

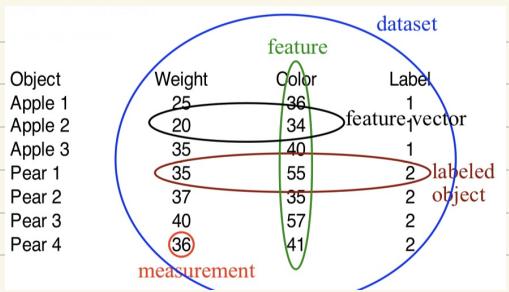
Generalization training/test set

Features

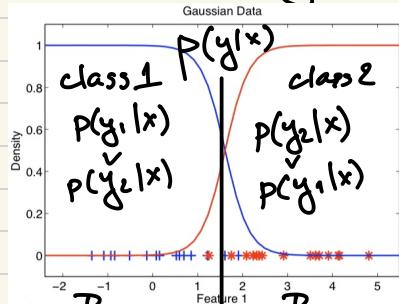
measurable data encoded
in vector $x = [x_1, x_2, \dots, x_n]^T$

Dataset

features \times label table



$p(x, y)$ - probability density over feature space



$$\sum p(y_i|x) = 1$$

Bayes Theorem

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

(conditional) class distribution
(unc.) data distribution \sim class prior

decision boundary $p(y_1|x) = p(y_2|x)$

Error of decision boundary

type I

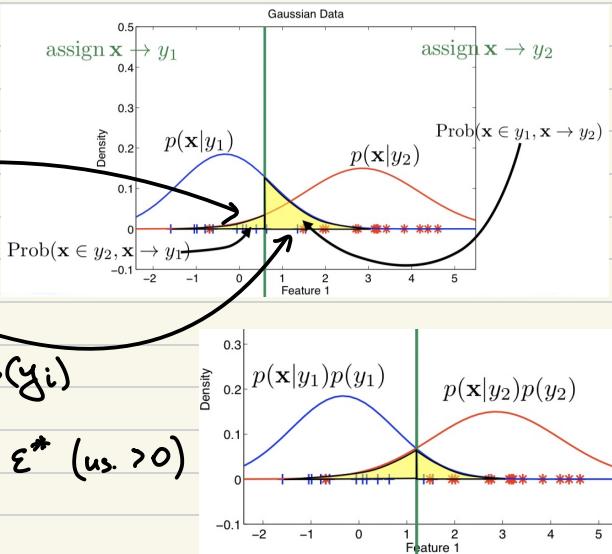
$$\varepsilon_1 = \int_{R_2} p(x|y_1) dx$$

type II

$$\varepsilon_2 = \int_{R_1} p(x|y_2) dx$$

$$p(\text{error}) = \sum_{i=1}^2 p(\text{error}|y_i) p(y_i)$$

Bayes (minimum) error ε^* (≈ 0)



Misclassification

mistake (true) y_i for $\hat{y}_i \rightarrow \hat{\alpha}_{ij}$
dataset $D = \{(x_i, y_i)\}_{i=1}^N$

estimated label \hat{y}_i

empirical risk $R = \frac{1}{N} \sum_{i=1}^N \lambda_{y_i \neq \hat{y}_i}$

conditional risk $\hat{L}(x) = \sum_{j=1}^C \hat{\alpha}_{ji} p(y_j|x)$

(assign x to class y_j)

average risk over region $r_i = \int_{R_i} \hat{L}(x) p(x) dx$

overall risk $r = \sum_{i=1}^C r_i = \sum_{i=1}^C \int_{R_i} \sum_{j=1}^C \hat{\alpha}_{ji} p(y_j|x) p(x) dx$

$x \in R_i$ if $\sum_{j=1}^C \hat{\alpha}_{ji} p(y_j|x) \leq \sum_{j=1}^C \hat{\alpha}_{jk} p(y_j|x) \quad k=1, \dots, C \quad \lambda_{y_i, y_j} = 0$

two-class problem $i=1 \quad \hat{\alpha}_{11} p(y_1|x) \quad i=2 \quad \hat{\alpha}_{12} p(y_2|x)$

4. Parametric Density-based Classifiers

$$\text{Bayes: } p(y_1|x) = \frac{p(x|y_1)p(y_1)}{p(x)} = \frac{p(x|y_2)p(y_2)}{p(x)} = p(y_2|x)$$

Model - approximate $p(y|x)$ or $p(x|y)$, $p(y)$

parametric vs. non-parametric

Discriminative

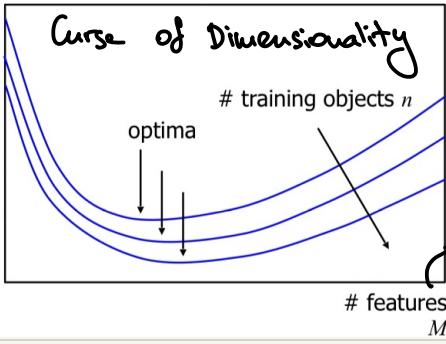
vs.

Generative

$\hat{p}(y|x) \rightarrow$ directly approximate

$$p(y|x) \propto p(x|y)p(y)$$

test error



Gaussian Distribution

$$\mu = 0, \sigma^2 = 1, [\mu - 2\sigma, \mu + 2\sigma] = 95\% \text{ data}$$

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

$$\Sigma^2 = \Sigma^1 \text{ (covariance matrix)}$$

$$p(x) = \frac{1}{\sqrt{(2\pi)^n \det(\Sigma)}} \exp\left(-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$M = \mathcal{N}(x|\mu, \Sigma)$$

2D:

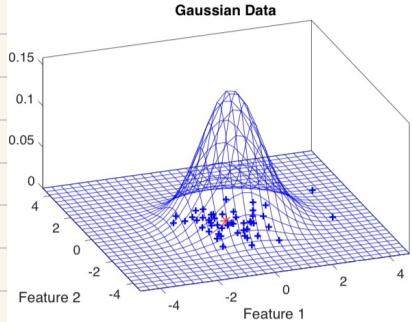
$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 \\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

