

supervised learning:

given: training dataset: $\{\vec{x}_1, \dots, \vec{x}_N\}$

target dataset: $\{\vec{t}_1, \dots, \vec{t}_N\}$

learn: a function y st. $y(\vec{x}_i) \approx \vec{t}_i \forall i$

test: check y on new datapoints

ex: classification, regression

unsupervised learning

given: training data

ex: clustering, density estimation, visualization

Probability.

sum rule - $p(x) = \sum_y p(x,y)$ if discrete; $p(x) = \int_y p(x,y) dy$ if continuous

product rule - $p(x,y) = p(x|y)p(y) = p(y|x)p(x) = p(y,x)$

Bayes' rule - $p(y|x) = \frac{p(x|y)p(y)}{\sum_y p(x|y)p(y)}$

$p(y|x)$ = posterior probability

$p(y)$ = prior probability

given a probability distribution function $p(x)$, the expectation or average of a function $f(x)$ in $p(x)$ is defined as

$E[f(x)] = \sum_x p(x)f(x)$ if discrete; $E[f(x)] = \int p(x)f(x) dx$ if continuous

note: $E[f(x)] = \frac{1}{N} \lim_{N \rightarrow \infty} \sum_{n=1}^N f(x_n)$

$E_x[f(x,y)]$ = expectation of $f(x,y)$

wrt marginal dist of x

$E[f(x)|y] = \sum_x p(x|y)f(x)$

= conditional expectation on $f(x)$

the variance, a measure of how "variable" a random variable is, is defined as

$$\text{var}[f] = E[(f(x) - E[f(x)])^2] = E[f(x)^2] - E[f(x)]^2$$

the covariance of two rv's is

$$\text{cov}[x,y] = E_{x,y}[(x - E[x])(y - E[y])] = E_{x,y}[x,y] - E[x]E[y]$$

$\text{cov}[x,y] = 0$ if x and y are independent

Gaussian distributions

a Gaussian (or normal) r.v. has a probability density function

$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2\sigma^2}(x-\mu)^2\right\}$$

random variables $\vec{x} = (x_1, \dots, x_N)$ are iid (independent and identically distributed)

in this case, the joint density is $p(\vec{x}) = \prod_{n=1}^N p(x_n)$

maximum likelihood of normally distributed iid samples

$$\mu_{ml} = \underset{\mu}{\operatorname{argmax}} p(\vec{x}|\mu, \sigma^2) = \underset{\mu}{\operatorname{argmax}} \ln p(\vec{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

set the derivative to 0 to find the max

$$\frac{1}{2\sigma^2} \sum_{n=1}^N 2(x_n - \mu) = 0 \Rightarrow \sum_{n=1}^N (x_n - \mu) = 0 \Rightarrow \sum_{n=1}^N x_n = N\mu \Rightarrow \mu_{ml} = \frac{1}{N} \sum_{n=1}^N x_n = \text{sample mean}$$

$$\sigma_{ml}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ml})^2 = \text{sample variance}$$

bias - sample variance underestimates

$$E[\mu_{ml}] = \mu$$

$$E[\sigma_{ml}^2] = \left(\frac{N-1}{N}\right)\sigma^2 \Rightarrow \text{unbiased estimator for variance is } \tilde{\sigma}^2 = \frac{N}{N-1} \sigma_{ml}^2$$

ex: input $\vec{x} = (x_1, \dots, x_n)^\top$, target $\vec{t} = (t_1, \dots, t_n)^\top$

assume $t_n \sim N(y(x_n, \vec{w}), \beta^{-1})$, i.e. the targets approximate $y(\vec{x}, \vec{w})$ with a little bit of noise

MLE approach

$$\text{maximize } p(\vec{t} | \vec{x}, \vec{w}, \beta) = \prod_{n=1}^N N(t_n | y(x_n, \vec{w}), \beta^{-1})$$

$$\text{log likelihood: } -\frac{\beta}{2} \sum_{n=1}^N [y(x_n, \vec{w}) - t_n]^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

$$\text{MLE: } \arg \min_{\vec{w}} E(\vec{w}) = \arg \min_{\vec{w}} \frac{1}{2} \sum_{n=1}^N [y(x_n, \vec{w}) - t_n]^2$$

$$\text{precision (average error)} = \frac{1}{B_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N [y(x_n, \vec{w}_{\text{ML}}) - t_n]^2$$

Bayesian approach (MAP)

assume prior $p(\vec{w} | \alpha) = N(\vec{w} | \vec{O}, \alpha^{-1} I) = (\frac{\alpha}{2\pi})^{(m+1)/2} \exp(-\frac{\alpha}{2} \vec{w}^\top \vec{w})$ where α is a hyperparameter that acts like β

by Baye's theorem, $p(\vec{w} | \vec{x}, \vec{t}, \alpha, \beta) \propto p(\vec{t} | \vec{x}, \vec{w}, \beta) p(\vec{w} | \alpha)$

log of posterior $p(\vec{w} | \vec{x}, \vec{t}, \alpha, \beta) \Rightarrow \frac{\beta}{2} \sum_{n=1}^N [y(x_n, \vec{w}) - t_n]^2 + \frac{\alpha}{2} \vec{w}^\top \vec{w}$ = regularized sum of squares

multivariate Gaussian: $N(\vec{x} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^\top \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\}$ Mahalanobis distance

$\vec{\mu}$ = mean vector, Σ = DxD covariance matrix, $|\Sigma| = \det(\Sigma)$

$$\begin{aligned} E[\vec{x}] &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{\mu})^\top \Sigma^{-1} (\vec{x} - \vec{\mu}) \right\} \vec{x} d\vec{x} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \int \exp \left\{ -\frac{1}{2} \vec{z}^\top \Sigma^{-1} \vec{z} \right\} (\vec{z} + \vec{\mu}) d\vec{z} \\ &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \left[\vec{\mu} \int \exp \left\{ -\frac{1}{2} \vec{z}^\top \Sigma^{-1} \vec{z} \right\} d\vec{z} + \underbrace{\int \vec{z} \exp \left\{ -\frac{1}{2} \vec{z}^\top \Sigma^{-1} \vec{z} \right\} d\vec{z}}_0 \right] \\ &= \vec{\mu} = \text{sample mean} \end{aligned}$$

$$E[\vec{x} \vec{x}^\top] = \vec{\mu} \vec{\mu}^\top + \Sigma \Rightarrow \text{cov}[\vec{x}] = E[(\vec{x} - \vec{\mu})(\vec{x} - \vec{\mu})^\top] = E[\vec{x} \vec{x}^\top] - 2\vec{\mu} \vec{\mu}^\top + \vec{\mu} \vec{\mu}^\top = \vec{\mu} \vec{\mu}^\top + \Sigma - \vec{\mu} \vec{\mu}^\top = \Sigma = \text{sample covariance}$$

suppose $\vec{x} \sim N(\vec{x} | \vec{\mu}, \Sigma)$, $\vec{x} = \begin{bmatrix} \vec{x}_a \\ \vec{x}_b \end{bmatrix}$, $\vec{\mu} = \begin{bmatrix} \vec{\mu}_a \\ \vec{\mu}_b \end{bmatrix}$, $\Sigma = \begin{bmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{bmatrix}$, $\Lambda = \Sigma^{-1} = \begin{bmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{bmatrix}$

conditional Gaussian: $p(\vec{x}_a | \vec{x}_b) = N(\vec{x}_a | \vec{\mu}_{a|b}, \Lambda_{aa})$

$$\Lambda_{aa}^{-1} = \Sigma_{bb} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$$

$$\vec{\mu}_{a|b} = \vec{\mu}_a - \Lambda_{aa}^{-1} \Lambda_{ab} (\vec{x}_b - \vec{\mu}_b)$$

marginal Gaussian: $p(\vec{x}_a) = N(\vec{x}_a | \vec{\mu}_a, \Sigma_{aa})$

Dual Representation and Kernels

kernel function: $k(\vec{x}, \vec{x}') = \vec{\Phi}(\vec{x})^\top \vec{\Phi}(\vec{x}')$; $\vec{\Phi}(\vec{x})$ is a feature space mapping

Mercer's Theorem - a symmetric function $k(\vec{x}, \vec{y})$ is a kernel function iff the Gram matrix K is positive semidefinite for any collection of data $\{\vec{x}_n\}$

