

CLASSIFIERS

$\{x : f(x) = 0\}$; isosurface for 0

Decision boundary: $(d-1)$ dim surface

linear C: decision b. line/plane

LINALG

unit vec: $\frac{x}{\|x\|}$

$$f(x) = w \cdot x + \alpha, H = \{w^T x = -\alpha\}$$

$$\vec{w} \perp H, w^T(y-x) = 0, y \in H$$

\vec{w} unit vector $\Rightarrow f(x)$ signed dist from x to H

α : distance from origin

CENTROID METHOD: mean of μ_c, μ_{not-c}

$$\hookrightarrow f(x) = (\mu_c - \mu_{not-c}) \cdot x - (\mu_c - \mu_{not-c}) \cdot \frac{\mu_c + \mu_{not-c}}{2}$$

PERCEPTRON METHOD: for linearly separable

\hookrightarrow uses gradient descent $O(r^2/y^2)$ iter
 $y_i = \begin{cases} 1 & \text{in } c \\ -1 & \text{not in } c \end{cases}$
 $\max \text{ margin } \frac{w^T x_i + \alpha}{\|w\|}$

$x_i \cdot w \geq 0, y_i = 1$
 $x_i \cdot w \leq 0, y_i = -1$
 $y_i \cdot x_i \cdot w \geq 0$

LOSS FUNCTION: $L(z, y_i) = \begin{cases} 0 & z \geq 0 \\ -y_i z & \text{otherwise} \end{cases}$

$$\text{RISK: } \frac{1}{n} \sum L(x_i \cdot w, y_i) = \frac{1}{n} \sum -y_i x_i \cdot w$$

$$\nabla R(w) = -\sum y_i x_i$$

while some $y_i x_i \cdot w < 0$, $w \leftarrow w + \epsilon y_i x_i$

Fictitious dim: Add α to w , 1 to x

SVM

margin: dist from decision b to nearest training point

$$\text{dist from } x_i: \frac{w \cdot x_i + \alpha}{\|w\|}$$

$$\text{margin: } \min \frac{1}{\|w\|} |w \cdot x_i + \alpha| \geq \frac{1}{\|w\|}$$

HARD MARGIN:

$$\min \|w\|^2 \text{ (smooth func)} \\ \text{s.t. } y_i(x_i \cdot w + \alpha) \geq 1$$

$$\text{Optimal: margin} = \frac{1}{\|w\|}$$

$$\text{Slab: } \frac{2}{\|w\|}$$

linearly separable, outliers \therefore

SOFT MARGIN:

$$\min \|w\|^2 + C \sum \xi_i$$

$$\text{s.t. } y_i(x_i \cdot w + \alpha) \geq 1 - \xi_i \\ \xi_i \geq 0$$

small C: max margin $\frac{1}{\|w\|}$, min $\|w\|$, underfit
 less sensitive, more flat

big C: $\xi_i \rightarrow 0$, small margin, overfit
 very sensitive, approach hard margin

Parabolic Lifting Map: $\Phi(x) = \begin{bmatrix} x \\ \|x\|^2 \end{bmatrix}$

Ellipsoid: axis aligned, no cross terms

$$Ax_1^2 + Bx_2^2 + Cx_3^2 + Dx_1x_2 + \dots + Ix_3 + \alpha = 0$$

$$f(x) = [A \ B \ C \ \dots \ I] \cdot \Phi(x) + \alpha$$

quadratic \Leftrightarrow conic in 2D

Raising: linearly separable \checkmark wider boundary robust

Edge Detection: Find regions with high gradients

OPTIMIZATION

\hookrightarrow convex func: no min, 1 min, ∞ min

\hookrightarrow strictly convex: 1 global min

\hookrightarrow high ellipticity of contours ($\lambda_{\max} \gg \lambda_{\min}$), no good learning rate

\hookrightarrow convex polytope satisfies linear program (\exists non empty, linearly separable)

DECISION THEORY

Bayes Theorem:

$$P(Y=1|x) = \frac{P(x|Y=1)P(Y=1)}{P(x)}$$

$$R(x) = E[L(r(x), Y)]$$

$$= \sum_x [(L(r(x), 1)P(Y=1|x) + L(r(x), -1)P(Y=-1|x))] P(x)$$

$$r''(x) = \begin{cases} 1 & \text{if } L(-1, 1)P(Y=1|x) > L(1, -1)P(Y=-1|x) \\ -1 & \text{otherwise} \end{cases}$$

When L symmetric, pick class with biggest posterior

$$E[g(x)] = \int_{-\infty}^{\infty} g(x) f(x) dx \quad \text{Var}(x) = E[(x - E(x))^2]$$

$$\mu = E[x] = \int_{-\infty}^{\infty} x f(x) dx \quad \sigma^2 = E[(x - \mu)^2] = E[x^2] - \mu^2$$

Bayes Decision Boundary:

$$P(x|Y=1)P(Y=1) = P(x|Y=0)P(Y=0)$$

Bayes Risk: Area under min of functions

GAUSSIAN DIS. ANALYSIS

$$X \sim N(\mu, \Sigma^2): f(x) = \frac{1}{(\sqrt{2\pi} \cdot \sigma)^d} \exp\left(-\frac{\|x - \mu\|^2}{2\sigma^2}\right)$$