Quadratic Classifier



$$f(x) = \log p(y_1|x) - \log p(y_2|x) = x^T \omega x + \omega^T x + \omega_0$$

$$\log(\hat{p}(y_j|x)) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\det \Sigma_j) - \frac{1}{2} (x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j) + \log p(y_j) - \log p(x)$$

$$\omega = \frac{1}{2} (\Sigma_1^{-1} - \Sigma_2^{-1})$$

$$\omega^T = \mu_1^T \Sigma_1^{-1} - \mu_2^T \Sigma_2^{-1}$$

$$\omega_0 = -\frac{1}{2} \log \det \Sigma_1 - \frac{1}{2} \mu_1^T \Sigma_1^{-1} \mu_1 + \log p(y_1) + \frac{1}{2} \log \det \Sigma_2 + \frac{1}{2} \mu_2^T \Sigma_2^{-1} \mu_2 - \log p(y_2)$$

$$\hat{\Sigma} = \frac{1}{C} \sum_{i=1}^C \hat{\Sigma}_k \leftarrow \text{assume shared covariance}$$

$$g_i(x) = -\frac{1}{2} \log(\det(\hat{\Sigma})) - \frac{1}{2} (x - \mu_i)^T \hat{\Sigma}^{-1} (x - \mu_i) + \log p(y_i)$$

$$\hookrightarrow g_i(x) = -\frac{1}{2} \mu_i^T \hat{\Sigma}^{-1} \mu_i + \mu_i^T \hat{\Sigma}^{-1} x + \log p(y_i)$$

LDA

$$f(x) = w^T x + w_0 \quad \text{linear normal-based classifier}$$

$$w = \hat{\Sigma}^{-1} (\hat{\mu}_1 - \hat{\mu}_2)$$

$$w_0 = \frac{1}{2} \hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_2 - \frac{1}{2} \hat{\mu}_1^T \hat{\Sigma}^{-1} \hat{\mu}_1 + \log p(y_1) - \log p(y_2)$$

$$x \rightarrow y_1 \text{ if } w^T x + w_0 \geq 0$$

$$\hookrightarrow y_2 \text{ if } w^T x + w_0 < 0$$

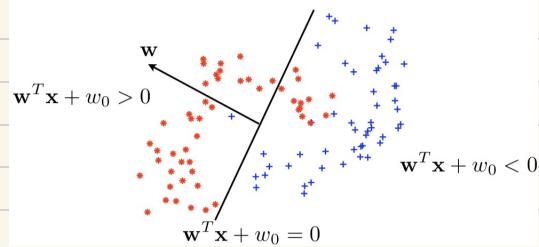
Nearest Mean

$$\hat{\Sigma} = \sigma^2 I \quad \text{meaning } \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

$$g_i(x) = -\frac{1}{\sigma^2} \left(\frac{1}{2} \hat{\mu}_i^T \hat{\mu}_i - \hat{\mu}_i^T x \right) + \log(p(y_i))$$

$$w = \hat{\mu}_1 - \hat{\mu}_2$$

$$w_0 = \frac{1}{2} \hat{\mu}_2^T \hat{\mu}_2 - \frac{1}{2} \hat{\mu}_1^T \hat{\mu}_1 - \sigma^2 \log \frac{p(y_1)}{p(y_2)}$$



4. Non-parametric Density-based Classifiers

histogram density estimation $\hat{p}(x) = \frac{1}{w} \frac{k_B}{N} \rightarrow$ objects in bin width of bin $\leftarrow w \frac{k_B}{N} \rightarrow \# \text{ of objects}$

Parzen (window) / Kernel Density Estimation

Parzen probability density $\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x-x_i}{w}\right)$

$(x|y_i) = \frac{1}{w} \sum_{j=1}^{n_i} N(x|x_j^{(i)}, w\Sigma) \leftarrow \text{Gaussian kernel, } \S \text{ as covariance}$

K-Nearest Neighbours

find k nearest neighbours to point

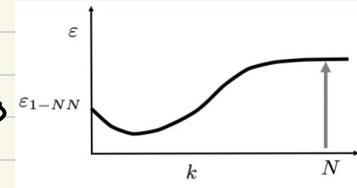
$$\hat{p}(x|y_i) = \frac{k_i}{n_i V_k}$$

\rightarrow volume of sphere of k nearest neighbours

$$\hat{p}(y_i) = \frac{n_i}{n}$$

$$\Rightarrow \hat{p}(x|y_i) \hat{p}(y_i) > \hat{p}(x|y_j) \hat{p}(y_j) \equiv k_i > k_j$$

$k \xrightarrow{\text{large}} \text{classify as most probable class}$
 $\xrightarrow{\text{small}} \text{unstable decision boundaries}$



ties \rightarrow odd k (solution for 2-class)

\rightarrow random

\rightarrow greater prior

\rightarrow nearest (1-nn)

Distance

$$\text{euclidean } D(x, x') = \sqrt{\sum_d |x_d - x'_d|^2}$$

$$\text{manhattan } D(x, x') = \sum_d |x_d - x'_d|$$

$$\text{hamming } D(x, x') = \sum_d \frac{1}{x_d \neq x'_d}$$

Naive Bayes

$$p(x_1, x_2, \dots, x_d | y)$$

assumption: independent features

conditional independence: $p(B, S | X) = p(B|X)p(S|X)$

$$\Rightarrow p(x|y) = \prod_{i=1}^d p(x_i|y)$$

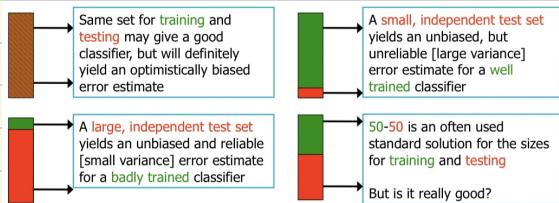
$$p(x|y) \neq 0 \text{ instead } \frac{\text{num}(x,y) + \epsilon}{\text{num}(y) + K\epsilon} \rightarrow \text{error}$$

$\# \text{ classes}$

Classifier Evaluation

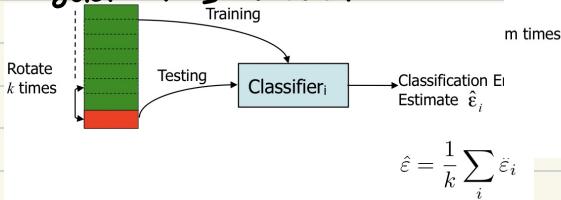
$$\hat{\epsilon} = \frac{1}{N} \sum_{i=1}^N z_i \quad z_i = \begin{cases} 0 & \text{if } x_i \text{ correctly classified} \\ 1 & \text{if } x_i \text{ incorrectly classified} \end{cases}$$

$$\sigma_{\hat{\epsilon}}^2 = \text{Var}(\hat{\epsilon} | \text{test set size } N) = \frac{\hat{\epsilon}(1-\hat{\epsilon})}{N} \quad \sigma_{\hat{\epsilon}} = \sqrt{\frac{\hat{\epsilon}(1-\hat{\epsilon})}{N}}$$