Linear Basis Function Models

linear regression - $y(\vec{x}, \vec{w}) = w_0 + w_1 x_1 + \dots + w_D x_D = \vec{w}^\top \vec{\Phi}(\vec{x}) \approx t$; where $\vec{\Phi}(\vec{x}) = \begin{bmatrix} 1 \\ \vec{x} \end{bmatrix}$

polynomial fitting - $y(\vec{x}, \vec{w}) = w_0 + w_1 x_1 + \dots + w_m x^m = \vec{w}^\top \vec{\Phi}(\vec{x}) \approx t$; where $\vec{\Phi}(\vec{x}) = \begin{bmatrix} 1 \\ \vec{x} \\ \vec{x}^2 \\ \vdots \\ \vec{x}^m \end{bmatrix}$

linear basis function model - $y(\vec{x}, \vec{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\vec{x}) = \vec{w}^\top \vec{\Phi}(\vec{x}) \approx t$; where $\vec{\Phi}(\vec{x}) = \begin{bmatrix} \phi_0(\vec{x}) \\ \phi_1(\vec{x}) \\ \vdots \\ \phi_{M-1}(\vec{x}) \end{bmatrix}$, $\phi_0(\vec{x}) = 1$

ex: let $t = y(\vec{x}, \vec{w}) + \epsilon$, $\epsilon \sim N(\epsilon | 0, \sigma^2)$

input: $X = \{\vec{x}_1, \dots, \vec{x}_N\}$, target: $\{t_1, \dots, t_N\}$, samples are iid

$$p(\vec{t} | X, \vec{w}, \sigma^2) = \prod_{n=1}^N N(t_n | \vec{w}^\top \vec{\Phi}(\vec{x}_n), \sigma^2) = \prod_{n=1}^N N(t_n | \vec{w}^\top \vec{\Phi}(\vec{x}_n), \sigma^2)$$

$$\text{maximum likelihood} - \ln p(\vec{t} | X, \vec{w}, \sigma^2) = -\frac{N}{2} \ln(2\pi) + \frac{N}{2} \ln \sigma^2 - \frac{\sigma^2}{2} \sum_{n=1}^N (t_n - \vec{w}^\top \vec{\Phi}(\vec{x}_n))^2$$

$$\vec{w}_{ML} = (\vec{\Phi}^\top \vec{\Phi})^{-1} \vec{\Phi}^\top \vec{t} = \vec{\Phi}^\dagger \vec{t}$$

$$\text{design matrix } \vec{\Phi} = \begin{bmatrix} \vec{\Phi}(\vec{x}_1) & \vec{\Phi}(\vec{x}_2) & \dots & \vec{\Phi}(\vec{x}_N) \\ \vec{\Phi}(\vec{x}_1) & \vec{\Phi}(\vec{x}_2) & \dots & \vec{\Phi}(\vec{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \vec{\Phi}(\vec{x}_1) & \vec{\Phi}(\vec{x}_2) & \dots & \vec{\Phi}(\vec{x}_N) \end{bmatrix} = \begin{bmatrix} \vec{\Phi}(\vec{x}_1) \\ \vec{\Phi}(\vec{x}_2) \\ \vdots \\ \vec{\Phi}(\vec{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times M}$$

$\vec{\Phi}^\dagger$ is the Moore-Penrose pseudoinverse of $\vec{\Phi}$

calculating the pseudoinverse is at least $O(\min\{M, N\}^3)$

$$\frac{1}{B_{ML}} = \frac{1}{N} \sum_{n=1}^N (t_n - \vec{w}_{ML}^\top \vec{\Phi}(\vec{x}_n))^2 = \text{sample variance}$$

$$w_0 = -\frac{\sigma^2}{2} \sum_{n=1}^N (t_n - \vec{w}^\top \vec{\Phi}(\vec{x}_n))^2 = -\frac{\sigma^2}{2} \sum_{n=1}^N (t_n - \sum_{j=1}^M w_j \vec{\Phi}_j(\vec{x}_n))^2 = \text{bias}$$

$$\text{maximized } w_0 = \bar{t} - \sum_{j=1}^M w_j \bar{\Phi}_j, \text{ where } \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n \text{ and } \bar{\Phi}_j = \frac{1}{N} \sum_{n=1}^N \vec{\Phi}_j(\vec{x}_n)$$

sequential learning - stochastic gradient descent

general setup - minimize $E(\vec{w}) = \sum_{n=1}^N E_n(\vec{w})$

$$\Rightarrow \vec{w}^{(k+1)} = \vec{w}^{(k)} - \eta \nabla E_n(\vec{w}^{(k)}) \text{ for a random } n, \text{ learning rate } \eta$$

least-mean-squares algorithm

$$\text{minimize } \sum_{n=1}^N \{t_n - \vec{w}^\top \vec{\Phi}(\vec{x}_n)\}^2 = \|\vec{t} - \vec{\Phi} \vec{w}\|^2$$

$$\Rightarrow \vec{w}^{(k+1)} = \vec{w}^{(k)} + \eta [t_n - \vec{w}^{(k)\top} \vec{\Phi}(\vec{x}_n)] \vec{\Phi}(\vec{x}_n)$$

Bayesian Linear Regression

likelihood - $p(\vec{t} | \vec{w}) = \prod_{n=1}^N N(t_n | \vec{w}^\top \vec{\Phi}(\vec{x}_n), \sigma^2)$

prior - $p(\vec{w}) = N(\vec{w} | \vec{m}_0, S_0)$

$$\begin{aligned} \text{posterior} - p(\vec{w} | \vec{t}) &= N(\vec{w} | (S_0^{-1} + \vec{\Phi}^\top \vec{\Phi})^{-1} (\vec{\Phi}^\top \vec{t} + S_0^{-1} \vec{m}_0), (S_0^{-1} + \vec{\Phi}^\top \vec{\Phi})^{-1}) \\ &= N(\vec{w} | \vec{m}_N, S_N) \end{aligned}$$

$$\text{where } \vec{m}_N = S_0^{-1} \vec{m}_0 + \vec{\Phi}^\top \vec{t}$$

$$S_N = (S_0^{-1} + \vec{\Phi}^\top \vec{\Phi})^{-1}$$

maximum posterior \vec{w} is the mode of $p(\vec{w} | \vec{t})$

$$\Rightarrow \vec{w}_{MAP} = \vec{m}_N$$

ex: $S_0 = \alpha^{-1} I$ is "infinitely broad"/uniform over the parameter space and offers no information

sequential method

- begin with prior $p(\vec{w})$

- update the posterior $p(\vec{w} | \vec{t})$ w/ data points

- use posterior $p(\vec{w} | \vec{t})$ as prior $p(\vec{w})$ for the next iteration, and repeat

special case: mean 0, isotropic Gaussian prior $p(\vec{w} | \alpha) = N(\vec{w} | \vec{0}, \alpha^{-1} I)$

then $p(\vec{w} | \vec{t}) = N(\vec{w} | \vec{m}_N, S_N)$

$$\text{where } \vec{m}_N = S_N (\alpha \vec{I} \vec{0} + \vec{\Phi}^\top \vec{t}) = \vec{\Phi} S_N \vec{\Phi}^\top \vec{t}$$

$$S_N = (\alpha \vec{I} + \vec{\Phi}^\top \vec{\Phi})^{-1}$$

$$\ln p(\vec{w} | \vec{t}) = -\frac{\sigma^2}{2} \sum_{n=1}^N (t_n - \vec{w}^\top \vec{\Phi}(\vec{x}_n))^2 - \frac{\alpha}{2} \vec{w}^\top \vec{w} + \text{const} = \text{regularized least-squares objective}$$

note: isotropic means cov is

in the form $\alpha^{-1} I$

predictive distribution: want to determine \vec{t} given a new input \vec{x}

$$p(\vec{t} | \vec{w}, \alpha, \beta) = \int p(\vec{t} | \vec{w}, \beta) p(\vec{w} | \vec{t}, \alpha, \beta) d\vec{w}$$

which is given by

$$p(\vec{t} | \vec{x}, \vec{t}, \alpha, \beta) = \mathcal{N}(\vec{t} | \vec{m}_n \vec{\Phi}(\vec{x}), \sigma_n^2(\vec{x}))$$

$$\sigma_n^2(\vec{x}) = \frac{1}{\beta} + \vec{\Phi}(\vec{x})^T S_n \vec{\Phi}(\vec{x})$$

