

Linear decision boundary is a hyperplane
select hyperparameters with validation
sample point = feature vector
decision fn: $f(x)$ that maps a
sample point x to a scalar
 $f(x) > 0$, if $x \in$ class c
 $f(x) \leq 0$, if $x \notin$ class c

decision boundary is $f(x) = 0$, $x \in \mathbb{R}^d$
 $\{x: f(x) = 0\} :=$ isosurface of f for
isovalue 0.

Euclidean norm: $\|x\| = \sqrt{x \cdot x} = \sqrt{x_1^2 + \dots + x_d^2}$
Normalize by: $\frac{x}{\|x\|}$

Hyperplane $\in \mathbb{R}^d$ has dim $d-1$.

$H = \{x: w \cdot x = -\alpha\}$, given linear
decision fn: $f(x) = w \cdot x + \alpha$
 w is normal vector of H .

Centroid! ~~doesn't require linear separability~~
 $\mu_c :=$ mean of all vectors $\in C$

$\mu_x :=$ mean of all $\in C$.

$f(x) = (\mu_c - \mu_x) \cdot x - (C\mu_c - \mu_x) \cdot \frac{\mu_c + \mu_x}{2}$
- decision boundary is hyperplane
that bisects line segment w
endpoints μ_c, μ_x

Perceptron

$y_i = \begin{cases} 1 & \text{if } x_i \in \text{class } C \\ -1 & \text{if } x_i \notin C \end{cases}$ $x_i \cdot w \geq 0, y_i = 1$
 $x_i \cdot w \leq 0, y_i = -1$

Goal: $y_i x_i \cdot w \geq 0$, i.e. $x_i \cdot w \geq 0, y_i = 1$
 $x_i \cdot w \leq 0, y_i = -1$

$LC(x_i, y_i) = \begin{cases} 0, & y_i x_i \geq 0 \\ -y_i x_i, & \text{otherwise} \end{cases} \Rightarrow$ want same sign

$R(w) = \sum_{i \in V} -y_i x_i \cdot w, V = \{i \text{ where } y_i x_i \cdot w < 0\}$

Goal: minimize risk $= \sum_{i \in V} LC(x_i, y_i, w)$

gradient of R wrt w is direction of steepest ascent
SGD: pick 1 misclassified x_i , do GD on $LC(x_i, y_i)$

simulate general hyperplane in d dim by using hyperplane
thru origin in $d+1$ dimensions.

Perceptron algorithm will find a soln if possible
Margin! distance from decision boundary to nearest
sample point

$w \cdot x + \alpha = 1$ $w \cdot x + \alpha = 0$ (decision boundary)
 \Rightarrow slab of width $\frac{2}{\|w\|}$ containing no sample points

Opt! Find w, α that minimize $\|w\|^2$ subject to
 $y_i(x_i \cdot w + \alpha) \geq 1$ where $y_i(x_i \cdot w + \alpha) = 1$

is a quadratic program "on margin"

This is a maximum margin classifier, aka hard
margin SVM w is \perp to decision boundary

Signed distance from hyperplane to x_i is $\frac{w}{\|w\|} \cdot x_i + \frac{\alpha}{\|w\|}$

Hard margin SVMs fail if data not linearly separable
sensitive to outliers

Soft Margin SVMs:
- allow some points to violate the margin, w slack
variables

$y_i(x_i \cdot w + \alpha) \geq 1 - \xi_i, \xi_i \geq 0$ constraints

Point i has nonzero ξ_i iff it violates the margin

Opt! minimize $\|w\|^2 + C \sum \xi_i$ C affects bias/variance tradeoff

large maximize margin keep misclassification var zero or small

large underfitting overfitting any point on slab/margin considered support vector

outliers less sensitive very sensitive

boundary more "flat" more "curved"

Nonlinear decision boundaries - nonlinear features that
lift points into higher d space.
- add features

- can lift points to d space to make them
linearly separable
- margin tends to get wider as degree increases
higher degree \rightarrow overfitting.