Bayes maximises $f(x|Y=c)\pi_c$

Equivalently, $\pi_c(x) = -\frac{\|x - \mu\|^2}{2\sigma^2} - d \ln \pi_c + \ln \pi_c$
 Models $P(Y=c|x)$ as a logistic

$$QDA \quad \text{SAMPLE variance} \quad S = \frac{1}{1+e^{-x}}$$

\hookrightarrow Pick class to maximise π_c

\hookrightarrow 2 class: $\pi_c(x) - \pi_d(x) = 0$ (DB)

$$P(Y=c|x) = \frac{f_{x|Y=c} \pi_c}{f_{x|Y=c} \pi_c + f_{x|Y=d} \pi_d} = s(\pi_c - \pi_d)$$

LDA

All gaussians have same variance

$$\max \frac{\mu_c - \mu}{\sigma^2} - \frac{\|\mu_c - \mu\|^2}{2\sigma^2} + \ln \pi_c$$

2 class: DB $\Leftrightarrow w \cdot \mu + \alpha = 0, P(Y=c|x) = s(\text{DB})$

MLE

$$\hat{\mu} = \frac{1}{n} \sum x_i, \quad \hat{\sigma}_2 = \frac{1}{n} \sum \|x_i - \hat{\mu}\|^2$$

LDA: Same means and priors, one variance

$$\hat{\sigma}_2 = \frac{1}{n} \sum_c \sum \|x_i - \hat{\mu}_c\|^2$$

QDA: Same mean $\hat{\mu}$, calc $\hat{\sigma}_c^2, \tilde{\pi}_c = \frac{n_c}{n}$

$$L(\mu, \sigma^2; x_1, \dots, x_n) = f(x_1) f(x_2) \dots f(x_n) \Sigma^{-n}$$

$$L(\dots) = \ln f(x_1) + \dots + \ln f(x_n)$$

$$= \sum \left(-\frac{\|x_i - \mu\|^2}{2\sigma^2} - d \ln \sqrt{2\pi} - d \ln \sigma \right)$$

EIG

$$\vec{v}_i, \lambda_i \text{ for } A \rightarrow \vec{v}_i, \lambda_i^k \text{ for } A^k$$

$$\rightarrow \vec{v}_i, \frac{1}{\lambda_i} \text{ for } A^{-1}$$

Spectral Theorem: Real, symmetric M has real λ_i , n eigenvectors which are mutually orthogonal

$$M = V \Lambda V^T \quad (\text{UNIT eigenvectors})$$

PSD:

$$x^T A x \geq 0, \lambda_i \geq 0, A = U \Lambda U^T$$

$A^{1/2}$ exists

VISUAL

$x^T A^{-1} x = \|A^{-1/2} x\|_2^2$ is an ellipsoid
 axes $\vec{v}_i \Lambda \vec{v}_i^T$, radii $\lambda_i^{-1/2} \Lambda$

A diagonal: axis aligned

MULTIV. GAUSSIAN

$$X \sim N(\mu, \Sigma)$$

$$f(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

Σ : covariance, Σ^{-1} : precision

Isocountours of $(x - \mu)^T \Sigma^{-1} (x - \mu)$
 determined by eigenpairs of $\Sigma^{-1/2}$

Covariance:

$$\text{Cov}(R, S) = E[(R - E[R])(S - E[S])^T] = E[RS^T] - \mu_R \mu_S^T$$

$\text{Var}(R)$: PSD matrix

R_i, R_j ind $\Rightarrow \text{Cov}(R_i, R_j) = 0$

$\text{Cov}(R_i, R_j) = 0$ and multivariate norm $\Rightarrow R_i, R_j$ independent

All features ind $\Rightarrow \text{Var}(R)$ diagonal

$\text{Var}(R)$ diagonal + joint normal

$$\Leftrightarrow f(x) = f(x_1) \dots f(x_n)$$

\Leftrightarrow ellipsoids axis aligned,
 squared radii on Σ diagonal

ANISOTROPIC

$$\text{LDA: } \frac{1}{n} \sum_c \sum_{y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

$$\text{QDA: } \frac{1}{n} \sum_{y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

$\hookrightarrow \pi_c - \pi_d$ is quadratic

\hookrightarrow DB quadric

$$-\frac{1}{2} (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c) - \frac{1}{2} \ln |\Sigma_c| + \ln \pi_c$$