Evidence Approximation

maximize evidence function $p(\vec{t} | \alpha, \beta) = \int p(\vec{t} | \vec{w}, \beta) p(\vec{w} | \alpha) d\vec{w}$

$$\text{where } p(\vec{t} | \vec{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \vec{w}^T \vec{\Phi}(\vec{x}_n), \beta^{-1})$$

$$p(\vec{w} | \alpha) = \mathcal{N}(\vec{w} | \vec{0}, \alpha^{-1} I)$$

$$\begin{aligned} p(\vec{t} | \alpha, \beta) &= \left[\left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{1}{2\pi} \right)^{M/2} |\alpha^{-1} I|^{-1/2} \exp \left\{ -\frac{\beta}{2} \sum_n ((t_n - \vec{w}^T \vec{\Phi}(\vec{x}_n))^2) \right\} \exp \left\{ -\frac{\alpha}{2} \vec{w}^T \vec{w} \right\} d\vec{w} \right] \\ &= \left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2} \int \exp \left\{ -\frac{\beta}{2} \| \vec{t} - \vec{\Phi} \vec{w} \|^2 - \frac{\alpha}{2} \vec{w}^T \vec{w} \right\} d\vec{w} \\ &= \left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2} \int \exp \left\{ -\frac{1}{2} (\vec{w} - \vec{m}_n)^T A (\vec{w} - \vec{m}_n) - \frac{\beta}{2} \| \vec{t} - \vec{\Phi} \vec{m}_n \|^2 - \frac{\alpha}{2} \vec{m}_n^T \vec{m}_n \right\} d\vec{w} \\ &= \left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2} \exp \left\{ -\frac{\beta}{2} \| \vec{t} - \vec{\Phi} \vec{m}_n \|^2 - \frac{\alpha}{2} \vec{m}_n^T \vec{m}_n \right\} \int \exp \left\{ -\frac{1}{2} (\vec{w} - \vec{m}_n)^T A (\vec{w} - \vec{m}_n) \right\} d\vec{w} \\ &= \left(\frac{\beta}{2\pi} \right)^{N/2} \left(\frac{\alpha}{2\pi} \right)^{M/2} \exp \left\{ -\frac{\beta}{2} \| \vec{t} - \vec{\Phi} \vec{m}_n \|^2 - \frac{\alpha}{2} \vec{m}_n^T \vec{m}_n \right\} (2\pi)^{M/2} |A|^{-1/2} \end{aligned}$$

$$\ln p(\vec{t} | \alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - \frac{N}{2} \ln (2\pi) - \frac{\beta}{2} \| \vec{t} - \vec{\Phi} \vec{m}_n \|^2 - \frac{\alpha}{2} \vec{m}_n^T \vec{m}_n - \frac{1}{2} \ln |A|$$

suppose $(\beta \vec{\Phi}^T \vec{\Phi})^{-1} = \lambda_i \vec{v}_i$, i.e. λ_i is an eigenvalue of $\beta \vec{\Phi}^T \vec{\Phi}$

$$\boxed{\alpha = \sum_{i=1}^M \frac{\lambda_i}{\vec{m}_n^T \vec{m}_n (\alpha + \lambda_i)}} \text{ maximizes } \ln p(\vec{t} | \alpha, \beta) \quad \frac{1}{B} = \frac{1}{N - \sum_{i=1}^M \frac{\lambda_i}{\alpha + \lambda_i}} \sum_{n=1}^N \left\{ t_n - \vec{m}_n^T \vec{\Phi}(\vec{x}_n) \right\}^2$$

procedure

1. choose α ,

2. find $\vec{m}_n = \beta S_n \vec{\Phi}^T \vec{t}$

3. estimate $\alpha_2 = \sum_{i=1}^M \frac{\lambda_i}{\vec{m}_n^T \vec{m}_n (\alpha + \lambda_i)}$

4. set $\alpha = \alpha_2$ and repeat

1. choose B ,

2. find $\vec{m}_n = \beta S_n \vec{\Phi}^T \vec{t}$

3. estimate $\frac{1}{B} = \frac{1}{N - \sum_{i=1}^M \frac{\lambda_i}{\alpha + \lambda_i}} \sum_{n=1}^N \left\{ t_n - \vec{m}_n^T \vec{\Phi}(\vec{x}_n) \right\}^2$

4. set $B_i = B_2$ and repeat

Discriminant Functions

goal: classify data, assuming classes are disjoint and therefore separable by decision boundaries

linear classification: data is linearly separable (decision surfaces are $(D-1)$ -dim hyperplanes)

setup: K-coding scheme for K classes (one-hot encoding)

$$\text{ex: } K=2 \Rightarrow \vec{t} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

discriminant function - simplest approach, assigns \vec{x} to a class

generalized linear model - $y(\vec{x}) = f(\vec{w}^T \vec{x} + w_0)$, $f()$ is a nonlinear function known as an activation function

K=2 case

$y(\vec{x}) = 0$ is the decision boundary

distance from origin to surface is determined by w_0

distance of point to surface is $\frac{|y(\vec{x})|}{\|\vec{w}\|}$

K>2 case

one-versus-one approach: use $K(K-1)/2$ binary classifiers

problem: ambiguous regions

$$y_k(\vec{x}) = \vec{w}_k^T \vec{x} + w_{k0}, \quad \vec{x} \in C_j \text{ where } j = \operatorname{argmax} y_i(x)$$

$y_i(x) = y_j(x)$ is the decision boundary between C_i and C_j

least squares: $\hat{y}(\vec{x}) = \tilde{\vec{w}}^T \vec{x}$

$$\tilde{\vec{w}} = \begin{bmatrix} \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_K \\ \vec{w}_1 & \vec{w}_2 & \dots & \vec{w}_K \end{bmatrix} \in \mathbb{R}^{(D+1) \times K}, \quad \vec{x} = \begin{bmatrix} 1 \\ \vec{x} \end{bmatrix}, \quad T = \begin{bmatrix} -\vec{t}_1^T \\ -\vec{t}_2^T \\ \vdots \\ -\vec{t}_K^T \end{bmatrix}$$

$$\vec{x} \in C_j \text{ where } j = \operatorname{argmax}_{\text{ith entry of } \hat{y}(\vec{x})} (\hat{y}(\vec{x}))_i$$

$$\text{error } E_o(\tilde{\vec{w}}) = \frac{1}{2} \operatorname{Tr} \{ (\vec{x} \tilde{\vec{w}} - T)^T (\vec{x} \tilde{\vec{w}} - T) \}$$

$$\text{closed form solution } \boxed{\tilde{\vec{w}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T T = \vec{X}^+ T}$$

problem: weak against outliers

Probabilistic Discriminative Models

$$\text{logistic sigmoid } \sigma(a) = \frac{1}{1 + \exp(-a)}$$

for $K=2$,

$$p(C_1 | \vec{x}) = \frac{p(\vec{x}|C_1)p(C_1)}{p(\vec{x}|C_1)p(C_1) + p(\vec{x}|C_0)p(C_0)} = \frac{1}{1 + \exp(-a)} = \sigma(a) \text{ where } a = \ln \frac{p(\vec{x}|C_1)p(C_1)}{p(\vec{x}|C_0)p(C_0)}$$

for $K>2$

$$p(C_k | \vec{x}) = \frac{p(\vec{x}|C_k)p(C_k)}{\sum_j p(\vec{x}|C_j)p(C_j)} = \underbrace{\frac{\exp(a_k)}{\sum_j \exp(a_j)}}_{\text{normalized exponential aka softmax}}$$

logistic regression

$$\text{posterior distribution } p(C_i | \vec{\Phi}) = y_i(\vec{\Phi}) = \sigma(\vec{w}^T \vec{\Phi})$$

$$\text{note: } \frac{d}{da} \sigma(a) = \sigma(a)(1 - \sigma(a))$$

$$\begin{aligned} p(\vec{t} | \vec{w}) &= \prod_{n=1}^N p(t_n | \vec{w}) \\ &= \prod_{n=1}^N (1 - p(C_n | \vec{\Phi}_n))^{1-t_n} p(C_n | \vec{\Phi}_n)^{t_n} \end{aligned}$$

