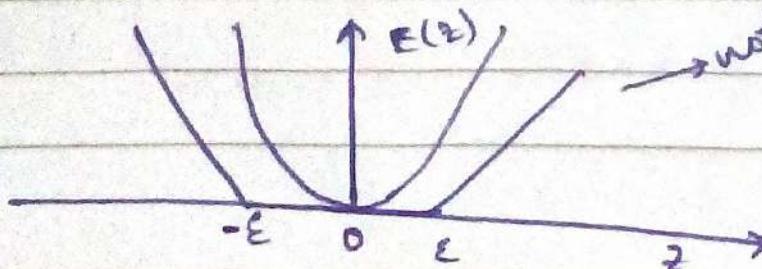
→ computation of K requires $|D|^2$ operations if K is not sparse.

Consider,

$$J(w) = C \sum_{n=1}^N E_\epsilon(y_n, t_n) + \frac{1}{2} \|w\|^2$$

with C inverse of λ and an ϵ -insensitive error function.

$$E_\epsilon(y, t) = \begin{cases} 0 & \text{if } |y-t| < \epsilon \\ |y-t| - \epsilon & \text{otherwise} \end{cases}$$



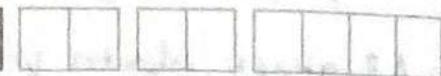
not differentiable
difficult to solve

Introduce slack variables $\xi_n^+, \xi_n^- \geq 0$

$$t_n \leq y_n + \varepsilon + \xi_n^+$$

$$t_n \geq y_n - \varepsilon - \xi_n^-$$

DATE



Points inside the ε -tube $y_n - \varepsilon \leq t_n \leq y_n + \varepsilon \Rightarrow \xi_{n+} = 0$

$$\xi_{n+}^+ > 0 \Rightarrow t_n > y_n + \varepsilon \quad \left. \right\} y_n = y(x_n; w)$$

$$\xi_n^- > 0 \Rightarrow t_n < y_n - \varepsilon$$

losses remain

$$J(w) = C \sum_{n=1}^N (\xi_{n+}^+ + \xi_n^-) + \frac{1}{2} \|w\|^2$$

subj. to constraints

$$t_n \leq y(x_n; w) + \varepsilon + \xi_{n+}^+$$

$$t_n \geq y(x_n; w) - \varepsilon - \xi_n^-$$

$$\xi_{n+}^+ > 0, \xi_n^- > 0$$

→ This is a standard quadratic program, can be easily solved.

Lagrangian Problem

$$\tilde{L}(a, a') = \dots \sum_{n=1}^N \sum_{m=1}^N a_n a_m \dots k(x_n, x_m) \dots$$

compute \hat{a}_n, \hat{a}_m (sparse) and

$$\hat{w} = t_n - \varepsilon - \sum_{m=1}^M (\hat{a}_m - \hat{a}'_m) k(x_n, x_m)$$

for some data point x_n such that $0 < a_n < C$

Pred.

$$y(x) = \sum_{n=1}^N (\hat{a}_n - \hat{a}'_n) K(x, x_n) + \hat{w}$$

From Karush-Kuhn-Tucker (KKT) condition,

Support vectors contribute to predict

$$\hat{a}_n > 0 \rightarrow \varepsilon + \xi_n + y_n - t_n = 0$$

↳ data point lies on or above the upper boundary of ε -tube

$a_n' > 0 \Rightarrow \epsilon + \epsilon_{n-1} - y_n + t = 0$
 ↳ data point lies on or below lower boundary of ϵ -tube.

→ All other data points inside the ϵ -tube have $\hat{a}_n = 0$, and $\hat{a}_n' = 0$, and thus do not contribute to pred.