Multi-LDA:

$$\max \mu_c^T \Sigma^{-1} x + \frac{\mu_c^T \Sigma^{-1} \mu_c}{2} + \ln \pi_c$$

LDA underfits, QDA overfits

$$\dot{x} = x - \mu$$

$$\text{Var}(x) = \frac{1}{n} \dot{x}^T \dot{x}$$

$$Z = V^{-1} X$$

Decorrelate X : $Z = X V$, $\text{Var}(R) = V \Lambda V^T$

$$X \sim N(\mu, \Sigma) \quad Z \sim N(0, \Lambda)$$

Whitening: center + spherizing
Spherizing $X: W = X \text{Var } R^{-1/2}$

REGRESSION

1. Linear: $w \cdot x + \alpha$
2. poly
3. logistic: $s(w \cdot x + \alpha)$

Loss functions

- A. $(z - y)^2$ squared
- B. $|z - y|$ absolute
- C. $-y \ln z - (1-y) \ln(1-z)$ logistic (cables)

Cost functions

- a. $\frac{1}{n} \sum L(h(x_i), y_i)$ mean
- b. $\max L(h(x_i), y_i)$ max
- c. $\sum w_i L(h(x_i), y_i)$ weighted
- d. $a/b/c + \lambda \|w\|^2$ L_2
- e. $a/b/c + \lambda \|w\|_1$, L_1

LEAST SQUARES

$$\min \sum_i (x_i \cdot w + \alpha - y_i)^2$$

$$\min \|Xw - y\|^2$$

$$\nabla_{\text{RSS}}: X^T X w = X^T y. \quad X^T X \text{ always PSD}$$

$$w = \underbrace{(X^T X)^{-1} X^T}_{\text{pseudo-inverse}} y$$

$$g = Xw = X^T y = Hy. \quad \text{Ideal } H = I$$

Unique when $X^T X$ singular ✓

Sensitive to outliers

LOGISTIC REGRESSION

$$J = \sum_i L(s(x_i \cdot w), y_i)$$

$$= -\sum_i (y_i \ln s_i + (1-y_i) \ln(1-s_i))$$

$$s'(x) = s(x)(1-s(x))$$

$$\nabla_w J = -X^T(y - s) \quad \text{PSD}$$

$$\nabla_w^2 J = X^T \cdot 2X \cdot \underbrace{\sqrt{2}}_{\text{diag}(s_i(1-s_i))}$$

$$\text{BGD: } w \leftarrow w + \epsilon X^T(y - s(Xw))$$

$$\text{SGD: } w \leftarrow w + \epsilon (y_i - s_i(x_i \cdot w)) x_i;$$

separates linearly separable points

WEIGHTED LEAST SQUARES

$$\min (Xw - y)^T \Sigma^{-1} (Xw - y)$$

$$\Rightarrow X^T \Sigma^{-1} X w = X^T \Sigma^{-1} y$$

NEWTON'S METHOD

Iterative optimisation method for smooth $J(w)$

Works when $J(w)$ is quadratic

$$w \leftarrow w - (\nabla^2 J)^{-1}(\nabla J)$$

$$w \leftarrow w + \epsilon, \quad (X^T \Sigma X) \epsilon = X^T(y - s)$$

small weight in $W \Rightarrow$ more emphasis

faster than SGD, esp when ϵ small

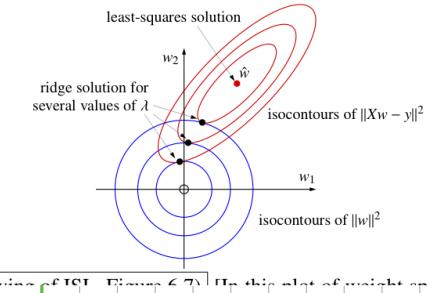
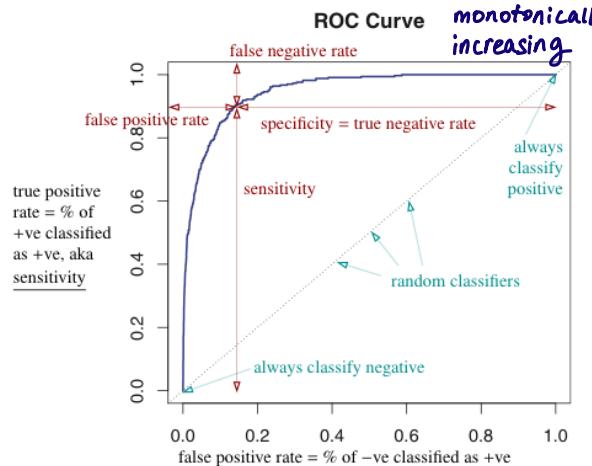
NOT on perceptron

$$\nabla_x \bar{y} = (\nabla_x \alpha) \bar{y} + \alpha \nabla_x \bar{y}$$

$$\nabla_x f(y) = (\nabla_x \bar{y})(\nabla_{\bar{y}} f(y))$$

$$\nabla_x (\bar{y} \cdot \bar{z}) = (\nabla_x \bar{y}) \bar{z} + (\nabla_x \bar{z}) \bar{y}$$

$$\nabla_x g(\bar{y}) = \nabla_x \bar{y} \quad \nabla_{\bar{y}} g(\bar{y})$$



Bayes classifier has nothing to do w sample points, needs to know abt distribution