cross-entropy error function

$$\begin{aligned} E(\vec{w}) &= -\ln p(\vec{t} | \vec{w}) = -\sum_{n=1}^N \{ t_n \ln p(C_n | \vec{\Phi}_n) + (1-t_n) \ln (1 - p(C_n | \vec{\Phi}_n)) \} \\ &= -\sum_{n=1}^N \{ t_n \ln \sigma(\vec{w}^T \vec{\Phi}_n) + (1-t_n) \ln (1 - \sigma(\vec{w}^T \vec{\Phi}_n)) \} \end{aligned}$$

$$\begin{aligned} \nabla E(\vec{w}) &= -\sum_{n=1}^N \left\{ t_n \frac{1}{\sigma(\vec{w}^T \vec{\Phi}_n)} \sigma(\vec{w}^T \vec{\Phi}_n) (1 - \sigma(\vec{w}^T \vec{\Phi}_n)) \vec{\Phi}_n + (1-t_n) \frac{1}{1-\sigma(\vec{w}^T \vec{\Phi}_n)} [-\sigma(\vec{w}^T \vec{\Phi}_n) (1 - \sigma(\vec{w}^T \vec{\Phi}_n))] \vec{\Phi}_n \right\} \\ &= -\sum_{n=1}^N \left[(t_n - \sigma(\vec{w}^T \vec{\Phi}_n)) \vec{\Phi}_n - (\sigma(\vec{w}^T \vec{\Phi}_n) - t_n) \sigma(\vec{w}^T \vec{\Phi}_n) \vec{\Phi}_n \right] \\ &= -\sum_{n=1}^N (t_n - \sigma(\vec{w}^T \vec{\Phi}_n)) \vec{\Phi}_n \\ &= \sum_{n=1}^N (\sigma(\vec{w}^T \vec{\Phi}_n) - t_n) \vec{\Phi}_n \end{aligned}$$

use iterative algorithm

Maximum Margin Classifier

support vector machines - choose \vec{w} and b to maximize the margin (dist b/w decision boundary and points)

note: distance of \vec{x}_n to decision surface is $\frac{t_n y(\vec{x}_n)}{\|\vec{w}\|} = \frac{t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b)}{\|\vec{w}\|}$

maximum margin solution: $\underset{\vec{w}, b}{\operatorname{argmax}} \left\{ \frac{1}{\|\vec{w}\|} \min_n [t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b)] \right\}$

note: rescaling \vec{w} doesn't change the distance

WLOG, let margin=1, i.e. $t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b) \geq 1$ for $n=1, \dots, N$

$t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b) = 1$ if \vec{x}_n is closest to decision surface

maximum margin solution: $\underset{\vec{w}, b}{\operatorname{argmin}} \frac{1}{2} \|\vec{w}\|^2 \text{ st. } t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b) \geq 1 \text{ for } n=1, \dots, N$

use Lagrange multipliers

$$L(\vec{w}, b, \alpha) = \frac{1}{2} \|\vec{w}\|^2 - \sum_{n=1}^N \alpha_n \{ t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b) - 1 \}$$

$$\frac{\partial}{\partial \vec{w}} L(\vec{w}, b, \alpha) = 0 \Rightarrow \vec{w} = \sum_{n=1}^N \alpha_n t_n \vec{\Phi}(\vec{x}_n)$$

$$\frac{\partial}{\partial b} L(\vec{w}, b, \alpha) = 0 \Rightarrow 0 = \sum_{n=1}^N \alpha_n t_n$$

$$\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m k(\vec{x}_n, \vec{x}_m)$$

where $\alpha_n \geq 0$ for $n=1, \dots, N$ and $\sum_{n=1}^N \alpha_n t_n = 0$

solve \tilde{L} for α then recover \vec{w} and b

predict with $y(\vec{x}) = \sum_{n=1}^N \alpha_n t_n k(\vec{x}, \vec{x}_n) + b$

note: the Karush-Kuhn-Tucker (KKT) conditions for \tilde{L} are

$$\begin{aligned} \alpha_n &\geq 0 \\ t_n y(\vec{x}_n) - 1 &\geq 0 \\ \alpha_n \{ t_n y(\vec{x}_n) - 1 \} &= 0 \end{aligned} \quad \Rightarrow \text{either } \alpha_n = 0 \text{ or } t_n y(\vec{x}_n) = 1$$

any point where $\alpha_n = 0$ doesn't contribute to predictions

points where $t_n y(\vec{x}_n) = 1$ are support vectors, i.e. they lie in the maximum margin hyperplane

hard margin - data must be linearly separable

minimize error function $\sum_{n=1}^N E_\infty(y(\vec{x}_n) t_n - 1) + \lambda \|\vec{w}\|^2$

$$\text{where } E_\infty(z) = \begin{cases} 0 & \text{if } z \geq 0 \text{ and } \lambda > 0 \\ \infty & \text{otherwise} \end{cases}$$

soft margin - introduce slack variable $\xi_n \geq 0$ to allow some "incorrect" points

$\xi_n = 0$ for correctly classified data, $\xi_n = |t_n - y(\vec{x}_n)|$ for other points

minimize $C \sum_{n=1}^N \xi_n + \frac{1}{2} \|\vec{w}\|^2$ st. $t_n (\vec{w}^\top \vec{\Phi}(\vec{x}_n) + b) \geq 1 - \xi_n$

Perceptron

perceptron: a two class model of the form $y(\vec{x}) = f(\vec{w}^\top \vec{\Phi}(\vec{x}))$ where $f(a) = \begin{cases} +1 & a \geq 0 \\ -1 & a < 0 \end{cases}$

minimize perceptron criterion $E_p(\vec{w}) = -\sum_{n \in M} \vec{w}^\top \vec{\Phi}(\vec{x}_n) t_n$ where M is the set of misclassified points

SGD on E_p

1. initialize $\vec{w}^{(0)}$

2. for $t=1, \dots, K$ randomly sample n

if $\vec{w}^{(t-1)\top} \vec{\Phi}(\vec{x}_n) < 0$, update $\vec{w}^{(t)} = \vec{w}^{(t-1)} + \eta \vec{\Phi}(\vec{x}_n) t_n$

else, $\vec{w}^{(t)} = \vec{w}^{(t-1)}$

perceptron convergence theorem - if data is linearly separable, perceptron algorithm guaranteed to converge in finite # of steps

Feed-forward Neural Networks

suppose D input vars x_1, \dots, x_D ; let $x_0 = 1$

each layer consists of multiplication by weight matrix and a nonlinear activation function

$$\text{in layer 1, } a_j^{(1)} = \sum_{i=0}^D w_{ji}^{(1)} x_i = \vec{w}_j^{(1)} \vec{x} \quad \text{for } j=1, \dots, M$$

activations weights

use a nonlinear activation function $h()$, hidden units $z_j = h(a_j)$

$$\text{in layer 2, } a_2 = \sum_{j=0}^M w_{kj}^{(2)} z_j$$

output unit activations

overall 2-layer network function, sometimes known as multilayer perceptron

$$y(\vec{x}, \vec{w}) = \sigma\left(\sum_{j=0}^M w_{kj}^{(2)} h\left(\sum_{i=0}^D w_{ji}^{(1)} x_i\right)\right)$$

evaluating the function is forward propagation

note: neural networks are universal approximators

a 2-layer NN can uniformly approximate any cont. func on compact input data to arbitrary accuracy w/ a sufficient number of hidden units and a broad range of activation functions

note: weight-space symmetry guarantees that no NN is unique

ex: if $h()$ is $\tanh()$, and we negate all weights in layer 1, we can negate all the weights in layer 2

to produce the same output (since \tanh is odd)

$$\sum_{j=0}^M w_{kj}^{(2)} \tanh\left(\sum_{i=0}^D w_{ji}^{(1)} x_i\right) = -\sum_{j=0}^M w_{kj}^{(2)} \tanh\left(-\sum_{i=0}^D w_{ji}^{(1)} x_i\right) = \sum_{j=0}^M (-w_{kj}^{(2)}) (-\tanh\left(\sum_{i=0}^D w_{ji}^{(1)} x_i\right)) = \sum_{j=0}^M (-w_{kj}^{(2)}) \tanh\left(\sum_{i=0}^D (-w_{ji}^{(1)}) x_i\right)$$