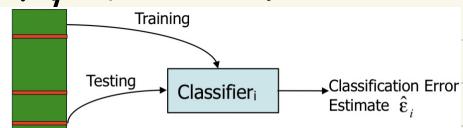


Bootstrapping
randomly draw N/N objects with replacement
left-over \rightarrow testing
repeat m times

k-fold cross-validation



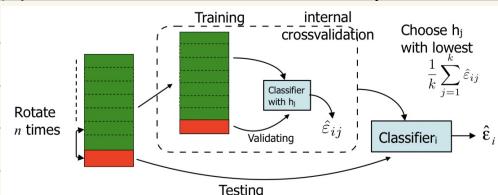
$$\hat{\epsilon} = \frac{1}{k} \sum_i \hat{\epsilon}_i$$



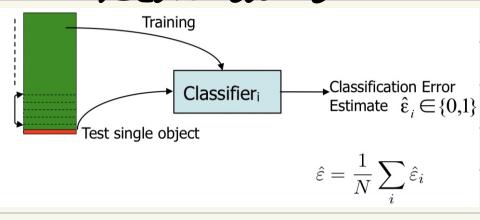
$$\hat{\epsilon} = \frac{1}{m} \sum_i \hat{\epsilon}_i$$

Leave-one-out

Double cross-validation

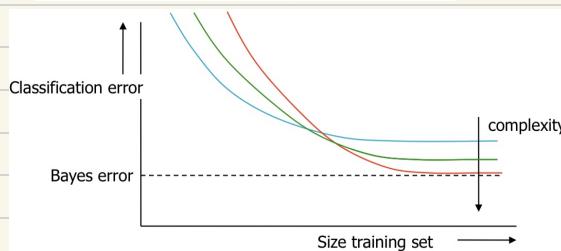
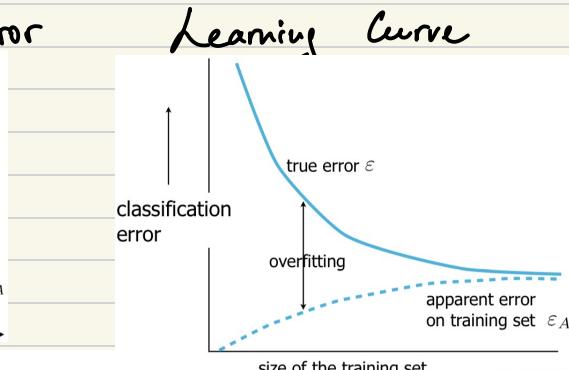
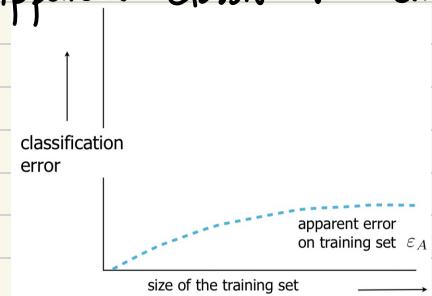


$$\hat{\epsilon} = \frac{1}{k} \sum_{j=1}^k \hat{\epsilon}_{ij}$$



$$\hat{\epsilon} = \frac{1}{N} \sum_i \hat{\epsilon}_i$$

Apparent Classification error



- Larger training set \rightarrow better classifier
- independent test set \rightarrow unbiased error estimates
- Larger test set \rightarrow more accurate error estimates
- more complex classifiers \uparrow larger training set
- larger feature set size

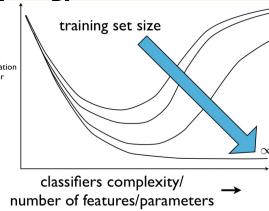
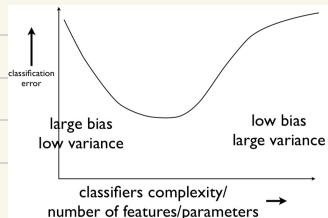
$$\text{Squared error: } L(\omega) = E \left[\|g(x) - y\|^2 \right]$$

$\mathcal{D} = \{(y_i, x_i); i = 1, \dots, N\}$

$$MSE = ED[(g(x; D) - E_D[g(x; D)])^2] + ED[(ED[g(x; D)] - E[g|x])^2]$$

variance bias

variance - how much \hat{y} vary over different training sets
bias - how much average \hat{y} differ from true output



2-class classification

$$\text{standard error } \varepsilon = \varepsilon_1 p(y_1) + \varepsilon_2 p(y_2)$$

$$\text{weighted error } \varepsilon = \lambda_{11} \varepsilon_1 p(y_1) + \lambda_{21} \varepsilon_2 p(y_2)$$

$$\text{F1-score (harmonic mean)} \quad F_1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

error - probability of erroneous classification
performance / accuracy = $1 - \text{error}$

sensitivity of class - performance of objects from class
specificity - performance for all objects outside class

precision (of class) - fraction of correctly assigned to class

recall - fraction of correctly classified objects

true positive rate - sensitivity

false positive rate - error for all objects outside class

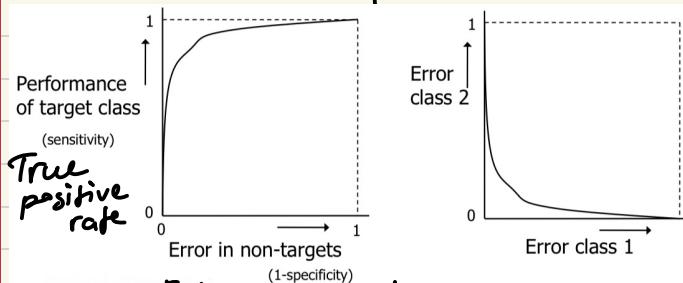
Confusion Matrix

$$A = \begin{bmatrix} \lambda_1 \\ \dots \\ \lambda_N \end{bmatrix} \quad \text{real labels}$$

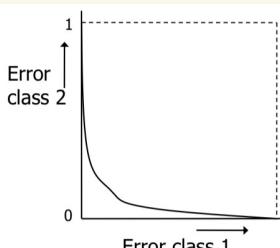
$$L = \begin{bmatrix} l_1 \\ \dots \\ l_N \end{bmatrix} \quad \text{predicted labels}$$

$$C = \begin{bmatrix} c_{11} & \dots & c_{1k} \\ \dots & \dots & \dots \\ c_{k1} & \dots & c_{kk} \end{bmatrix} \quad c_{ij} = N \cdot P[l_j | \lambda_i]$$