~~Dragon Decision Document~~

Dimensionality Reduction - L #15

For data with structure*

- we expect fewer distortions than dimensions
 - Data live in a lower dimension manifold

Conclusion: Deal with high dimensional data by looking for lower dimensional embedding

Principal Component Analysis (PCA)

↳ widely used technique for various tasks such as:

- dimensionality reduction
 - data compression
 - data visualization
 - feature extraction

Given data $\{x_n\} \in \mathbb{R}^d$

Goal: Minimize data variance after projection
to some direction \mathbf{w}_1

Projected points: $U_1 x_n^T$

NOTE: $u_i^T u_i = 1$

$$\text{Mean value} \Rightarrow \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

Data centred matrix (\mathbf{x}) ($N \times D$):

$$X = \begin{bmatrix} (x_1 - \bar{x})^T \\ (x_2 - \bar{x})^T \\ \vdots \\ (x_N - \bar{x})^T \end{bmatrix}$$

Mean of projected points: $u_1^T \bar{x}$
 Variance of projected points: $\frac{1}{N} \sum_{n=1}^N [u_1^T x_n - u_1^T \bar{x}]^2 = u_1^T S u_1$
 with $S(D \times D)$ covariance matrix
 of the data.

DATE

--	--	--	--	--	--

$$S = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})(x_n - \bar{x})^T = \frac{1}{N} X^T X$$

Maximize the projected variance: $\max_{u_1} u_1^T S u_1$

subj. to constraint $u_1^T u_1 = 1$

Equivalent to unconstrained maximization
 with a Lagrange multiplier

$$\rightarrow \max_{u_1} u_1^T S u_1 + \lambda_1 (1 - u_1^T u_1)$$

Setting derivative w.r.t u_1 to zero we have

$$\rightarrow u_1^T S u_1 = \lambda_1 u_1^T u_1$$

u_1 must be an eigenvector of S

Multiply u_1^T

$$u_1^T S u_1 = \lambda_1 \underbrace{u_1^T u_1}_{=1}$$

$$\boxed{u_1^T S u_1 = \lambda_1} \rightarrow \text{variance after projection}$$

\rightarrow Variance is maximal when u_1 is the eigenvector corresponding to largest eigenvalue λ_1 .

This is called the first principal component.

\Rightarrow Repeat to find other directions which

- Maximize variance of projected data
- Are orthogonal to the previous directions

Steps: To perform PCA in an M -dimensional projection space with $M < D$

• compute \bar{x} : mean of data

• compute S : covar. mat of data

• compute M eigenvectors of S corresponding to M largest eigenvalues

PCA Error Minimization

Consider a complete orthonormal D-dimensional basis

$$u_i^T u_j = \delta_{ij}$$

with $\delta_{ij} \begin{cases} 0 & \text{otherwise} \\ 1 & \text{if } i=j \end{cases}$

DATE						
------	--	--	--	--	--	--

Each data point can be written as:

$$x_n = \sum_{i=1}^D a_{ni} u_i$$

Using the orthonormality prop. we have,

$$\alpha_{nj} = x_n^T u_j$$

$$x_n = \sum_{i=1}^D (\alpha_{ni} u_i) u_i$$

Goal: Approx. x_n using a lower dimensional rep.

$$\tilde{x}_n = \sum_{i=1}^M z_{ni} u_i + \sum_{i=M+1}^D b_i u_i$$

Evaluate approx error as:

$$J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$$

Minimizing w.r.t z_{nj} we get

$$z_{nj} = x_n^T u_j, j = 1, \dots, M$$

Minimizing w.r.t b_j we get

$$b_j = \tilde{x}_n^T u_j, j = M+1, \dots, D$$

Using these expressions we get,

$$\Rightarrow x_n - \tilde{x}_n = \sum_{i=M+1}^D [(x_n - \tilde{x}_n)^T u_i] u_i$$

The overall approx. error becomes

$$J = \frac{1}{N} \sum_{n=1}^N \sum_{i=M+1}^D (x_n^T u_i - \tilde{x}_n^T u_i) = \sum_{i=M+1}^D s_{ui}^2$$

minimize the approx. error subj. to $u_i^T u_i = 1$

$$\rightarrow J = \sum_{i=m+1}^D u_i^T S u_i + \lambda_i (1 - u_i^T u_i)$$

Get $d(u_i^T)' = 0$

DATE							
------	--	--	--	--	--	--	--

$S u_i = \lambda_i u_i$

u_i is an eigen vector of S with eigenvalue λ_i .