training the model

assume: target has Gaussian distribution, i.e. $p(t|y(\vec{x}, \vec{w})) = N(t|y(\vec{x}, \vec{w}), \beta^{-1})$

training data is i.i.d., i.e. likelihood is $p(t|X, \vec{w}, \beta) = \prod_n p(t_n|y(\vec{x}, \vec{w}), \beta) = \prod_n N(t_n|y(\vec{x}, \vec{w}), \beta^{-1})$

goal: maximize the log likelihood

$$-\frac{\beta}{2} \sum_n \{y(\vec{x}, \vec{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

equivalent to minimizing the sum of squares error

$$E(\vec{w}) = \sum_n (y(\vec{x}_n, \vec{w}) - t_n)^2$$

use SGD to find a local minimum

note: $E(\vec{w})$ is not convex, so we might find a local minimum instead of the global minimum

weight-space symmetry \Rightarrow multiple local minima \Rightarrow train the model multiple times and select the best

note: $E(\vec{w}) = \sum_n E_n(\vec{w})$ where $E_n(\vec{w}) = (y(\vec{x}_n, \vec{w}) - t_n)^2$

update weights:

1) initialize $\vec{w}^{(0)}$, could be random

2) pick a random x_n

3) compute $\nabla E_n(\vec{w}^{(0)})$

$$4) \vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla E_n(\vec{w}^{(t)})$$

Back Propagation

1 layer: $E_n = \frac{1}{2} \sum_{k=1}^K (y_{nk} - t_{nk})^2 = \frac{1}{2} \sum_{k=1}^K (y_k(\vec{x}_n, \vec{w}) - t_{nk})^2 = \frac{1}{2} \sum_{k=1}^K \left(\sum_{i=0}^D w_{ki} x_{ni} - t_{nk} \right)^2$

 $\frac{\partial E_n}{\partial w_{ji}} = (y_{nj} - t_{nj}) x_{ni}$

2 layer: $E_n = \frac{1}{2} \sum_{k=1}^K (y_{nk} - t_{nk})^2 = \frac{1}{2} \sum_{k=1}^K (y_k(\vec{x}_n, \vec{w}) - t_{nk})^2 = \frac{1}{2} \sum_{k=1}^K \left(\sum_{i=0}^M h(a_i) - t_{nk} \right)^2$

 $\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_i} \frac{\partial a_i}{\partial w_{ji}}$
 $= \delta_i \frac{\partial a_i}{\partial w_{ji}}$
 $= \boxed{\delta_i Z_i}$

note: $a_i = \sum_{j=0}^D w_{ji} Z_j$

 $\frac{\partial a_i}{\partial w_{ji}} = Z_j$

compute δ_i :

for output units (last layer): $\delta_k = y_k - t_k$

for hidden units, use chain rule:

$$\begin{aligned}\delta_i &= \frac{\partial E_n}{\partial a_i} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_i} \quad \text{for all } a_k = f(a_i) \\ &= \sum_k \delta_k \vec{w}_{ki} h'(a_i) \\ &= h'(a_i) \sum_k \delta_k \vec{w}_{ki}\end{aligned}$$

$$\begin{aligned}a_k &= \sum_{j=0}^D \vec{w}_{kj} h(a_j) \\ \frac{\partial a_k}{\partial a_j} &= \vec{w}_{kj} h'(a_j)\end{aligned}$$

then update weights with $\vec{w}^{(t+1)} = \vec{w}^{(t)} - \eta \nabla E_n(\vec{w}^{(t)})$

Principal Component Analysis (PCA)

applications: dimension reduction

data visualization

feature extraction

topic modeling

etc

other names: Karhunen-Loeve theorem (signal processing)

proper orthogonal decomposition (mechanical engineering)

Eckart-Young theorem (linear algebra)

empirical orthogonal functions (meteorological science)

idea: orthogonally project data onto a lower dimensional subspace

maximize variance of projections, or
minimize average projection distance

same algorithm

setup: data $\{x_i\}_{i=1}^n \in \mathbb{R}^D$

goals: project data onto an M-dimensional subspace, $M < D$, called the principal subspace

assumptions: $M=1$ (generalize with induction)

data is centered, i.e. $\bar{x} = \frac{1}{N} \sum_{i=1}^N \bar{x}_i = \vec{0}$ where \bar{x} is the sample mean

1) write errors between points and their projections onto principal subspace (a line through the origin since $M=1$)

projection = $(\vec{v}_1^\top \vec{x}_n) \vec{v}_1$.

mean of projections is $\vec{0}$

the error between points and projections

$$\begin{aligned}\|\vec{x}_n - (\vec{v}_1^\top \vec{x}_n) \vec{v}_1\|^2 &= [\vec{x}_n - (\vec{v}_1^\top \vec{x}_n) \vec{v}_1]^\top [\vec{x}_n - (\vec{v}_1^\top \vec{x}_n) \vec{v}_1] \\ &= \vec{x}_n^\top \vec{x}_n - 2(\vec{v}_1^\top \vec{x}_n) \vec{v}_1^\top \vec{x}_n + (\vec{v}_1^\top \vec{x}_n)^2 \vec{v}_1^\top \vec{v}_1 \\ &= \|\vec{x}_n\|^2 - 2(\vec{v}_1^\top \vec{x}_n)^2 + (\vec{v}_1^\top \vec{x}_n)^2 \|\vec{v}_1\|^2 \\ &= \|\vec{x}_n\|^2 - (\vec{v}_1^\top \vec{x}_n)^2\end{aligned}$$

2) choose basis vector(s) to minimize the mean error

$$\text{mean error} = \frac{1}{N} \sum_{n=1}^N \| \vec{x}_n - (\vec{v}_i^\top \vec{x}_n) \vec{v}_i \|_2^2 = \frac{1}{N} \sum_{n=1}^N [\| \vec{x}_n \|_2^2 - (\vec{v}_i^\top \vec{x}_n)^2]$$

note: minimize mean error $\frac{1}{N} \sum_{n=1}^N [\| \vec{x}_n \|_2^2 - (\vec{v}_i^\top \vec{x}_n)^2]$

$$\Leftrightarrow \text{maximize } \frac{1}{N} \sum_{n=1}^N (\vec{v}_i^\top \vec{x}_n)^2 = \frac{1}{N} \sum_{n=1}^N (\vec{v}_i^\top \vec{x}_n - \underbrace{\vec{v}_i^\top \vec{x}_n}_{\text{"sample covariance matrix"}})^2$$

\Leftrightarrow maximize sample variance of projections

$$\text{note: } \frac{1}{N} \sum_{n=1}^N (\vec{v}_i^\top \vec{x}_n - \vec{v}_i^\top \bar{\vec{x}}_n)^2 = \vec{v}_i^\top \left[\frac{1}{N} \sum_{n=1}^N (\vec{x}_n - \bar{\vec{x}})(\vec{x}_n - \bar{\vec{x}})^\top \right] \vec{v}_i$$

"sample covariance matrix"

constrain $\| \vec{v}_i \|_2 = 1$ so objective is bounded

$$\text{let } f(\vec{v}_i) = \vec{v}_i^\top S \vec{v}_i \text{ and } g(\vec{v}_i) = \vec{v}_i^\top \vec{v}_i - 1$$

$$\nabla f(\vec{v}_i) = \lambda \nabla g(\vec{v}_i)$$

$$2S\vec{v}_i = 2\lambda \vec{v}_i$$

$$S\vec{v}_i = \lambda \vec{v}_i$$

to maximize, λ has to be the largest eigenvalue

PCA algorithm

1) form sample covariance matrix $S = \frac{1}{N} \sum_{n=1}^N (\vec{x}_n - \bar{\vec{x}})(\vec{x}_n - \bar{\vec{x}})^\top$

2) compute eigenvectors associated to the M largest eigenvalues of S , $\{ \vec{v}_i \}_{i=1}^M$