ROC (Receiver-Operator Characteristic) Analysis



False positive rate



Parzen / Kernel density estimation

- doesn't assume known distribution
- estimates probability densities using kernel function
- uses kernel function of fixed shape and width
- width matters

K - Nearest Neighbour

Pros

- Simple and flexible
- good performance
- simple to adapt complexity

Cons

- large training sets (stored)
- compute all distances
- sensibly scaled features
- optimized K

Naive Bayes

Pros

- handle high dimensional feature spaces
- fast training time
- handle continuous and discrete data

Cons

- can't deal with correlated features

4. Linear Discriminative Classifiers

$L(h(x), y) \rightarrow$ loss / cost function
 ↴ classifier

$$\min_{\mathbf{w} \in \mathbb{R}^{D+1}} \frac{1}{N} \sum_{i=1}^N L(h(\mathbf{x}_i), y_i)$$

$$h(\mathbf{x}) = \mathbf{x}^T \mathbf{w} + w_0$$

$y = \begin{cases} c_1 & \text{if } h(\mathbf{x}) > 0 \\ c_0 & \text{if } h(\mathbf{x}) \leq 0 \end{cases}$ decision boundary $h(\mathbf{x}) = 0 = \mathbf{w}^T \mathbf{x} + w_0$

$$\min_{\mathbf{w}, w_0 \in \mathbb{R}^{D+1}} \frac{1}{N} \sum_{i=1}^N L(\mathbf{x}_i^T \mathbf{w} + w_0, y_i)$$

$L:$ absolute $|x_i^T \mathbf{w} + w_0 - y_i|$
 squared $(x_i^T \mathbf{w} + w_0 - y_i)^2$

$$\min_{\mathbf{w}, w_0 \in \mathbb{R}^{D+1}} \frac{1}{N} \sum_{i=1}^N (x_i^T \mathbf{w} + w_0 - y_i)^2$$

$$= \min_{\mathbf{w}, w_0} J(\mathbf{w}, w_0)$$

$$\mathbf{w}_j^{t+1} = \mathbf{w}_j^t - \alpha \frac{\partial J(\mathbf{w}, w_0)}{\partial w_j} \quad |_{\mathbf{w}, w_0 = \mathbf{w}^t, w_0^t}$$

$$\frac{\partial J(\mathbf{w}, w_0)}{\partial w_j} = \frac{1}{N} \sum_{i=1}^N \lambda(x_i^T \mathbf{w} + w_0 - y_i) x_i^{(j)}$$

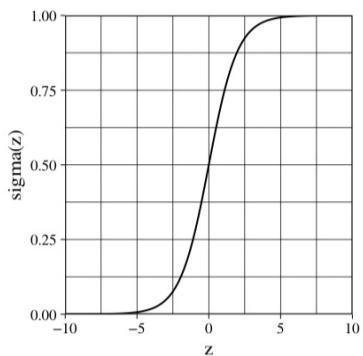
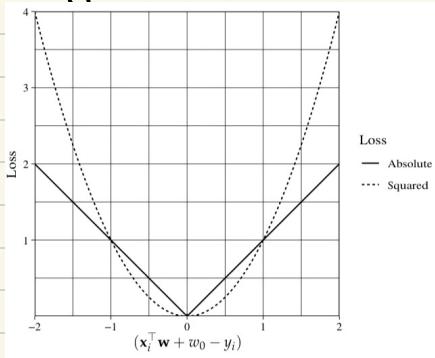
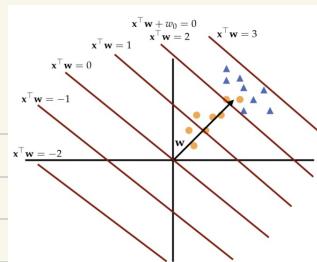
$$\frac{\partial J(\mathbf{w}, w_0)}{\partial w_0} = \frac{1}{N} \sum_{i=1}^N \lambda(w_0 + x_i^T \mathbf{w} - y_i)$$

$$h(\mathbf{x}) \approx P(y=1|\mathbf{x})$$

$$\text{logistic function: } \sigma(z) = \frac{1}{1 + e^{-z}}$$

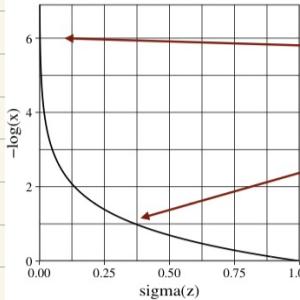
$$\text{label } y_i \text{ for } i = \begin{cases} \sigma(h(\mathbf{x}_i)), & y_i = 1 \\ 1 - \sigma(h(\mathbf{x}_i)), & y_i = 0 \end{cases}$$

$$L(h) = \prod_{i=1}^N \sigma(h(\mathbf{x}_i))^{y_i} (1 - \sigma(h(\mathbf{x}_i)))^{1-y_i}$$

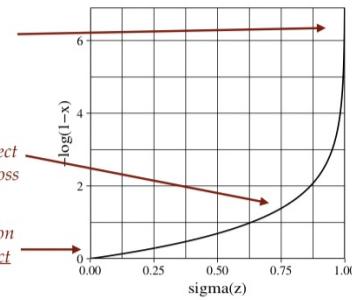


$$\mathcal{J}(h) = -\sum_{i=1}^N y_i \log(\sigma(h(x_i))) + (1-y_i) \log(1-\sigma(h(x_i)))$$

For $y_i = 1$



For $y_i = 0$



$$\min_{w, w_0 \in \mathbb{R}^{D+1}} -\sum_{i=1}^N y_i \log(\sigma(x_i^T w + w_0)) + (1-y_i) \log(1-\sigma(x_i^T w + w_0))$$

$$w_j^{\text{new}} = w_j - \alpha \frac{\partial \mathcal{J}(w, w_0)}{\partial w_j} \quad \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$\nabla_w \mathcal{J}(w) = -\sum_{i=1}^N (y_i - \sigma(x_i^T w)) x_i$$

Support Vector Machines

$$x_i^T w + w_0 \geq M \quad y_i = +1 \quad \frac{1}{x_i^T w + w_0} = \frac{1}{\|w\|} \Rightarrow \text{margin} = \frac{2}{\|w\|}$$