Ellipse eqn requires quad coeffs

$$P(x \geq t) \leq \frac{E[x]}{t}$$

$$\Sigma = E[(z - \mu)(z - \mu)^T]$$

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

$$R(r(n)) = i | n$$

$$= \sum L(r(n)=i, y=j) P(Y=j | n)$$

$$\text{bias} = E[\hat{s}] - \theta$$

Knob: posterior prob threshold for GDA/Log Reg
Based on real test data

BIAS - VARIANCE

Bias: Inability of h fitting g

Variance: Random noise in data

Overfitting: Bias \downarrow Variance \uparrow

Underfitting: Bias \uparrow Variance \downarrow

+ Feature: Variance \uparrow Bias $= \downarrow$

$$E[L(h(z)), y] = E[(h(z) - y)^2]$$

$$= E[h(z)^2] + E[y^2] - 2E[yh(z)]$$

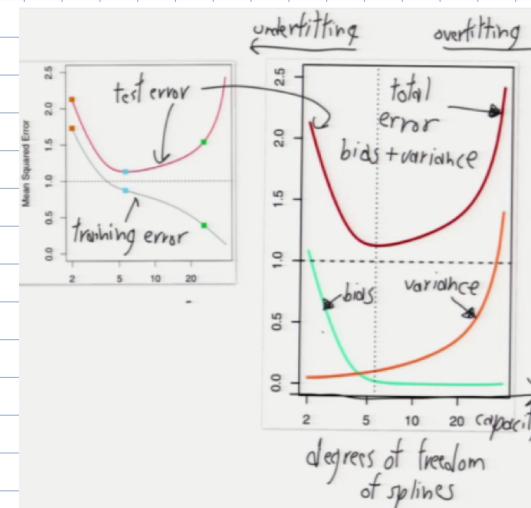
$$= \text{Var}(h(z)) + E[h(z)]^2 + \text{Var}(y) + E[y]^2 - 2E[yh(z)]$$

$$= (E[h(z)] - g(z))^2 + \text{Var}(h(z)) + \text{Var}(y)$$

bias

var

noise



Variance more harmful to test error
hyperparameters (λ, C, d) affect b-v

RIDGE

$$\min \|Xw - y\|^2 + \lambda \|w\|^2$$

ALWAYS unique optimum, ONLY one

$\lambda \rightarrow \infty, \text{ var} \rightarrow 0, \text{ bias} \uparrow, w^* \rightarrow 0$

$$(X^T X + \lambda I) w = X^T y$$

PD

Bayesian Justification: $w \sim N(0, c^2)$

$$f(w) \propto e^{-\|w\|^2 / 2c^2}$$

$$\max \ln f(y|x) + \ln f(w) - \text{const}$$

$$= -\text{const} \|Xw - y\|^2 - \text{const} \|w\|^2$$

$$= \min \|Xw - y\|^2 + \lambda \|w\|^2$$

Gaussian w $\propto cI$ prior

LASSO

$$\min \|Xw - y\|^2 + \lambda \|w\|, \quad \text{quadratic prog sparsity (intersects w axes)}$$

FEATURE

feature \uparrow , variance \uparrow

FORWARD: Start with 0 features, keep few :: adding until val error \uparrow
won't find best 2-f model

BACKWARD: Start with d, remove if more :: removing \downarrow val error

Train data \downarrow , Train Accuracy \uparrow

DECISION TREES

- ↳ classification + regression
- ↳ Tree with 2 nodes (internal nodes / feature + leaf)

$$\text{Entropy} : H(S) = -\sum_c p_c \log_2 p_c \quad \text{all same} \Rightarrow 0$$

$$\text{Weighted Avg.} : H_a = |S_l|H(S_l) + |S_r|H(S_r) \quad n \text{ dif classes} = \log_2 n$$

$$\text{Entropy} : H_a = |S_l| + |S_r|$$

Info gain : $H(S) - H_a$ Choose split to MAX info gain

Test : $\leq D(d)$, usually $\leq O(\log N)$

Train : $\leq O(nd)$ fast, bias \downarrow , translation invariant, robust to irrelevant features

VARIATIONS ↗ variance

Regression : Average sample points at each node

↳ Piecewise constant regression fn

↳ Choose split to MIN weighted avg. Variance

Stopping Early : ↗ limits tree depth, tree size

↳ Stops pure tree overfitting, good for overlapping

↳ Return posterior prob. at each leaf dist b/w classification(majority) / regression(coverage)

Pruning : Grow big tree, kill split if validation ↘

↳ More reliable than stopping early

Multivariate Split : Non-axis aligned splits w/ other algs. Classification ↗ Building Speed ↘

Ensemble Learning : Reduces variance in DTs

↳ Bagging : Same learning alg. on random

subsamples (WITH replacement), T learners

↳ Random Forests : Bagging + feature subset m

$m \downarrow \Rightarrow$ randomness T, bias T, tree correlation ↘

KERNEL $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ d input features

Degree p-polynomials blow up to $O(p^d)$ features

In many algs., $w = X^T a = \sum_i a_i x_i$ (linear combo of pts)

and can use inner product of $\Phi(x)$ only.