Edge detection: collect line orientations in local
histograms, use histograms as features

Applications/Data

Model ~~Problem~~

Optimization ~~Algorithm~~

Optimization Algorithm

Type of Optimization Problems:
- find w that minimizes/maximizes
a continuous objective fn $f(w)$

Finding global minimum is hard

For smooth f :
- Gradient descent
- blind
- w likelihood
- stochastic

nonlinear conjugate gradient

Newton's method

Non-smooth f :
- Gradient descent
- blind
- direct likelihood search

likelihood finds a local minimum
by solving an optimization problem in
1D

Constrained Optimization
- w that optimizes $f(w)$ subject to
 $g(w) = 0$, g is a smooth fn

\rightarrow Lagrange multipliers, transform
constrained to unconstrained
optimization.

Linear Program - $Aw \leq b$
- linear objective fn + linear
inequality constraints

Set of points that satisfy all
constraints is a convex polytope
called feasible region

optimum $\in F$ is point furthest
in direction c .

Active Constraints: constraints
for which optimum achieves
equality.

Set of optimal pts always convex

Linearly separable iff feasible
region isn't empty set

Hard figuring out active constraints
- involve more discrete than
continuous unconstrained

Quadratic Program
- quadratic, convex objective
fn + linear inequality constraints

$f(w) = w^T Q w + c^T w$, subject
to $Aw \leq b$

Q is PD \Rightarrow 1 local minimum

Convex Program
- convex objective fn + convex
inequality constraints

Posterior: $P(Y=y|X)$
Prior: $P(Y=y)$

Loss fn specifies badness
for each incorrect prediction.

$R(x) = E[L(C(x), y)]$

Bayes Decision Rule! - assumes no loss

$R(x) = \begin{cases} 1 & \text{if } (C(x)=1) \text{ and } P(Y=1|X) < 0.5 \\ 0 & \text{otherwise} \end{cases}$

When L is symmetric, pick class w
biggest posterior probability

PDF: $P(x)$

$E[f(x)] = \int_{-\infty}^{\infty} f(x) P(x) dx$

$\text{Var} = E[(X - \mu)^2] = E[X^2] - \mu^2$

$\mu = E[X] = \int_{-\infty}^{\infty} x P(x) dx$

Bayes decision rule

$P(X|Y=1)P(Y=1)$

$P(X|Y=1)P(Y=1)$

look up x , pick curve w highest
probability

Classifiers

1) generative models e.g. LDA
assume points come from probability
distr, different for each class

- fit distr param to class c
points, giving $P(X|Y=c)$
- estimate $P(Y=c)$
- gives $P(Y|X)$

0-1 loss - pick class c that maximizes
 $P(Y=c|X=x)$

\rightarrow maximizes $P(X=x|Y=c)P(Y=c)$

2) discriminative models e.g. logistic regression
model $P(Y|X)$ directly

3) Find decision boundary
- model $r(x)$ directly (no posterior)

1 & 2
 $P(Y|X)$ tells you probability your guess is
wrong

1
- can diagnose outliers: $P(X)$ very small

- hard to estimate distr accurately

Generative Model
- most popular when you have phenomena
well approximated by normal distr
+ lots of sample points

Gaussian Discriminant Analysis as logistic
each class comes from normal distr

$X \sim N(\mu, \sigma^2)$: $P(X) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu)^2}{2\sigma^2})$

Bayes decision rule returns class c that
maximizes $P(X=x|Y=c)P(Y=c)$

$Q_c(x) = R_c(C - \ln P_c(X)) =$
 $-\ln P_c(X) - \ln R_c$

$Q_c(x) = -\ln P_c(X) - \ln R_c$

QDA allows you to determine
probability that your classification
is correct

$P(Y=c|X=x) = \frac{S(Q_c(x) - Q_p(x))}{S}$
 $S = \frac{1}{1 + e^{-S}}$

LDA:
- linear decision boundaries, less likely
to overfit

Assume all Gaussians have same var σ^2
 $Q_c(x) - Q_p(x) = \frac{C(\mu_c - \mu_p) \cdot x}{\sigma^2}$