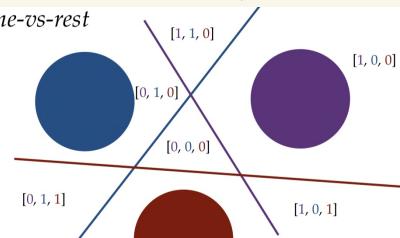
$$x_i^T w + w_0 \leq -M \quad y_i = -1$$

$$\max_{w, w_0} \frac{2}{\|w\|} \rightarrow \min_{w, w_0} \frac{1}{2} \|w\|^2$$

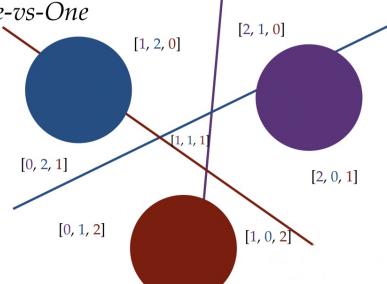
$$\text{Soft-Margin} \quad \min_{w, w_0} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^N (1 - y_i (x_i^T w + w_0))_+$$

$$P(y=k|x) = \frac{e^{-x^T w_k}}{\sum_{k=1}^K e^{-x^T w_k}}$$

One-vs-rest



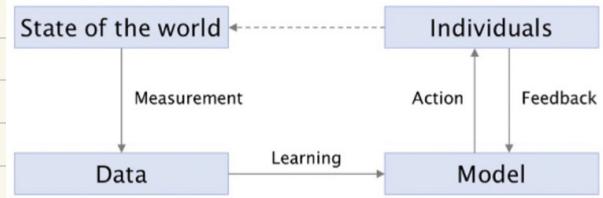
One-vs-One



Measurement
race:

- social construct
- changes over time
- different terms

Subjective judgements inherit stereotypes



Learning

Algorithms extract stereotypes, discrimination from data
remove stereotype-identifying information → what's left?

Action

Model reproduces patterns in data

Automated decision → only outcome is important

Ethics → decision process is important

Feedback

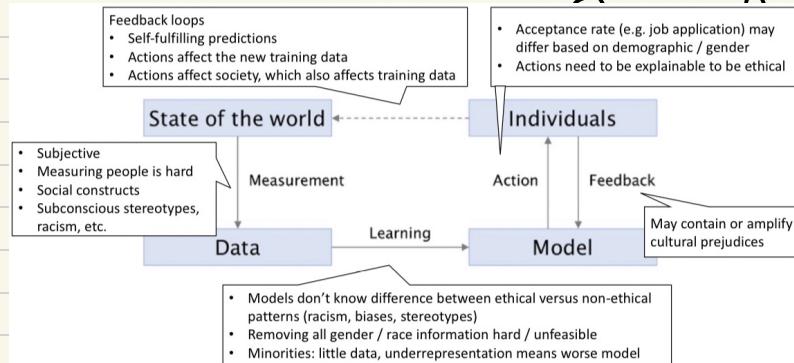
Feedback can reflect or amplify cultural prejudices

Loop

→ Self-fulfilling predictions

→ Predictions that affect the training data

→ Predictions that affect society at large



A - sensitive attribute

X - other features

Y - classification target

R - classifier output $R=0 \rightarrow \text{reject}$ $R=1 \rightarrow \text{accept}$

Criteria 1: Independence

$$P(R=1 | A=a) = P(R=1 | A=b)$$

Acceptance rate for each group \rightarrow same

Classification rule \rightarrow may differ for each group

Criteria 2: Separation

Fraction of mistakes for each group \rightarrow same

$$\epsilon_+ = \frac{FN}{TP+FN} = \frac{FP}{FP+TN} = \epsilon_-$$

Ethics - moral standard of 'right' and 'wrong' that prescribe how we ought to act

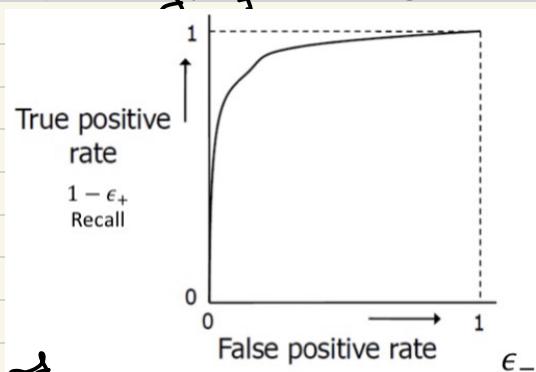
Descriptive - explanations of, and factual statements about moral systems people subscribe to, how we act and understand ourselves as moral beings (biology/psychology - anthropology)

Normative - prescriptive accounts on how we ought to act and understand ourselves as moral beings (philosophy)

\rightarrow Values (abstract goals or ideas)

\rightarrow Norms (rules about actions)

\rightarrow Virtues (habits of character)



Normative Ethical Theory - formulation of an account or set of principles that prescribe us how to act

- consequences of action
- action itself and norm it's based on
- agent

Consequentialism

action - morally required is consequences → maximize happiness
optimize best outcome
quantification and justification of trade-offs between alternatives
↳ minimize suffering

Duty Ethics

never use a person as a mere means to some desired end
person - 'ends in themselves'; always respected
rules for right action regardless of outcome
absoluteness of basic moral principles

Bias - preference or inclination for or against something

→ Natural / Evolutionary

→ Social / Historical

→ Technologically mediated

Bias → Prejudice → Discrimination

↳ preconceived judgements based on inadequate knowledge
discr. → action that treats someone exclusively based on group

Fairness - given everyone equal consideration and opportunity to benefit in relation to some good or value