Optimize in dual weights a instead of primal w

Solve $n \times n$ linear system instead of $d \times d$

Kernel Ridge Regression :

$$(X^T X + \lambda I)^{-1} w = X^T y$$

$$(X^T X + \lambda I) a = y \quad [\text{works if you center } X]$$

Then $X^T y = (X^T X + \lambda I) X^T a$ [and y no $E(\text{bias}) \approx 0$]

$\Rightarrow w = X^T a$ is a solution.

$$\text{Obj: } \min_a \|X^T X a - y\|^2 + \lambda \|X^T a\|^2$$

Train: Solve $(X^T X + \lambda I) a = y$

$$\text{Test: } h(z) = w^T z = a^T X z = \sum a_i (x_i^T z)$$

Kernel fn: $k(x, z) = x^T z$. $K = X X^T (n \times n)$

$$K_{ij} = k(x_i, x_j)$$

$$\text{Alg: } K_{ij} \leftarrow k(x_i, x_j) \quad O(n^2 d)$$

$$\text{Solve } (K + \lambda I) a = y \quad O(n^3)$$

$$\text{Test: } h(z) = \sum a_i k(x_i, z) \quad O(nd)$$

Polynomial Kernel : Degree p in x

$$k(x, z) = (x^T z + 1)^p = \phi(x)^T \phi(z)$$

Compute $\phi(x)^T \phi(z)$ in $O(d^p)$ instead of $O(d^2)$

SGD dual: $a_i \leftarrow a_i + \epsilon (y_i - s(k_{ii}))$

Gaussian Kernel : generate feature vec in ∞ space

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^{\infty}. k(x, z) = \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right) O(d)$$

inner products of similar features large, small

for dissimilar features

Key: $h(z) = \sum a_i k(x_i, z)$ is linear combo

of Gaussians centered at training pts

Dual weight coefficients of linear combo

Why? Very smooth h , like K-NN but

smoother, oscillates less than polynomials

larger $\sigma \rightarrow$ wider gaussian \rightarrow bias \uparrow

smoother h variance \downarrow

NEURAL NETWORKS

- ↳ classification + regression
- ↳ Tree with 2 nodes (internal nodes / feature + leaf)

perceptrons failed @ XOR

$$S'(z) = S(z)(1 - S(z))$$

non-linear activation

$$S(z) = \frac{\# \text{points in } c}{\# \text{total pts}}$$

half C, half D = 1

all same = 0

Input: x_1, \dots, x_d

Hidden: h_1, \dots, h_k

Output: z_1, \dots, z_K

Training: SGD / Batch. Cost function NOT convex

Initialize random weights to break symmetry, $\frac{k=0}{\sqrt{m}}$

Backprop: $\Delta(\text{edges})$, dp alg

Vanishing Gradient: s close to 0 or 1, $s' \rightarrow 0$

Gradients change too slowly. Layers \uparrow Vanishing \uparrow

Sigmoid \rightarrow ReLU: max {0, y}, might explode

Output Unit

① Regression: linear output units + squared loss

② Classification: classes ≥ 2 , softmax units

$$z_i(t) = \frac{e^{t_i}}{\sum_{j=1}^K e^{t_j}}, \text{ cross entropy} = \sum_{i=1}^K y_i \ln z_i$$

linear / ReLU \rightarrow squared error

sigmoid / softmax \rightarrow cross entropy

OPTIMIZATIONS

Fix vanishing gradient \rightarrow SGD on large sets

Normalize training points: center, then scale

Use small rand. subsample to choose ϵ

Dif ϵ for dif layers of weights

Emphasize schemes (repeat rare exemplar)

Acceleration scheme for SGD

Add enough weights, test error \downarrow "double descent"

Heuristics:

① Weight decay: Add $\lambda \|w\|^2$ to loss

step $\Delta w_i = -\epsilon \frac{\partial J}{\partial w_i}$ has extra term $-2\lambda w_i$

② Ensemble of NNs or dropout in 1 NN (p=0.5)

CNN out = $1 + \lfloor \frac{\text{input} - \text{kernel height} + 2p}{\text{stride}} \rfloor$

① Local Connectivity \rightarrow hidden unit connects to

small patch from previous layer

② Shared Weights \rightarrow group of hidden units

share set of weights "Kernel / filter"

Convolution: same linear transf. applied to different

patches of input by shifting

NEAREST NEIGHBOURS

Given q , find K training points near q . dist metric up to you. Regression(avg label) or classification

(majority or histogram w probabilities)

nr \rightarrow underfitting | n $\rightarrow \infty$, 1-NN err $< 28 - B^2$

nr $\uparrow \rightarrow$ overfitting | 2 classes $\leq 28 - B^2$

draw pts independently from same prob. distribution.

n $\rightarrow \infty$, K $\rightarrow \infty$, K/n $\rightarrow 0$, err converge to B

Exhaustive K-NN Alg: Given q , scan through

all n points, computing squared distance to q. Maintain a max-heap with K shortest

distances. Query time $O(nd + n \log K)$

2-S d \rightarrow Voronoi

Medium d \rightarrow K-d trees

Large d \rightarrow exhaustive k-NN

can use PCA / rand. proj.