$-\frac{1}{2\sigma^2}(\mu_c^2 - \mu_p^2) + \frac{1}{\sigma^2}(\mu_c - \mu_p)x$

$= w \cdot x + \alpha$ } decision boundary

Choose class that maximizes
linear discriminant fn

$\frac{\mu_c \cdot x}{\sigma^2} - \frac{1}{2\sigma^2} \ln P_c$

MLE to use QDA
given x_1, \dots, x_n , find best fit
Gaussian

$L(\mu, \sigma; x_1, \dots, x_n) = P(x_1) \dots P(x_n)$

maximize log likelihood.

$R(\mu, \sigma) = \ln P(x_1) + \dots + \ln P(x_n)$

$= \sum_{i=1}^n \left(-\frac{1}{2\sigma^2} (x_i - \mu)^2 - \frac{1}{2} \ln \sqrt{2\pi} \right)$

$\frac{\partial R}{\partial \mu} = \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} = 0$

$\frac{\partial R}{\partial \sigma} = \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} - \frac{n}{\sigma} = 0$

Use mean + variance of points
in class c to estimate mean
+ variance of Gaussian for
class c .

QDA: separate conditional mean
+ variance for each class c

$\frac{\mu_c}{\sigma_c^2} = \frac{\sum_{i: Y=c} x_i}{\sum_{i: Y=c} 1}$

LDA: one variance for all classes
is within class

\rightarrow use each point's distance from its
class' mean

$\sigma^2 = \frac{1}{n} \sum_{i: Y=c} \|x_i - \mu_c\|^2$

$\hat{\mu}_c = \frac{1}{n} \sum_{i: Y=c} x_i$

$\hat{\sigma}^2 = \frac{1}{n} \sum_{i: Y=c} \|x_i - \hat{\mu}_c\|^2$

$\hat{\mu}_c = \frac{1}{n} \sum_{i: Y=c} x_i$

$\hat{\sigma}^2 = \frac{1}{n} \sum_{i: Y=c} \|x_i - \hat{\mu}_c\|^2$

$\hat{\mu}_c = \frac{1}{n} \sum_{i: Y=c} x_i$

Eigenvector $Av = \lambda v$

v points in same direction after
being multiplied by A

$A^k v = \lambda^k v$ $A^{-1} v = \frac{1}{\lambda} v$

Spectral Theorem: every real, symmetric
 $n \times n$ has real eigenvalues and n
eigvec that are mutually orthogonal

- if 2 eigvec have same eigenval,
every linear combination of those
eigvec are also an eigvec

Quadratic Form: shows how applying
the matrix affects length of a vector.
 $\|Ax\|^2 = x^T A^T A x$

Ellipsoids radii are the reciprocals of
eigenvalues, in matrix A maps
sphere to ellipsoid

bigger eigenval \Rightarrow steeper hill \Rightarrow shorter ellipsoid
radius.

A is diagonal \Leftrightarrow eigvec are coordinate
axes \Leftrightarrow ellipsoids are axis aligned

Positive definite: $w^T B w > 0 \forall w \neq 0$
 $\Leftrightarrow \lambda > 0$

PSD: $w^T B w \geq 0 \forall w \Leftrightarrow \lambda \geq 0$

Indefinite: $\lambda > 0, \lambda < 0$

Invertible: $\lambda \neq 0$

+ eigenval: curvature goes up
- eigenval: curvature goes down

$A = V \Lambda V^T = \sum_{i=1}^n \lambda_i v_i v_i^T$

diagonal \rightarrow rotate to be axis aligned

$A^2 = V \Lambda V^T V \Lambda V^T = V \Lambda^2 V^T$

$M^{1/2} = A$

1) compute eigenval of M

2) take square roots of eigenval

3) reassemble A w square roots as
eigenval, same eigvec

$P(x) = \frac{1}{(\sqrt{2\pi})^d} \frac{1}{|Z|} \exp(-\frac{1}{2}(x-\mu)^T Z^{-1}(x-\mu))$