VI. Non-Linear Discriminative Classifiers

Decision Tree
 split feature space one feature at a time, recursively
 → forms tree structure
 → partition space in "rectangles"
 each rectangle/leaf → value

$$h(x) = \sum_l c_l \mathbb{1}(x \in L)$$

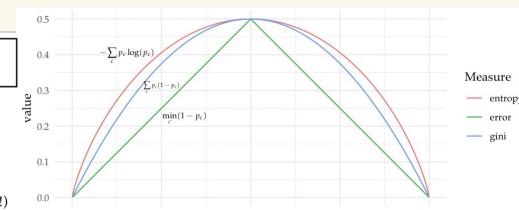
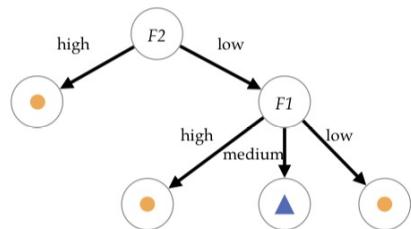
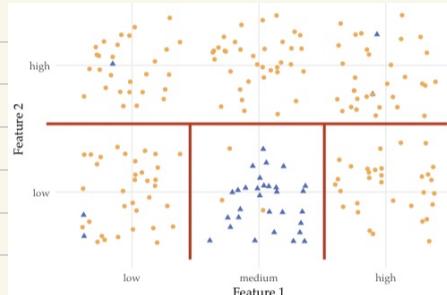
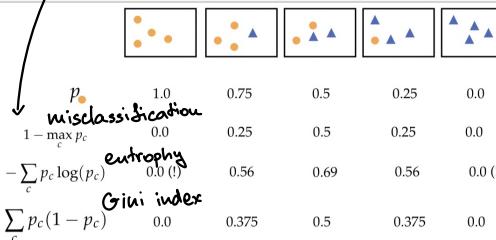
(L Leaves)

$$\min_{\mathcal{L}} \sum_{i=1}^N \mathbb{1}(h(x_i) \neq y_i)$$

objects in node

$$G(S) = - \sum_{V \in \text{values}(S)} \frac{|S_V|}{|S|} \mathbb{1}(S_V)$$

objects value V
feature



Note: Entropy scaled to have the same maximum

Information Gain

Overall:

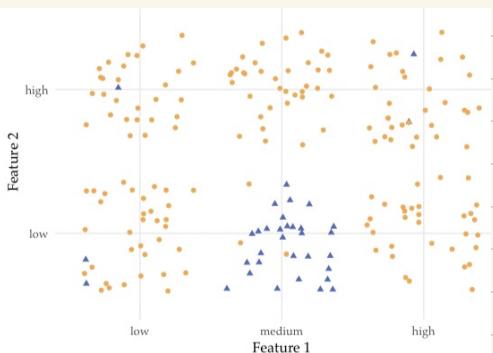
35/200

Feature 1:

Low (3/65), Medium (30/69), High (2/66)

Feature 2:

High (3/99), Low (32/101)



	Node Entropy	Information Gain
Overall: 35/200	$E(S) = -\sum_i p_i \log(p_i)$ $-(\frac{35}{200} \log(\frac{35}{200}) + \frac{165}{200} \log(\frac{165}{200})) = 0.464$	
Feature 1:		
Low (3/65)	$-(\frac{3}{65} \log(\frac{3}{65}) + \frac{62}{65} \log(\frac{62}{65})) = 0.187$	
Medium (30/69)	$-(\frac{30}{69} \log(\frac{30}{69}) + \frac{39}{69} \log(\frac{39}{69})) = 0.685$	$0.464 - (\frac{65}{200} \cdot 0.187 + \frac{69}{200} \cdot 0.685 + \frac{66}{200} \cdot 0.136) = 0.122$
High (2/66)	$-(\frac{2}{66} \log(\frac{2}{66}) + \frac{64}{66} \log(\frac{64}{66})) = 0.136$	
Feature 2:		
High (3/99)	$-(\frac{3}{99} \log(\frac{3}{99}) + \frac{96}{99} \log(\frac{96}{99})) = 0.136$	
Low (32/101)	$-(\frac{32}{101} \log(\frac{32}{101}) + \frac{69}{101} \log(\frac{69}{101})) = 0.624$	$0.464 - (\frac{99}{200} \cdot 0.136 + \frac{101}{200} \cdot 0.624) = 0.081$

$$G(S) = S(S) - \sum_{v \in \text{values}(F)} \frac{|S_v|}{|S|} S(S_v)$$

$$SV(S) = - \sum_{v \in \text{values}(F)} \log \frac{|S_v|}{|S|}$$

$$\text{GainRatio}(S) = \frac{G(S)}{SV(S)}$$

Pruning

$$C_2(\text{Tree}) = \sum_{\text{leaves}} |S_i| G(S_i) + \alpha |\text{leaves}|$$

- Unstable
- cannot model linear relationships efficiently
- greedy

Criteria

- min node size
- maximum tree depth
- minimum information gain

Regression

- squared loss function
- least value - average outcome
- splits minimize variance

Advantages

- "interpretable"
- automatic feature selection
- easy to incorporate discrete features and missing values
- fast

Combining Classifiers

Condorcet's Jury Theorem

If $p > 0.5$, adding voters to the jury increases probability of correct majority vote (assuming voting is independent)

Fixed Rules

hard (majority voting)

soft (mean product)

$$\max_y \prod_{j=1}^c h_j(y|x)$$

$$\max_y \prod_{j=1}^c h_j(y|x)$$

$$\min_{p(y|x)} \sum_j D(h_j(y|x), p(y|x))$$

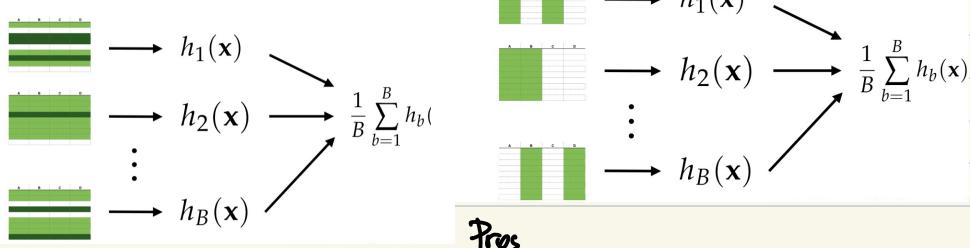
Learned Rules

decision based on set of base classifiers

Random Forest - large number of trees using randomly selected objects and features combined
bagging + random subspaces

Random Subspaces

Bagging (Bootstrap Aggregation)



Tree

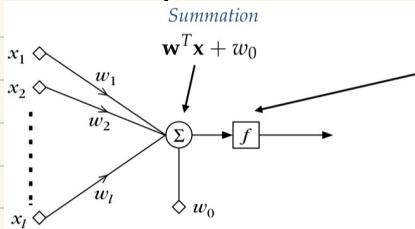
- depth
- # leaves
- min # object in node
- information criterion
- pruning

Forest

- # trees
- size of subspaces

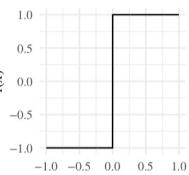
- Pros**
- flexible / low bias
 - many different types of data
 - embarrassingly parallel
 - out-of-bag (OOB) estimates
 - (scale) invariant
 - (relatively) few hyperparameters
- Cons**
- harder to implement than single trees
 - computationally expensive