Voronoi: X point set. Vor w = {l1p.w1 \leq l1p.w2 $\leq \dots \leq$ l1p.wn}

voronoi cell always convex polyhedron / polytope

Size $\propto O(n^{d+1})$, often $O(n)$

For 2d: $O(n \log n)$ to compute V.d. + trapezoidal map for pt. location. $O(\log n)$ query time.

d-d: Use binary space partition tree

K-d Trees:

1. Choose splitting feature w greatest width. ft.

2. in $\max(x_{ji} - x_{ki})$, make boxes cubical

3. Choose splitting value, median pt. for feature

or $(x_{ji} + x_{ki})/2$. median guarantees $O(\log n)$

depth. $O(\log n)$ building time. Each internal

node stores a training pb.

Goal: Given q , find w so $\|l_1 p.w1 \leq \dots \leq l_1 p.wn\| \leq 1 + \epsilon$

Maintain: K-nearest so far, binary min-heap of stored

UNSUPERVISED LEARNING

Sample points, no labels \rightarrow discover structure in data

PCA

Dimensionality reduction: find K directions that capture most of the variation

Why: reducing # dimensions, remove irrelevant dimensions, find small basis for repres

$X: n \times d$, assume $\sum_i x_i = 0$, $\text{proj}_w x = \frac{\|w\|^2}{\|w\|^2} w$

$X^T X$: square, symmetric, PSD, $(v_i, v_i) = \text{corresponding unit principal components}$

① Fit Gaussian to data with MLE

Choose K Gaussian axes of greatest variance

MLE estimates $\hat{\Sigma} = \frac{1}{n} X^T X$

1. Center X, normalize X

2. Compute unit eigenvectors of $X^T X$

3. Choose K, choose v_{d-K+1}, \dots, v_d

4. Compute $\bar{x} \cdot \bar{v}_i$ for each training point

② Find direction \bar{w} that maximizes sample variance of projected data

$$\max_w \frac{1}{n} \sum (X_i \cdot \frac{w}{\|w\|})^2 = \frac{1}{n} \frac{\|Xw\|^2}{\|w\|^2} = \frac{1}{n} \frac{w^T X^T X w}{w^T w}$$

$\Rightarrow v_d$ achieves max variance λ_d / n

③ Find direction \bar{w} that minimizes mean squared projection distance

(proj orthogonal to proj hyperplane)

$$\begin{aligned} \min_w \sum_i \|X_i - \bar{X}_i\|^2 &= \sum_i \left\| X_i - \frac{\sum a_i x_i}{\|w\|^2} w \right\|^2 \\ &= \sum_i \left(\|x_i\|^2 - \left(\frac{x_i \cdot w}{\|w\|} \right)^2 \right) \\ &= \text{constant} - n (\text{Var } w) \end{aligned}$$

SVD

Problem: Computing $X^T X$ takes $\Theta(nd^2)$ time, eigenvalues

$X = UDV^T$. When $n \geq d$

$$X = U = \sum_i \delta_i u_i v_i^T$$

δ_i rank 1 outer product matrix

$$U^T U = I = V^T V \text{ (orthonormal)}$$

D diagonal

D: $\delta_1, \dots, \delta_d \rightarrow$ Nonnegative, singular values of X

Non-zero $\delta_i = \text{rank}(X)$

$$X^T X = V D V^T U D U^T = V D^2 V^T$$

v_i has eigenvalue δ_i^2

Find K greatest singular values in $O(nk)$ time

Since $X = UDV^T \Rightarrow XV = UD \Rightarrow \text{Row}_i(VD) = x_i \cdot v_i$

principal coordinates

CLUSTERING

Partition data into clusters so points in clusters are more similar than across clusters

Why? Discovery, hierarchy, graph partitioning, quantization

K-means clustering : partition n points into K disjoint clusters

Cluster i mean: $\mu_i = \frac{1}{n_i} \sum X_j$

$$\min \sum_y \sum_{j:y=j} \|X_j - \mu_i\|^2$$

squared distances from points to means

NP-hard, Solve in $O(nk^n)$ time

K-means heuristic:

1. Fixed y_j , Update μ_i
2. Fixed μ_i , Update y_j

Halt when step 2 changes no labels.

Both steps minimize cost, but not all vars.

Alg never returns to prev assignment. MUST terminate. Finds local minimum.