Σ : covar matrix $(x-\mu)$, $d = \#$ features

$|Z|^{-1}$: precision matrix

$P(x) = n(C(x), \mu) = (x-\mu)^T Z^{-1}(x-\mu)$

$g(x)$ is signed distance from
 $Z^{-1/2} x$ to $Z^{-1/2} \mu$

Cov(C, S) = ~~covariance~~

$= E((C - E(C))(S - E(S))^T)$

$= E(C^T S^T) - \mu_C \mu_S^T$

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

$\Rightarrow R_i, R_j$ independent

All features pairwise independent
 $\Rightarrow \text{Var}(C)$ is diagonal

or!
 $\text{Var}(C) = \begin{bmatrix} \text{Var}(R_1) & \text{Cov}(R_1, R_2) & \dots & \text{Cov}(R_1, R_d) \\ \text{Cov}(R_2, R_1) & \text{Var}(R_2) & \dots & \text{Cov}(R_2, R_d) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(R_d, R_1) & \dots & \dots & \text{Var}(R_d) \end{bmatrix}$

$\Sigma = V \Lambda V^T$ eigenval of Σ are variances
along the eigvec, $\sigma_{ii} = \sigma_i^2$

$\Sigma^{1/2} = V \Lambda^{1/2} V^T$: maps spheres to ellipsoids

eigenval are Gaussian widths / ellipsoid
radii (standard deviations)

Isocountours of multivariate
normal distr are same as
isocountours of quadratic form
of Σ^{-1} .

Anisotropic Gaussian

QDA: $\hat{\Sigma}_c = \frac{1}{n_c} \sum_{i: Y=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$

For anisotropic! decision boundary could be hyperbola
 * need to apply logistic fn to find decision boundary

LDA: Decision fn is linear, decision boundary is hyperplane.

Maximize linear discriminant fn:

$$w^T X - \frac{1}{2} w^T \Sigma^{-1} w + \ln \pi_c$$

For 2 classes:

LDA has $d+1$ params

QDA has $\frac{d(d+3)}{2} + 1$ params

With features, LDA can give nonlinear boundaries, QDA using quadratic

LDA/QDA work well when data can only support simple decision boundaries such as linear/quadratic, bc Gaussian provide stable estimates

X : $n \times d$ design matrix of sample pts

Centering: subtracting μ^T from each row of X

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

Decomposing X : applying $Z = X V$, $\text{Var}(Z) = V \Delta V^T$

V transforms sample points to eigens coordinate system

Sphering X : applying transform $W = X \text{Var}(X)^{-1/2}$

Whitening X : centering + sphering, $X \rightarrow W$

W has covariance matrix I

Whitening: puts features on equal basis

Regression

given point x , predict a numerical value

Regression fns: linear: $h(x; w, d) = w \cdot x + d$

polynomial (equivalent to linear w poly features)

logistic $h(x; w, d) = \sigma(w \cdot x + d)$

Loss fns: z = prediction $h(x)$, y = true val

Squared Error: $LC_2(y) = (z - y)^2$

Absolute Error: $LC_1(y) = |z - y|$

Cross Entropy: $LC_2(y) = -y \ln z - (1-y) \ln (1-z)$

Cost fns: minimize

$J(w) = \frac{1}{n} \sum_{i=1}^n L(h(x_i; w), y_i)$ mean loss

$J(w) = \max_i L(h(x_i; w), y_i)$ max loss

$J(w) = \sum_{i=1}^n w_i L(h(x_i; w), y_i)$ weighted sum

$J(w) = \frac{1}{n} \sum_{i=1}^n L(h(x_i; w), y_i) + \lambda \|w\|^2$ L2 regularized

$J(w) = \frac{1}{n} \sum_{i=1}^n L(h(x_i; w), y_i) + \lambda \|w\|_1$ L1 regularized

Least Squares Linear Regression

Find w that minimizes $\|Xw - y\|^2 = \text{RSS}(w)$

Soln: $X^T X w = X^T y$

Let $X^+ = (X^T X)^{-1} X^T$

$\hat{w} = w \cdot x_i \Rightarrow \hat{y} = Xw = X X^+ y = I \cdot y$

Xw is subspace of \mathbb{R}^n spanned by columns

$X \in \mathbb{R}^{n \times d}$

$y \in n$ d's space

Minimizing $\|y - y^+\|^2$ finds y^+ nearest y in subspace

= orthogonal projection.