3) compute coordinates of the projection of the centered data $(\vec{x}_n - \bar{\vec{x}})$ onto the principal subspace
basis: $\{ \vec{v}_i \}_{i=1}^M$

coefficients: $\vec{v}_i^\top (\vec{x}_n - \bar{\vec{x}})$

the variance captured by the projection is $\sum_{i=1}^M \lambda_i$

the missed variance, or distortion, is $\sum_{i=M+1}^N \lambda_i$

Probabilistic PCA

let \vec{z} be a latent variable corresponding to the principal subspace

prior: $p(\vec{z}) = \mathcal{N}(\vec{z} | \vec{0}, I)$

conditional dist $p(\vec{x} | \vec{z}) = \mathcal{N}(\vec{x} | W\vec{z} + \vec{\mu}, \sigma^2 I)$

note: \vec{z} "generates" \vec{x} because $\vec{x} = W\vec{z} + \vec{\mu} + \vec{\epsilon}$ where $\vec{\epsilon} \sim \mathcal{N}(\vec{\epsilon} | \vec{0}, \sigma^2 I)$

find $p(\vec{x}) = \int p(\vec{x} | \vec{z}) p(\vec{z}) d\vec{z}$

$$\mathbb{E}[\vec{x}] = \mathbb{E}[W\vec{z} + \vec{\mu} + \vec{\epsilon}] = W\mathbb{E}[\vec{z}] + \vec{\mu} + \mathbb{E}[\vec{\epsilon}] = \vec{\mu}$$

$$\text{cov}[\vec{x}] = \mathbb{E}[(W\vec{z} + \vec{\epsilon})(W\vec{z} + \vec{\epsilon})^\top]$$

$$= W\mathbb{E}[\vec{z}\vec{z}^\top]W^\top + 2\mathbb{E}[\vec{z}W^\top \vec{\epsilon}] + \mathbb{E}[\vec{\epsilon}\vec{\epsilon}^\top]$$

$$= WW^\top + 2\mathbb{E}[\text{tr}(\vec{\epsilon}\vec{z}^\top W)] + \sigma^2 I$$

$$= WW^\top + 2\text{tr}(\mathbb{E}[\vec{\epsilon}\vec{z}^\top] W) + \sigma^2 I$$

$$= WW^\top + \sigma^2 I \quad \text{"0 b/c cov of 2 ind r.v.s"}$$

$$\Rightarrow p(\vec{x}) = \mathcal{N}(\vec{x} | \vec{\mu}, WW^\top + \sigma^2 I)$$

recall: when $p(\vec{x}) = \mathcal{N}(\vec{x} | \vec{\mu}, \Lambda')$ is a marginal Gaussian
and $p(\vec{y} | \vec{x}) = \mathcal{N}(\vec{y} | A\vec{x} + \vec{b}, L')$ is a conditional Gaussian
then $p(\vec{y}) = \mathcal{N}(\vec{y} | A\vec{\mu} + \vec{b}, L' + A\Lambda'A')$
and $p(\vec{x} | \vec{y}) = \mathcal{N}(\vec{x} | \Sigma \{ A^T L (\vec{y} - \vec{b}) + A\vec{\mu} \}, \Sigma)$
where $\Sigma = (\Lambda + A^T L A)^{-1}$

Find posterior: $p(\vec{z} | \vec{x}) = \mathcal{N}(\vec{z} | (I + W^T \sigma^{-2} I W)^{-1} \{ W \sigma^{-2} I (\vec{x} - \vec{\mu}) + I \vec{0} \}, (I + W^T \sigma^{-2} I W)^{-1})$
 $= \mathcal{N}(\vec{z} | M^{-1} W^T (\vec{x} - \vec{\mu}), \sigma^2 M^{-1})$
where $M = W^T W + \sigma^2 I$

likelihood $p(\vec{x} | \vec{\mu}, W, \sigma^2) = \prod_{n=1}^N \mathcal{N}(\vec{x}_n | \vec{\mu}, W W^T + \sigma^2 I)$

log likelihood $\ln p(\vec{x} | \vec{\mu}, W, \sigma^2) = \sum_{n=1}^N [-\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |WW^T + \sigma^2 I| - \frac{1}{2} (\vec{x}_n - \vec{\mu})^T (WW^T + \sigma^2 I)^{-1} (\vec{x}_n - \vec{\mu})]$

$$\vec{\mu}_{ML} = \bar{\vec{x}}$$

$$W_{ML} = U_M (L_M - \sigma^2 I)^{1/2}$$

\downarrow orthogonal matrix
 \downarrow diagonal matrix of M largest eigenvalues
 \downarrow cols are M eigenvectors

$$\sigma^2_{ML} = \frac{1}{D-M} \sum_{i=M+1}^D \lambda_i$$

$E[\vec{z} | \vec{x}] = M^{-1} W_{ML}^T (\vec{x} - \bar{\vec{x}})$ in latent space maps to $W_{ML} E[\vec{z} | \vec{x}] + \vec{\mu}_{ML}$ in data space
where $M = W^T W + \sigma^2 I$

Kernel PCA

given: $\{\vec{x}_i\} \in \mathbb{R}^D$ of N observations

principal components are given by the eigenvectors of the data covariance matrix $S = \frac{1}{N} \sum_{n=1}^N \vec{x}_n \vec{x}_n^T$

$$S \vec{v}_i = \lambda_i \vec{v}_i \quad \text{with } \|\vec{v}_i\| = 1$$

consider a nonlinear transformation $\vec{\Phi}: \mathbb{R}^D \rightarrow \mathbb{R}^M$

each data point \vec{x}_n is projected onto $\vec{\Phi}(\vec{x}_n)$ in feature space

$$\text{assume } \sum_n \vec{\Phi}(\vec{x}_n) = 0$$

sample covariance matrix in feature space is $C = \frac{1}{N} \sum_n \vec{\Phi}(\vec{x}_n) \vec{\Phi}(\vec{x}_n)^T$

principal components are given by the eigen vector equations

$$C \vec{v}_i = \lambda_i \vec{v}_i \quad (\text{but solving this directly would be } O(M^3))$$

note: $\lambda_i \vec{v}_i = \frac{1}{N} \sum_{n=1}^N \vec{\Phi}(\vec{x}_n) [\vec{\Phi}(\vec{x}_n)^T \vec{v}_i]$, so \vec{v}_i is a linear combination of $\vec{\Phi}(\vec{x}_n)$

$$\text{write } \vec{v}_i = \sum_{n=1}^N a_{in} \vec{\Phi}(\vec{x}_n)$$

$$\text{then } \frac{1}{N} \sum_{n=1}^N \vec{\Phi}(\vec{x}_n) \vec{\Phi}(\vec{x}_n)^T \vec{v}_i = \lambda_i \vec{v}_i$$

$$\frac{1}{N} \sum_{n=1}^N \vec{\Phi}(\vec{x}_n) \vec{\Phi}(\vec{x}_n)^T \sum_{m=1}^N a_{im} \vec{\Phi}(\vec{x}_m) = \lambda_i \sum_{m=1}^N a_{im} \vec{\Phi}(\vec{x}_m)$$

$$\vec{\Phi}(\vec{x}_i)^T \frac{1}{N} \sum_{m=1}^N \vec{\Phi}(\vec{x}_m) \vec{\Phi}(\vec{x}_m)^T \vec{\Phi}(\vec{x}_i) a_{im} = \vec{\Phi}(\vec{x}_i)^T \lambda_i \sum_{m=1}^N \vec{\Phi}(\vec{x}_m) a_{im}$$

$$\frac{1}{N} \sum_{m=1}^N k(\vec{x}_i, \vec{x}_m) k(\vec{x}_m, \vec{x}_i) a_{im} = \lambda_i \sum_{m=1}^N k(\vec{x}_i, \vec{x}_m)$$

$$\text{where } k(\vec{x}_i, \vec{x}_m) = \vec{\Phi}(\vec{x}_i)^T \vec{\Phi}(\vec{x}_m)$$

this is a system of linear equations, so we can write it as $K \vec{a}_i = N \lambda_i K \vec{a}_i$

find \vec{a}_i by solving $K \vec{a}_i = \lambda_i N \vec{a}_i$

$$\text{note: } \|\vec{v}_i\| = 1 \Rightarrow 1 = \vec{v}_i^T \vec{v}_i = (\sum_{n=1}^N a_{in} \vec{\Phi}(\vec{x}_n))^T (\sum_{m=1}^N a_{im} \vec{\Phi}(\vec{x}_m)) = \vec{a}_i^T K \vec{a}_i = \lambda_i N \vec{a}_i^T \vec{a}_i$$

note: we can recover \vec{v}_i with \vec{o}_i and get the principal components

note: to find the projection of a new point \vec{x}

$$\text{coefficients } y_i(\vec{x}) = \vec{\Phi}(\vec{x})^T \vec{v}_i = \sum_{n=1}^N a_{in} \vec{\Phi}(\vec{x})^T \vec{\Phi}(\vec{x}_n) = \sum_{n=1}^N a_{in} k(\vec{x}, \vec{x}_n)$$

applications of PCA

data compression - storing data directly = ND space

storing compressed data = D+NM+MD space

store $\sum_{i=1}^M (\vec{v}_i^T \vec{x}) \vec{v}_i$, \vec{x}_n , and \vec{v}_i for $i=1, \dots, M$