Perceptron



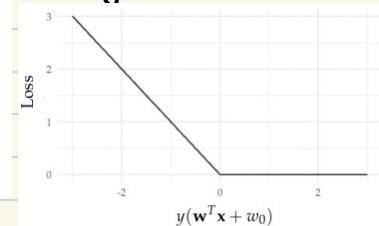
Activation function

Perceptrons use a step function
 $f(x) = 2(\mathbb{I}(x > 0) - 0.5)$



$$\sum_{i=1}^N \max(-y_i(\mathbf{w}^T \mathbf{x}_i + w_0), 0)$$

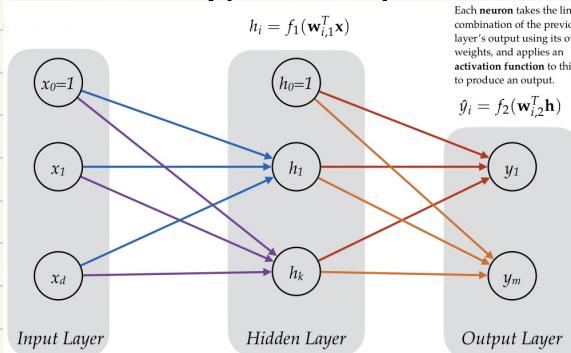
$$y_i \in \{-1, 1\}$$



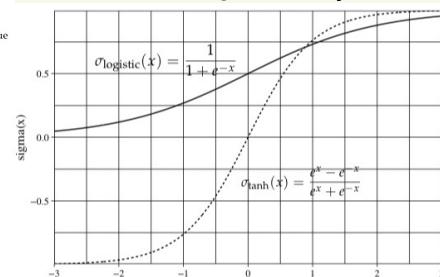
$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = \sum_{i=1}^N -y_i \mathbf{x}_i \quad y_i(\mathbf{w}^T \mathbf{x}_i) < 0$$

$$\mathbf{w}^{t+1} = \mathbf{w}^t + \alpha_t \sum_{i: y_i(\mathbf{w}^T \mathbf{x}_i) < 0} y_i \mathbf{x}_i$$

Multi-layer Perceptron



Activation Functions

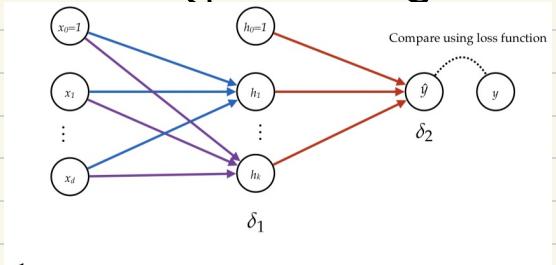


$$J(\mathbf{w}) = \sum_{i=1}^N L(f_2(\mathbf{w}_2^T f_1(\mathbf{w}_1^T \mathbf{x}_i)), y_i)$$

$$\mathbf{x} \rightarrow \mathbf{a}_1 = \mathbf{w}_1^T \mathbf{x} \rightarrow \mathbf{h} = f_1(\mathbf{a}_1) \rightarrow \mathbf{a}_2 = \mathbf{w}_2^T \mathbf{h} \rightarrow \hat{\mathbf{y}} = f_2(\mathbf{a}_2) \rightarrow L(\hat{\mathbf{y}}, \mathbf{y})$$

Back Propagation
 $\mathbf{w}^{t+1} = \mathbf{w}^t - \alpha_t \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^t}$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \begin{bmatrix} \frac{\partial J(\mathbf{w})}{\partial w_1} \\ \vdots \\ \frac{\partial J(\mathbf{w})}{\partial w_n} \end{bmatrix}$$



$$J(\omega_2, \omega_1) = L(f_2(\omega_2^T f_1(\omega_1^T x)), y)$$

$$(\nabla_{\omega_2} J(\omega_2, \omega_1))^T = \underbrace{\frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i}}_{\delta_2} \cdot \frac{\partial f_2(\omega_2^T f(\omega_1^T x))}{\partial \omega_2^T f(\omega_1^T x)} \cdot (f_1(\omega_1^T x))^T$$

$$(\nabla_{\omega_1} J(\omega_2, \omega_1))^T = \underbrace{\frac{\partial L(y_i, \hat{y}_i)}{\partial \hat{y}_i}}_{\delta_1 = \delta_2 \cdot \omega_2^T \cdot \delta_2} \cdot \underbrace{\frac{\partial f_2(\omega_2^T f(\omega_1^T x))}{\partial \omega_2^T f(\omega_1^T x)}}_{\delta_2} \cdot \omega_2^T \cdot \underbrace{\frac{\partial f_1(\omega_1^T x)}{\partial \omega_1^T x}}_{\delta_1}$$

$\begin{bmatrix} x^T \\ \sigma^T \\ \dots \\ 0^T \end{bmatrix}$

Advantages

- Flexible model class
- Good empirical performance on many (structured) problems
- Easily adapted to different learning settings

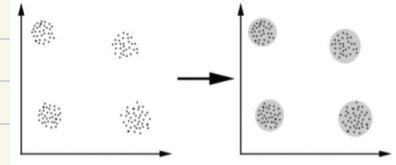
Disadvantages

- Computationally expensive
- Lots of hyperparameters
- No convergence to a unique optimum
- Hard optimization
- Hard to interpret

4. Unsupervised Learning

Clustering
natural groups in data:

- items within the group are close
- items between groups are far away



Proximity Measure

→ Similarity measure $s(x_i, x_k) \uparrow x_i \sim x_k$

→ Dissimilarity measure $d(x_i, x_k) \downarrow x_i \sim x_k$

→ Cosine similarity $\text{cos}(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$

→ Pearson's correlation coefficient $\text{Pearson}(x, y) = \frac{(x - \mu_x)^T (y - \mu_y)}{\|x - \mu_x\| \|y - \mu_y\|}$

Clustering Evaluation

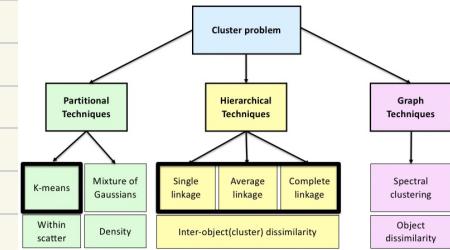
Intra-cluster cohesion (compactness)

→ how close data points are to cluster's mean

→ Sum of squared errors (SSE) vs.