Advantages:

- easy to compute; just solve a linear system

- unique, stable soln.

Disadvantages:

- sensitive to outliers

- fails if $X^T X$ is singular

Logistic Regression - discriminative

fits probabilities in range $(0,1)$

used for classification

Find w that minimizes:

$$J = - \sum_{i=1}^n (y_i \ln \sigma(Xw_i) + (1-y_i) \ln (1 - \sigma(Xw_i)))$$

$\sigma(x) = \frac{1}{1 + e^{-x}}$

$\nabla_w J = -X^T (y - \sigma(Xw))$

Gradient Descent Rule: $w \leftarrow w + \epsilon \nabla_w J$

SGD: $w \leftarrow w + \epsilon (y_i - \sigma(Xw_i)) X_i$

Least Squares Polynomial Regression

- replace each X_i with vector

$$\Phi(X_i) = [X_i^0 \ X_i^1 \ X_i^2 \ \dots \ X_i^d]^T$$

$X_i^0 \ X_i^1 \ X_i^2 \ \dots \ X_i^d$

Weighted Least Squares Regression

- assign trusted sample points

a higher weight w

Greater $w_i \rightarrow$ work harder to

minimize $\|y - y^+\|^2$

Find w that minimizes

$$(Xw - y)^T \Omega (Xw - y)$$

$$= \sum_{i=1}^n w_i (X_i \cdot w - y_i)^2$$

Solve for w in normal eqn:

$$X^T \Omega X w = X^T \Omega y$$

Newton's Method - iterative

iterative optimization (least

method for smooth fn $J(w)$)

- faster than gradient descent

Idea: at point v ,

Approximate $J(w)$ near v by

quadratic fn. Jump to critical

pt.

pick starting point w

repeat till convergence

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

$w \leftarrow w + e$

The closer J is to quadratic,

the faster Newton's converges.

- doesn't know difference b/w

optim, must start close enough to soln.

Advantages

- tries to find right step length +

reach min

- tries to choose better descent

direction than just steepest descent

Disadvantages

- computing Hessian is expensive

- doesn't work for nonsmooth fns

Newton's for logistic regression.

$$S_i = \sigma(X_i^T w)$$

$$\nabla_w J(w) = -X^T (y - S)$$

$$\nabla_w^2 J(w) = X^T \Omega X, \Omega = \begin{bmatrix} s_1(1-s_1) & 0 \\ 0 & s_n(1-s_n) \end{bmatrix}$$

$$e \leftarrow (X^T \Omega X)^{-1} e \leftarrow X^T (y - S)$$

Ω prioritizes points with $s_i \approx 0.5$,

tunes at pts near 0.1

\Rightarrow sample pts near decision boundary

have largest effect on iterations,

made by contribution to logistic fit.

If n very large, save time by

using a random subsample of pts

per iteration, increase sample size

as you go.

LDA vs Logistic Regression

LDA:

- stable for well-separated classes

- more accurate when classes

nearly normal

Logistic Regression

- more emphasis on decision boundary

CLDA gives equal weight)

- less sensitive to most outliers

- easy treatment of partial membership

- LDA pts all or nothing

- more robust on non-Gaussians

RUC curve evaluates classifier after

it's trained

- shows rate of false positive vs

negative true positive

x -axis: false positive rate

y -axis: true positive rate (sensitivity)

False-negative rate: vertical distance

curve to top $(1 - \text{sensitivity})$

true negative rate: horizontal distance

from curve to right 'specificity'

$(1 - \text{false positive rate})$

True positive rate

False positive rate

False positive rate

Upper Right Corner: always classify +

Lower Left: always classify -

diagonal: random classifier

Classifier's effectiveness = area under

curve

Always right = 1.