\rightarrow the error: $J = \sum_{i=M+1}^D \lambda_i$ \rightarrow This is minimized using u_i as the eigenvectors corresponding to $D-M$ smallest eigenvalues. choosing $D-M$ smallest eigenvalues of S corresponds to M highest eigenvalues of S in max. val. fnc.

PCA for high dimensionality ($N < D$)

At least $D-N+1$ eigenvalues of S are zero.

Ex: small-set of high-resolution images.

In this case finding eigenvalues of S ($D \times D$) matrix is inefficient.

For $N < D$:

Let X be the $N \times D$ centred data matrix

(i.e. $u - \bar{u}$ row is $(x_n - \bar{x})^T$) and corresponding covariance mat.

$$S = \frac{1}{N} X^T X$$

The corresponding eigenvectors e.g. are:

$$\frac{1}{N} X^T X u_i = \lambda_i u_i$$

Left-multip. by X ,

$$\frac{1}{N} XX^T(Xv_i) = \lambda_i(Xv_i)$$



$$v_i = Xv_i$$

$$\Rightarrow \frac{1}{N} XX^T v_i = \lambda_i v_i$$

- XX^T has same eigenvalues ($N-1$) as $X^T X$ (others are 0)
- XX^T is $N \times N$ matrix whose eigenvalues can be computed efficiently.

Given eigenvalues λ_i of XX^T , to find eigenvectors we left-multiply by X^T

$$\Rightarrow \left(\frac{1}{N} X^T X \right) (X^T v_i) = \lambda_i \underbrace{(X^T v_i)}_{\substack{\text{eigen} \\ \text{value}}} \text{ eigenvector}$$

$$u_i = \frac{1}{\sqrt{N\lambda_i}} X^T v_i$$

Summing up, $N < D$

- Consider the centered data matrix X
- Compute $N-1$ eigenvalues λ_i and eigenvectors v_i of XX^T
- $u_i = \frac{1}{\sqrt{N\lambda_i}} X^T v_i$

The so-obtained $N-1$ vectors u_i are eigenvectors of $S = X^T X$, with non-null eigenvalue λ_i .

Prob. PCA

- Rep. data x with lower dimensional latent variables z
 - Assume linear relationship DATE
- $$x = Wz + \mu$$

- Assume Gaussian dist of latent var.

$$P(z) = N(z; 0, I)$$

- Assume Linear-Gaussian relationship b/w latent vars & data.

$$P(x|z) = N(x; Wz + \mu, \sigma^2 I)$$

Marginal dist

$$P(x) = \int P(x|z) P(z) dz = N(x; \mu, C)$$

$$C = WW^T + \sigma^2 I$$

Posterior dist:

$$P(z|x) = N(z; M^{-1}W^{-T}(x - \mu), \sigma^2 M)$$

$$\rightarrow M = W W^T + \sigma^2 I$$

MLE: given data X

$$\Rightarrow \underset{W, M, \sigma^2}{\text{argmax}} \ln P(X|W, M, \sigma^2) = \sum_{n=1}^N \ln P(x_n|W, \sigma^2 M)$$

set derivatives to 0, we get.

$$\mu_{ML} = \bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$W_{ML} = \dots, \sigma^2_{ML} = \dots$$

ML \rightarrow PCA can be obtained by EM.

AEs for anomaly detect.

- Train an autoencoder (AE) with $\{x_n\}$ (learn a latent rep. for normal samples)
- compute final train loss
- Determine a threshold: $\delta = \text{mean}(\text{loss}) + \text{std}(\text{loss})$
- Given $x' \notin D$, reconstruct x' with AE and compute loss' .
- If $\text{loss}' < \delta$ return normal otherwise ~~or~~ abnormal.

→ Works better with ensemble of AEs.

VAEs → focus on learning latent space structure

GANs → focus on learning a distribution

CNNs

$$w_{out} = \frac{w_{in} - w_k + 2p}{s} + 1$$

$$h_{out} = \frac{h_{in} - h_k + 2p}{s} + 1$$

of trainable params of conv layer.

$$|\theta| = \underbrace{w_k \cdot h_k \cdot \text{ch} \cdot \text{dout}}_{\text{kernel weights}} + \underbrace{\text{dout}}_{\text{bias}}$$