"D space" "NM space" "MD space"

saves space if $N \gg M$ or $D \gg M$

data visualization - project onto \mathbb{R}^2 or \mathbb{R}^3 and plotting

data whitening - transform new data st. it has mean 0 and unit covariance

stack eigenvector equations into $SU=UL$

$U \in \mathbb{R}^{D \times D}$ has columns \vec{v}_i :

$L \in \mathbb{R}^{D \times D}$ is a diagonal matrix st. $L_{ii} = \lambda_i$

transform data using $\vec{y}_n = L^{-\frac{1}{2}} U^T (\vec{x}_n - \vec{x})$

$$\begin{aligned} \sum_n \vec{y}_n &= \sum_n L^{-\frac{1}{2}} U^T (\vec{x}_n - \vec{x}) \\ &= L^{-\frac{1}{2}} U^T \left(\sum_n \vec{x}_n - N\vec{x} \right) \\ &= L^{-\frac{1}{2}} U^T \vec{O} \\ \Rightarrow \text{mean} &= 0 \end{aligned}$$
$$\begin{aligned} \frac{1}{N} \sum_n (\vec{y}_n - \vec{O})(\vec{y}_n - \vec{O})^T &= \frac{1}{N} \sum_n L^{-\frac{1}{2}} U^T (\vec{x}_n - \vec{x})(\vec{x}_n - \vec{x})^T U L^{-\frac{1}{2}} \\ &= L^{-\frac{1}{2}} U^T \left[\frac{1}{N} \sum_n (\vec{x}_n - \vec{x})(\vec{x}_n - \vec{x})^T \right] U L^{-\frac{1}{2}} \\ &= L^{-\frac{1}{2}} U^T S U L^{-\frac{1}{2}} \\ &= L^{-\frac{1}{2}} U^T U L L^{-\frac{1}{2}} \\ &= L^{-\frac{1}{2}} L L^{-\frac{1}{2}} = I \\ \Rightarrow \text{unit covariance} \end{aligned}$$

PCA for higher-dimensional data

if $N < D$, we can compute eigenvalues and eigenvectors of $S \in \mathbb{R}^{N \times N}$ instead of $S \in \mathbb{R}^{D \times D}$

let $X \in \mathbb{R}^{N \times D}$ where the nth row is $(\vec{x}_n - \vec{x})^T$

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

the eigenvalue equations are $\frac{1}{N} X^T X \vec{v}_i = \lambda_i \vec{v}_i$

$$\frac{1}{N} X X^T (\vec{v}_i) = \lambda_i (X \vec{v}_i)$$

$$\frac{1}{N} X X^T (\vec{v}_i) = \lambda_i \vec{v}_i$$

so $X X^T \in \mathbb{R}^{N \times N}$ has the same eigenvalues λ_i as $X^T X \in \mathbb{R}^{D \times D}$

and the eigenvectors differ by a linear transformation $\vec{v}_i = X \vec{j}_i$

K-Means Clustering

$\{\vec{\mu}_k\}_{k=1}^K \in \mathbb{R}^P$ represents the center of each cluster

$r_{nk} \in \{0, 1\}$ is a binary indicator for if the nth point belongs to the kth cluster

i.e. $r_{nk} = \begin{cases} 1 & \text{if } \vec{x}_n \in \text{cluster } k \\ 0 & \text{otherwise} \end{cases}$

then the distortion (error measure) is $J = \sum_{n=1}^N \sum_k r_{nk} \|\vec{x}_n - \vec{\mu}_k\|^2$

K-means algorithm

1) fix $\hat{\mu}_k$ and optimize J wrt r_{nk}

$$r_{nk}=1 \text{ for } k = \arg\min_j \| \vec{x}_n - \hat{\mu}_j \|^2$$

$r_{nk}=0$ otherwise

i.e. assign each point to its nearest center

2) fix r_{nk} and optimize J wrt $\hat{\mu}_k$

$$\frac{\partial J}{\partial \hat{\mu}_k} = -2 \sum_{n=1}^N r_{nk} (\vec{x}_n - \hat{\mu}_k) = 0$$

$$\sum_{n=1}^N r_{nk} \vec{x}_n = \sum_{n=1}^N r_{nk} \hat{\mu}_k$$

$$\hat{\mu}_k = \frac{\sum_{n=1}^N r_{nk} \vec{x}_n}{\sum_{n=1}^N r_{nk}} = \frac{\sum_{n=1}^N \vec{x}_n}{|C_k|}$$

where $C_k = \{n \mid r_{nk}=1\} \Rightarrow |C_k| = \text{mean of points in } C_k$

3) repeat

every iteration decreases J and $J \geq 0$ so convergence is guaranteed

might converge to local min

K-medoids: $J = \sum_{n=1}^N \sum_{k=1}^K r_{nk} v(\vec{x}_n, \vec{\mu}_k)$

K-means application

image segmentation and compression: reduce an image into K colors using \mathbb{R}^3 vectors (RGB value of pixels)

not super great because it doesn't consider spatial proximity of pixels

lossy data compression: store the cluster of each data point and the k centers

often called vector quantization

$\vec{\mu}_k$ are codebook vectors

EM Algorithm (General)

X : set of all data

Z : set of all latent variables

Θ : set of all model parameters

log likelihood $\ln p(X|\Theta) = \ln \left\{ \sum_z p(X, Z|\Theta) \right\}$

we only observe incomplete data X and not complete data $\{X, Z\}$

cannot maximize likelihood directly \Rightarrow use iterative process

1. initialize Θ^{old}

2. evaluate $p(X|Z, \Theta^{old})$

3. evaluate $\Theta^{new} = \operatorname{argmax}_{\Theta} Q(\Theta, \Theta^{old})$

where $Q(\Theta, \Theta^{old}) = \sum_z p(X|Z, \Theta^{old}) \ln p(X, Z|\Theta)$

4. set $\Theta^{old} \leftarrow \Theta^{new}$, repeat until convergence

Mixtures of Gaussians

let \vec{z} be a 1-of-K coding (aka one-hot coding) vector

$$z_k \in \{0, 1\} \text{ for } k=1, \dots, K \text{ and } \sum_{k=1}^K z_k = 1$$

$$\text{suppose } p(z_k=1) = \pi_k$$

$$\text{then } 0 \leq \pi_k \leq 1, \sum_{k=1}^K \pi_k = 1, \text{ and } p(\vec{z}) = \prod_{k=1}^K \pi_k^{z_k}$$

assume the conditional is Gaussian: $p(\vec{x}|z_k=1) = \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$

$$p(\vec{x}|\vec{z}) = \prod_{k=1}^K \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)^{z_k}$$

$$\text{Gaussian mixture model: } p(\vec{x}) = \sum_{\vec{z}} p(\vec{x}|\vec{z}) p(\vec{z}) = \sum_{k=1}^K \pi_k \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)$$

$$\text{log likelihood: } \ln p(X|\vec{\pi}, \vec{\mu}, \Sigma) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(\vec{x}_n|\vec{\mu}_k, \Sigma_k) \right)$$

unable to solve via the usual ML approach bc of singularities

$$\text{define } \gamma(z_k) := p(z_k=1 | \vec{x}) = \frac{p(\vec{x}, \vec{z})}{p(\vec{x})} = \frac{p(z_k=1)p(\vec{x}|z_k=1)}{\sum_{j=1}^K p(z_j=1)p(\vec{x}|z_j=1)} = \frac{\pi_k \mathcal{N}(\vec{x}|\vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\vec{x}|\vec{\mu}_j, \Sigma_j)} = \text{"responsibility" that component } k \text{ takes for "explaining" } \vec{x}$$