Inter-cluster separation (isolation)

→ how far away are different cluster's means



K-means

n data points: $\{x_1, x_2, \dots, x_n\}$

p-dimensional feature vector: $x_i = \{x_{i1}, x_{i2}, \dots, x_{ip}\}$

k clusters

1. k random data points (seeds) for centroids (cluster's mean)
2. Each data point to closest centroid
3. Re-compute centroids using current cluster membership
4. Convergence criterion not met, repeat 2 & 3
 - minimum re-assignments of data points to different clusters
 - minimum change at centroids
 - minimum decrease in the sum of squared errors (SSE)

$$J(c, \mu) = \frac{1}{n} \sum_{i=1}^n \|x_i - \mu_c\|^2$$

Initial - run 1000 random, get lowest wss

Advantages

→ fast + simple

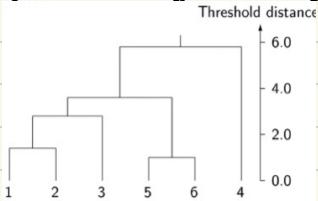
Disadvantages

→ only convex (round) clusters

→ sensitive to initialization

→ can get stuck in local minima

Hierarchical clustering
cluster algorithm can't pick ↴



top level - all points

bottom level - one cluster per data point

Divisive (top-down)

all points → 1 cluster

split cluster (sensibly)

↳ apply k-means recursively

run k-mean for $k=2$

for each cluster repeat)

Agglomerative (bottom-up)

each point → cluster

group two closest

sequence of clusters (inner size)

each level, merge two clusters

Merging rules

Single linkage - two nearest objects in clusters

$$g(R, S) = \min_{i,j} \{ d(x_i, x_j) : x_i \in R, x_j \in S \}$$

Complete linkage - two most remote objects in clusters

$$g(R, S) = \max_{i,j} \{ d(x_i, x_j) : x_i \in R, x_j \in S \}$$

Average linkage - cluster centers

$$g(R, S) = \frac{1}{|R||S|} \sum_{i,j} \{ d(x_i, x_j) : x_i \in R, x_j \in S \}$$

Pros

Cons

→ computationally expensive

→ clustering → "hierarchical nesting"

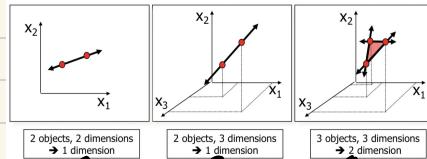
→ dendrogram - overview

→ linkage type - various shapes

→ different dissimilarities

Dimensionality Reduction

- storage / processing data - easier
- (visual) discovery of hidden structure
- redundant and noisy features
- smaller intrinsic dimensionality



Feature selection

subset of original dimensions $\{x_1, x_2, \dots, x_d\}$

domain knowledge + statistics-based selection methods

Feature extraction

new set of dimensions $E_i = f(x_1, \dots, x_d)$

(linear) combinations of original

Principal Component Analysis

linear subspace → maximized variance

principal components:

1 - direction of greatest variability

2 - perpendicular to 1 + greatest variability

d - original dimensionality

First m components → m new dimensions

$$\max_{\|\omega\|=1} \text{var}(\omega^T x) \quad \text{var}(x) = \frac{1}{N} \sum_{n=1}^N \left(x_n - \frac{1}{N} \sum_{n=1}^N x_n \right)^2$$

$$\text{mean}(x) \quad \text{covariance matrix } M = \begin{bmatrix} E[(x_1 - \mu_1)(x_1 - \mu_1)] & \dots & E[(x_1 - \mu_1)(x_d - \mu_d)] \\ \vdots & \ddots & \vdots \\ E[(x_d - \mu_d)(x_1 - \mu_1)] & \dots & E[(x_d - \mu_d)(x_d - \mu_d)] \end{bmatrix}$$

$$\text{zero-mean: } M = \frac{1}{n} X X^T$$

$$(\Rightarrow \text{assume: } \text{var}(\omega^T X) = [\omega^T X X^T \omega] = \omega^T M \omega)$$

$$\text{Lagrange multiplier: } \max_{\|\omega\|=1} \text{var}(\omega^T X) = \max_{\omega, \lambda} \omega^T M \omega - \lambda (1 - \omega^T \omega)$$

$$\text{gradient respect to } \omega: M \omega - \lambda \omega = 0 \Rightarrow M \omega = \lambda \omega$$

$M = n^{-1} \sum_{i=1}^n x_i x_i^T$ - M-square matrix, λ - constraint, e - non-zero column vector
 λ - eigenvalue of M, e - eigenvector of M for eigenvalue λ

1st principal component \rightarrow eigenvector with highest eigenvalue

$$\lambda e = \lambda e$$

$$(M - \lambda I)e = 0$$

$$\Rightarrow \det(M - \lambda I) = 0$$

\hookrightarrow n-degree polynomial \Rightarrow n eigenvalues of M

unit vector x_0

$$x_{k+1} := \frac{Mx_k}{\|Mx_k\|}$$

$$\text{converged: } \lambda_1 = x^T M x$$

$$\text{second pair: } M^* = M - \lambda_1 x x^T$$

PCA:

X matrix - (zero-mean) points in Euclidean space

covariance $X X^T$ + eigenpairs

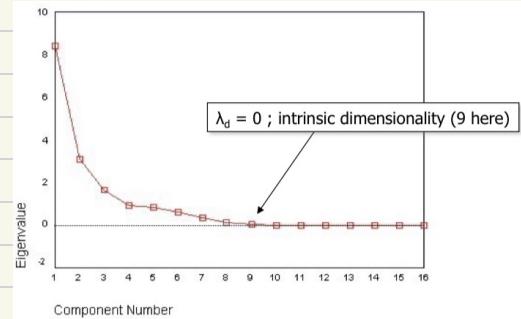
E matrix - eigenvector columns, largest eigenvalue first

$X^T E \rightarrow$ new coordinate space

1st axis (largest eigenvalue) - most significant

E_k - first k columns of E

$X^T E_k$ - k-dimensional X



Issues

Covariance - sensitive to large values

\rightarrow normalize dimension to zero mean $x' = \frac{x - \mu}{\sigma}$

Pros

\rightarrow reflects intuition about data

\rightarrow reduction in data size

Cons

\rightarrow expensive for many applications

\rightarrow assumptions behind methods