Random = $1/2$.

Model of Reality:

sample pts from unknown prob distr

y -values are sum of unknown non-zero

fn + random noise

$Y_i = f(X_i) + \epsilon_i, \epsilon_i \sim D, \mu = 0$

Risk for hypothesis h is expected loss

$$R(h) = E[L]$$

Empirical distr: discrete uniform distr

over sample pts.

Empirical Risk: expected loss under

empirical distr

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n L(h(X_i), Y_i)$$

Max likelihood explains where logistic

loss fn comes from

2 sources of error in h

bias: error due to instability of h to

fit & perfectly

variance: error due to fitting

random noise in data

$$R(h) = (E[h(z)] - t(z))^2$$

bias

variance

Underfitting: too much bias

Overfitting: too much variance

variance $\rightarrow 0$ as $n \rightarrow \infty$

- adding good feature reduces

bias; adding bad feature

reduces ~~variance~~ increases

variance

adding a feature increases

variance

noise in test set affects only

$\text{Var}(E)$

noise in training affects

bias + $\text{Var}(h)$

Ridge Regression

Find w that minimizes

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

w' is w w'd replaced by 0.

regularization terms guarantees

PD \Rightarrow unique soln

many minims: ill-posed

regularization reduces variance

Ideally: features "normalized" to

have same variance

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

Subset Selection

- all features increase variance,

but not all features reduce bias

- identify poorly predictive features,

ignore them

- less overfitting, smaller test error

Alg: try all $2^d - 1$ nonempty subsets

of features

\rightarrow train one classifier per subset +

choose best classifier by cross-validation.

Heuristic 1: Forward stepwise selection

- repeatedly add best feature until

validation error starts increasing. (Ocd)

Heuristic 2: Backward stepwise selection

- start w all d features, remove feature

whose removal gives best reduction in validation

error (Ocd)

Goal: try to remove features w small weights

Lasso

- often naturally sets some weights to zero

Find w that minimizes $\|Xw - y\|^2 + \lambda \|w\|_1$

$$\|w\|_1 = \sum_{i=1}^d |w_i|$$

$\lambda \uparrow$, more + more weights go to 0.

$$Tf(x) = \frac{\partial T}{\partial x} (A^T x)_i = \sum_{j=1}^d A_{ij}^T x_j$$

$$(A^T)_i = \sum_{j=1}^d A_{ij}^T x_j$$

$$f(x + \Delta) \approx f(x) + \frac{\partial f}{\partial x} \Delta + o(\|\Delta\|)$$

Gradient = $\frac{\partial f}{\partial x}$

Posterior: $P(Y = y | x)$

$$\Theta_{MLE} = \arg \max_{\Theta} P(X | \Theta)$$

$$= \arg \max_{\Theta} \prod_{i=1}^n P(X_i | \Theta)$$

$$\Theta_{MAP} = \arg \max_{\Theta} P(\Theta | X)$$

$$= \arg \max_{\Theta} P(X | \Theta) P(\Theta)$$

$$= \arg \max_{\Theta} \prod_{i=1}^n P(X_i | \Theta) P(\Theta)$$

Covariance matrix must be PSD

Gaussian prior = L2 reg.

Laplace prior = L1 reg

Cov = $\frac{1}{n} X^T X$

determinant = product of eigval

PDF of univariate Gaussian:

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

Bayes:

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

Maximum a posteriori

- maximizing posterior $P(w | X, y)$

length of ellipsoid axes for multivariate

Gaussian are $\sqrt{\lambda_i}$ of Σ

- train on more data to improve training

accuracy, less data to improve test

- Bayes risk is 0 when class distr

don't overlap, prior for one class = 1

- add λI to covariance matrix

causes isocountours to be more spherical

$\sigma^2 = 1/2$ is decision boundary

No covariance terms: isocountour

is axis aligned

- space between contours is

$\sqrt{\lambda_i}$ of Σ

- higher SD \Rightarrow larger gaps

between significant

isocountours.