Starting:

- ↳ Forgy Method: Choose K rand points to be initial $\mu_i \rightarrow$ Go to Step 2 better
- ↳ Random Partition: Randomly assign pts to a cluster \rightarrow Go to Step 1

Equivalent: Min within-cluster variation

$$\min \sum_{y_i} \sum_{i=1}^k \sum_{j \in i} \|x_j - \bar{x}_i\|^2$$

K-Medoids Clustering:

- ↳ Specify distance fn $d(x, y)$
- ↳ Replace mean with medoid - pt that min total distance to other pts in same cluster
- ↳ Medoid ALWAYS input point, more robust to outliers

HIERARCHICAL CLUSTERING

K no longer a hyperparameter $O(n^3)$ time

1 Point Set \rightarrow Bottom-up 'Agglomerative'

Start w/ a point, keep fusing upwards

2 Graphs \rightarrow Top-down 'Divisive'

Start w/ everything in a cluster, keep splitting

Complete Linkage: $\max \{d(a, b)\}$

Single Linkage: $\min \{d(a, b)\}$

Average Linkage: $\frac{1}{|A||B|} \sum_{a \in A} \sum_{b \in B} d(a, b)$

Centroid Linkage: $d(\mu_A, \mu_B)$ euclidean distance

HIGH DIMENSIONS

↳ Intuition wrong in higher dim

↳ $\|p\|^2 = p_1^2 + p_2^2 + \dots + p_d^2$ in a dim

$p_i \sim N(0, 1)$ $p_i^2 \sim \chi^2(1)$ $E[p_i^2] = 1$ $\text{Var}(p_i^2) = 2$

$\|p\|^2 \sim \chi^2(d)$

↳ For large d $\|p\|^2$ is concentrated in thin shell around radius \sqrt{d} with thickness $\approx \sqrt{2d}$

↳ In high dim, 1-nn and 1000-nn don't differ much \Rightarrow K-means & nn not as effective

$\cos \theta = \frac{p \cdot q}{\|p\|\|q\|}$ $E[\cos \theta] = 0$, $\text{SD}[\cos \theta] = \frac{1}{\sqrt{d}}$

↳ d large \rightarrow $\cos \theta$ tends to 0 / θ tends to 90°

↳ d large \Rightarrow 2 rand vectors almost always \perp

Random Projection: Alt to PCA, preserves dist.

b/w points. Project very high-d space to rand medium-d subspace

$$K = \left[\frac{2 \ln(1/\delta)}{\epsilon^2/2 - \epsilon^3/3} \right]. \text{For some point } q,$$

q is orthogonal proj of q onto S , multiplied by $\sqrt{\frac{d}{K}}$

Choosing projection direction: choose each component from a univariate Gaussian dist. then normalise vector to unit length. K rand directions, then Gram-Schmidt orthogon.

PSEUDOINVERSE

D: diagonal nxd matrix

Transpose D \rightarrow replace nonzero entry w reciprocal $-D^{-1}$

$$DD^T D = D \quad D^T DD^T = D^T$$

SVD of $X = UDV^T$. rank D = rank X

Moore-Penrose Pseudoinv. of X : $X^+ = VD^+U^T$

$$1. XX^+ = UDV^T D^+ U^T = U(DD^T)U^T, \text{ PSD}$$

$$2. X^T X = VD^T U^T UDV^T = V(D^T D)V^T, \text{ PSD}$$

$$3. D, D^+, DD^T, D^T D, X, X^T, XX^T, X^T X, \text{ ranks equal}$$

$$4. X \text{ has rank } n, \quad XX^+ = I_n, \quad X^+ \text{ right inv}$$

$$5. X \text{ has rank } d, \quad X^T X = I_d, \quad X^+ \text{ left inv}$$

$$6. XX^T X = X [UDV^T D^+ U^T UDV^T = UDV^T = X]$$

$$7. X^T X X^+ = X^+ [\text{symmetric proof } T]$$

Theorem: $X^T X w = X^T y$; $w = X^T y$

$$\text{Proof: } X^T X w = X^T X^T y = VD^T U^T UDV^T U^T y = VD^T D^+ V^T y = VD^T y = X^T y$$

If normal eqns. have multiple solns, then $w = X^T y$ is the least norm soln (minimizes $\|w\|$ among all solutions)

↳ Useful when $X^T X$ is singular

LEARNING THEORY

Range Space / Set System: Pair (P, H) where

P is set of all possible test/training pts

H is set of hypotheses/classifiers

↳ $h \subseteq P$ that specifies which points h predicts are in class C

1. Power-set classifier: H power set of P containing $2^{|P|}$ subsets of P

2. Linear classifier: $P = \mathbb{R}^d$, H is set of all halfspaces $\{x: w \cdot x \geq -\alpha\}$ / every possible linear classifier in d dimensions

Power-set classifier sucks at learning

Risk (Generalization Error): $R(h)$ of h is probability that h misclassifies a random point x drawn from D

∞ test data \Rightarrow risk = test error

Empirical Risk (Training Error): $\hat{R}(h)$ is 1/- of X misclassified by H

$n \rightarrow \infty$, $\hat{R}(h)$ approximates $R(h)$

Hoeffding's Inequality:

$$Pr(\hat{R}(h) - R(h) > \epsilon) \leq 2e^{-2\epsilon^2 n}$$

Choose $\hat{h} \in H$ that minimizes $\hat{R}(\hat{h})$

Dichotomies: Dich of X is $X \cap h$, $h \in H$

For n points, 2^n dichotomies.

Dichotomy assigns points to class C or not C

If H allows all 2^n , $\hat{R}(\hat{h}) = 0$

H induces Π dichotomies.