EM algorithm for GMM

1. initialize $\vec{\mu}_k, \Sigma_k, \pi_k$

$$2. E: \text{evaluate } \gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\vec{x}_n|\vec{\mu}_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\vec{x}_n|\vec{\mu}_j, \Sigma_j)}$$

3. M: re-estimate parameters

$$\boxed{\begin{aligned} \vec{\mu}_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \vec{x}_n \\ \Sigma_k &= \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\vec{x}_n - \vec{\mu}_k)(\vec{x}_n - \vec{\mu}_k)^T ; \quad N_k = \sum_{n=1}^N \gamma(z_{nk}) \\ \pi_k &= \frac{N_k}{N} \end{aligned}}$$

4. repeat until convergence

$$\text{note: } p(X, Z | \mu, \Sigma, \pi) = \prod_{n=1}^N \prod_{k=1}^K \pi_k^{z_{nk}} \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k)^{z_{nk}}$$

$$\Rightarrow \ln p(X, Z | \mu, \Sigma, \pi) = \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \}$$

by Bayes' Theorem, $p(Z|X, \mu, \Sigma, \pi) \propto p(X, Z | \mu, \Sigma, \pi)$

$$\Rightarrow \ln p(Z|X, \mu, \Sigma, \pi) \propto \ln p(X, Z | \mu, \Sigma, \pi) \\ \propto \sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \}$$

$$\mathbb{E}_z [\ln p(X, Z | \mu, \Sigma, \pi)] = \mathbb{E}_z \left[\sum_{n=1}^N \sum_{k=1}^K z_{nk} \{ \ln \pi_k + \ln \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \} \right] \\ = \sum_{n=1}^N \sum_{k=1}^K \mathbb{E}_z [z_{nk}] \{ \ln \pi_k + \ln \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \} \\ = \sum_{n=1}^N \sum_{k=1}^K \gamma(z_{nk}) \{ \ln \pi_k + \ln \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k) \}$$

maximizing this expectation wrt μ_k, Σ_k , and π_k gives the M step (step 3)

Committee Methods

idea: improve accuracy by aggregating the results of many models

bootstrapping - given a training set $\{X_i\}_{i=1}^N$

generate a bootstrap training set by randomly sampling N points with replacement
train L models on each of L bootstrap training set

boosting - weight each data point, giving more weight to points labeled incorrectly by the previous training iteration

1. initialize weights to be equal, $w_n^{(1)} = \frac{1}{N}$ for $n=1, \dots, N$

2. fit a classifier $y_m(\vec{x})$ to the data

$$\text{minimize weighted error } J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n)$$

where $I(y_m(\vec{x}_n) \neq t_n) = \begin{cases} 1 & \text{if } y_m(\vec{x}_n) \neq t_n \\ 0 & \text{otherwise} \end{cases}$

3. evaluate

$$E_m = \frac{\sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n)}{\sum_{n=1}^N w_n^{(m)}}$$

then evaluate $\alpha_m = \ln \left\{ \frac{1-E_m}{E_m} \right\}$

4. update weights

$$w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(\vec{x}_n) \neq t_n) \}$$

5. repeat 2-4 for $m=1, \dots, M$

6. final prediction

$$Y_m(\vec{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m y_m(\vec{x}) \right)$$

note: initially $0 \leq E \leq \frac{1}{2}$

$$\Rightarrow \alpha_0 \geq \alpha \geq 0$$

many errors, update lots of weights slightly

Fewer errors, update weights aggressively

equivalently, minimize exponential error function wrt weights α_m and base classifiers y_m

$$E = \sum_{n=1}^N \exp \{ -t_n f_m(\vec{x}_n) \}$$

where $f_m(\vec{x}_n) = \frac{1}{2} \sum_{k=1}^d \alpha_k y_k(\vec{x}_n)$

note: $E = \sum_{n=1}^N \exp \{ -t_n f_m(\vec{x}_n) \}$

$$= \sum_{n=1}^N \exp \{ -t_n f_{m-1}(\vec{x}_n) - \frac{1}{2} t_n \alpha_m y_m(\vec{x}_n) \}$$

$$= \sum_{n=1}^N w_n^{(m)} \exp \{ -\frac{1}{2} t_n \alpha_m y_m(\vec{x}_n) \}$$

where $w_n^{(m)} = \exp \{ -t_n f_{m-1}(\vec{x}_n) \}$

when minimizing wrt α_m and y_m , we can view $w_n^{(m)}$ as constants

let $T_m = \{ \text{data points correctly classified by } y_m \}$

$M_m = \{ \text{data points incorrectly classified by } y_m \}$

then $E = \sum_{n=1}^N w_n^{(m)} \exp \{ -\frac{1}{2} t_n \alpha_m y_m(\vec{x}_n) \}$

$$= \sum_{n \in T_m} w_n^{(m)} \exp \left\{ -\frac{\alpha_m}{2} t_n y_m(\vec{x}_n) \right\} + \sum_{n \in M_m} w_n^{(m)} \exp \left\{ -\frac{\alpha_m}{2} t_n y_m(\vec{x}_n) \right\}$$

↑ b/c classified correctly

↑ b/c classified incorrectly

$$= \exp \left\{ -\frac{\alpha_m}{2} \sum_{n \in T_m} w_n^{(m)} \right\} + \exp \left\{ \frac{\alpha_m}{2} \sum_{n \in M_m} w_n^{(m)} \right\}$$

$$= \exp \left\{ -\frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) = t_n) \right\} + \exp \left\{ \frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n) \right\}$$

$$= \exp \left\{ -\frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} - \exp \left\{ -\frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n) \right\} + \exp \left\{ \frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n) \right\} \right\}$$

$$= \exp \left\{ -\frac{\alpha_m}{2} \sum_{n=1}^N w_n^{(m)} + \left(\exp \left\{ \frac{\alpha_m}{2} \right\} - \exp \left\{ -\frac{\alpha_m}{2} \right\} \right) \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n) \right\}$$

when minimizing E wrt y_m , $\exp\left\{-\frac{\alpha_m}{2}\right\} \sum_{n=1}^N w_n^{(m)}$ is constant

so we can minimize $J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(\vec{x}_n) \neq t_n)$ (see boosting 2)

when minimizing E wrt α_m ,

we get $\alpha_m = \ln\left\{\frac{1-E_m}{E_m}\right\}$ (see boosting 3)

when updating the weights,

$$\begin{aligned} w_n^{(m+1)} &= w_n^{(m)} \exp\left\{-\frac{1}{2} t_n \alpha_m y_m(\vec{x}_n)\right\} \\ &= w_n^{(m)} \exp\left\{-\frac{\alpha_m}{2} (1 - 2I(y_m(\vec{x}_n) \neq t_n))\right\} \\ &= w_n^{(m)} \exp\left\{-\frac{\alpha_m}{2}\right\} \exp\{\alpha_m I(y_m(\vec{x}_n) \neq t_n)\} \end{aligned}$$

since $\exp\left\{-\frac{\alpha_m}{2}\right\}$ is independent of n , it would affect all weights equally and so can be discarded

thus we get $w_n^{(m+1)} = w_n^{(m)} \exp\{\alpha_m I(y_m(\vec{x}_n) \neq t_n)\}$ (see boosting 4)

note: $t_n y_m(\vec{x}_n) = \begin{cases} 1 & \text{if } t_n = y_m(\vec{x}_n) \\ -1 & \text{if } t_n \neq y_m(\vec{x}_n) \end{cases}$

so $t_n y_m(\vec{x}_n) = 1 - 2I(y_m(\vec{x}_n) \neq t_n)$

where $I(y_m(\vec{x}_n) \neq t_n) = \begin{cases} 1 & \text{if } y_m(\vec{x}_n) \neq t_n \\ 0 & \text{if } y_m(\vec{x}_n) = t_n \end{cases}$