- eigenvectors form

orthonormal basis

- tell of something is

Equivalent Objective F_n :

$$\min \sum_{y \in \mathcal{Y}} \frac{1}{n_y} \sum_{x_i \in \mathcal{X}_y} |x_i - x_j|^2$$

→ average squared distance from point to all other pts in cluster

K-Medians clustering: only requires pairwise distances!

Euclidean distance (means not good)

specify distance $f_n d(x, y)$ between points

replace mean w median, sample pt that minimizes total distance to other pts in same cluster

medial always a sample pt.

Hierarchical Clustering: only requires pairwise distances

creates a tree, every subtree is a cluster

bottom up; agglomerative clustering: start with each point a cluster, repeatedly fuse pairs

top-down; divisive clustering: start with all points in one cluster; repeatedly split it

Distance f_n for clusters A, B

complete linkage: max distance between any point in A, any pt in B

max $\{d(w, x) : w \in A, x \in B\}$

single linkage: min distance

min $\{d(w, x) : w \in A, x \in B\}$

average linkage: average over all pairs $w \in A, x \in B$

$d(A, B) = \frac{1}{|A||B|} \sum_{w \in A, x \in B} d(w, x)$

centroid linkage: distance between means

$d(\mu_A, \mu_B)$

Greedy Agglomerative: repeatedly fuse two clusters that minimize $d(A, B)$. $O(n^3)$

Dendrogram: Illustration of cluster hierarchy in which vertical axis encodes all linkage distances

Spectral Graph Clustering

Input: weighted undirected graph $G=(V, E)$, every pair of vertices has weight

edge weights = similarity measure

big weight means two vertices want to be in same cluster

Goal: cut G into 2 pieces G_1, G_2 of similar size

e.g. minimize sparsity $\frac{Cut(G_1, G_2)}{Mass(G_1)Mass(G_2)}$

$Cut(G_1, G_2)$ = total weight of cut edges

$Mass(G_i) = \# \text{ of vertices in } G_i$

Denominator penalizes imbalanced cuts

$n=|V|, y_i = \begin{cases} 1, & \text{vertex } i \in G_1 \\ -1, & \text{vertex } i \in G_2 \end{cases}$

$w_{ij} = \frac{(y_i - y_j)^2}{4} = \begin{cases} w_{ij} & \text{if } i, j \text{ is cut} \\ 0 & \text{if } i, j \text{ isn't cut} \end{cases}$

$Cut(G_1, G_2) = \sum_{(i,j) \in E} w_{ij} \frac{(y_i - y_j)^2}{4}$

$= y^T L y$

$L_{ij} = \begin{cases} -w_{ij}, & i \neq j \\ \sum_k w_{ik}, & i = j \end{cases}$ - off diagonal

L is symmetric, $n \times n$: matrix rep of G

Bi-section: $1^T y = 0$

Find y that minimizes $y^T L y$ subject to

$y_i, y_i = -1$ and

$1^T y = 0$ - balance

lets relax binary constraint to allow fractional vertices

Now: y must lie on hypersphere of radius \sqrt{n}

minimize $y^T L y$ - an and $1^T y = 0$

→ minimize $\frac{y^T L y}{y^T y}$ = Rayleigh quotient

The y that minimizes this is eigenvector of smallest eigenval.

λ_2 = second smallest eigenval of L

eigenvector v_2 = Fiedler vector

Spectral partitioning alg:

compute Fiedler vector $v_2 \in \mathbb{R}^n$

sort components of v_2

the $n-1$ cuts between successive components. choose min sparsity cut

Vertex Masses:

assign masses to vertices

lets L be diagonal matrix with vertex masses on diagonal

new balance constraint: $1^T M y = 0$

minimize $y^T M y$ - Mass $(G_i) = \sum M_{ii}$

Want Fiedler vector of generalized eigensystem $L v = \lambda M v$

Greedy Divisive Clustering

partition G into 2 subgraphs; recursively cluster them

Normalized Cut:

vertex i 's mass $M_{ii} = L_{ii}$ - sum of edge weights adjoining vertex i

Image Segmentation

$w_{ij} = \exp(-\frac{(w_i - w_j)^2}{\sigma^2} - \frac{(b_i - b_j)^2}{\beta^2})$

location w_i , brightness b_i

→ cluster pixels with similar color

have greater weight.