$Pr(\text{at least one dich has } |\hat{R} - R| > \epsilon) \leq \delta$,

$$\delta = 2\Pi e^{-2\epsilon^2 n}$$

$$|\hat{R}(h) - R(h)| \leq \epsilon = \sqrt{\frac{1}{2n} \ln \frac{2\Pi}{\delta}}$$

Smaller Π , larger n \Rightarrow training error \cong true risk

smaller $\Pi \Rightarrow$ less likely to overfit (variance \downarrow)

Sample complexity: # training pts needed for prob(ϵ) \leq

$$N \geq \frac{1}{2\epsilon^2} \ln \frac{2\Pi}{\delta}$$

SHATTER FUNCTION

of dichotomies: $\Pi_H(X) = |\{X \cap h : h \in H\}|$, $n \in |X|$

Shatter function: $\Pi_H(n) = \max \Pi_H(X)$ $\Pi_H(2) = 8$ $\Pi_H(4) = 16$

For all range spaces, $\Pi_H(n)$ is polynomial in n, or $\Pi_H(n) = 2^n$

Cover's Theorem: linear classifiers in \mathbb{R}^d allow $\Pi_H(n) = 2 \sum_{i=0}^d \binom{n-1}{i}$ dichotomies of n points

$n \leq d+1$, $\Pi_H(n) = 2^n$ $n \geq d+1$, $\Pi_H(n) = 2 \binom{c(n-1)}{d}$

To achieve $R(h) \leq R(\hat{h}) + \epsilon \leq R(\hat{h}) + 2\epsilon$ w.p. $\geq 1-\delta$.

$n \geq \frac{1}{2\epsilon^2} (d \ln \frac{n-1}{\delta} + d + \ln \frac{n}{\delta})$. Linear classifiers only need $n \in \Theta(d)$ training points for train err \approx test err

VC DIM. $VC(H) = \max \{n : \Pi_H(n) = 2^n\}$

H shatters set X of n pts if $\Pi_H(X) = 2^n$. $VC(H)$ is largest X H can shatter. (H allows all 2^n dichotomies)

$\Pi_H(n) \leq \sum_{i=0}^n \binom{n}{i}$. $n \geq VC(H)$, $\Pi_H(n) \leq \left(\frac{en}{VC}\right)^n$

VC-dim upper bound on exponent of polynomial $O(VC(H))$ pts. suffice for accuracy

VC finite \Rightarrow sample complexity grows w. # features

VC $\infty \Rightarrow$ no amt. of training data will make it generalise well

$VC(H) = 3$. $\Pi_H(n) \leq \frac{e^3}{27} n^3$. $O(1)$ sample complexity

ADA-BOOST

ensemble, classif + regression reduces bias

Train T classifiers G_1, \dots, G_T correct $\cong w_t \downarrow$ wrong $\cong w_t \uparrow$

Weight for x_i in G_t grows according to how many classifiers misclassified it

Train G_t to try harder to correctly classify x_i with w_t

Metalearner is linear combo. of learners

Test point z, $M(z) = \sum B_t G_t(z)$. G_t is ± 1 , M continuous. Retain sign of M(z).

Iteration T: Find G_T, B_T so

$$\min \text{Risk} = \frac{1}{n} \sum_i l(M(x_i), y_i), \quad l(x_i) = \sum_{t=1}^T B_t G_t(x_i)$$

Metalearner uses exponential loss:

$$l(p, t) = e^{-p} \quad \begin{cases} e^{-p} & t=+1 \\ e^p & t=-1 \end{cases}$$

cost fn. for G_t : # misclassified points

$$n \cdot \text{RISK} = \sum_i e^{-y_i M(x_i)} = \sum_i \Pi e^{-B_t G_t(x_i)}$$

$$= e^{-B_t} \sum_i w_i + e^{B_t} \sum_i w_i$$

$$w_t^{T+1} = w_t \cdot e^{-B_t y_i} G_t(x_i) \quad \begin{cases} w_t e^{-B_t} & \text{correct} \\ w_t e^{B_t} & \text{misclassified} \end{cases}$$

$$B_T \text{ Set } \frac{\partial}{\partial B_T} (\text{RISK}) = 0$$

$$B_T = \frac{1}{Z} \ln \left(\frac{1 - \text{err}_T}{\text{err}_T} \right) \quad \begin{cases} \text{err}_T = 0, B_T = 0 & \text{(perfect)} \\ \text{err}_T = 1/2, B_T = 0 & \text{(no variance)} \end{cases}$$

Algorithm:

1. Initialize $w_1 \leftarrow 1/n$

2. For $t = 1$ to T. 1) Train G_t with weights w_t 2) calc. err

3) Reweight points $w_i \leftarrow w_i \cdot e^{B_t}$ or $w_i \cdot e^{-B_t}$

3. Return M $m(z) = \text{sign} \left(\sum B_t G_t(z) \right)$

Why?: Reduces bias reliably, no hyperparameter search needed, fast, short trees + AdaF = feature linear decision boundaries don't boost well

More: $P(Y=1|n) \approx 1/(1+e^{-2M(n)})$

Exponential loss vulnerable to outliers

If all learners beat error $< 50%$, train err = 0 wally