each pixel is vertex in graph

Clustering w Multiple Eigenvectors

for k clusters, compute first k eigenvectors of $L v = \lambda M v$

$V = [v_1 \dots v_k]$

Row V_i is spectral vector for vertex i

cluster vertices together if spectral vectors point in similar directions

use k -means to cluster together vector separated by small angle

Latent Factor Analysis

X_{ij} = occurrence of term j in doc i

SVD $X = U D V^T = \sum \sigma_i u_i v_i^T$

diagonals greatest to smallest

greatest σ_i, v_i lists terms in a given cluster of documents

u_i lists documents in a given cluster with similar terms, i.e. make similar documents.

σ_i is a score.

game is hidden; nothing explicitly states it.

LFA = clustering that allows clusters to overlap.

Applications:

- Fuzzy Search (related words)

- Denoising

- matrix compression

- collaborative filtering: fill in unknown values

Geometry of High Dimensional spaces

shell between spheres of radii r_1, r_2

random points from uniform disk in \mathbb{R}^2 : nearly all in outer shell

from Gaussian: nearly all in some shell

variations in distance less extreme in higher dimensional space

→ clustering less effective in higher dimensions

Random Projection: dimensionality reduction

Adaptive Filtering: random values

cheap alternative to PCA: don't have to compute eigenv.

random subspace $S \subset \mathbb{R}^d$ of dim k

$K = \frac{2 \ln(1/\delta)}{\epsilon^2}$ random $6, \delta$

For any pt g , let $\hat{g} = \frac{1}{K} \sum_{i=1}^K g_i$

orthogonal projection of g onto S

With probability $\geq 1 - 2\delta$, distance between two points after projecting is within some range of original distance

For $g, h \in \mathbb{R}^d, (1-\epsilon) \|g-h\|^2 \leq \|\hat{g}-\hat{h}\|^2 \leq (1+\epsilon) \|g-h\|^2$

Boosting

ADA BOOST: adaptive boosting

ensemble method

train multiple learners on weighted sample points

use different weights for each learner

increase weights of misclassified sample points

gives bigger votes to more accurate learners

weight for sample point grows according to how many it's misclassified

$M(z) = \sum_{i=1}^n p_i G_i(z)$

M 's continuous

Risk = average loss

$R(M) = \frac{1}{n} \sum_{i=1}^n L(M(x_i), y_i)$

Metalearner uses exp. loss for $L(p, k) = e^{-p k} \sum_{i=1}^n e^{p k_i}$

exp loss p pushes hard against badly misclassified point

Optimal learner G_T minimizes risk by minimizing sum of weights over all misclassified pts x_i

$w_i^{(t+1)} = w_i^{(t)} e^{-\eta y_i G_t(x_i)}$

$B_T = \frac{1}{T} \ln \left(\frac{1 - \text{err}_T}{\text{err}_T} \right)$

$\text{err}_T = G_T$'s weighted error rate

$= \sum_{i=1}^n w_i^{(T)} \frac{y_i - G_T(x_i)}{2}$

→ more accurate data learners get bigger votes in metalearner

AdaBoost:

1. initialize weights $w_i \leftarrow \frac{1}{n}$

2. for $t \in 1$ to T

a. Train G_t w weights w_i

b. Compute error, B_t

c. reweight points

3. return metalearner

$\text{Sign}(\sum_{t=1}^T B_t G_t(x))$

Boost decision trees b.c.:

- Fast

- no hyperparameter search

- easy to note tree best for training accuracy

- easy bias variance control

- AdaBoost's weak learner is form of robust rejection

- linear decision boundaries don't boost well.

posterior is approximated

nearest Neighbors

given query point g , find k sample pts nearest g

- distance metric of your choice

- optional return average label of k points

- classification: return class w most votes from k pts or return histogram of class probabilities

Exhaustive KNN

- maintain max heap w k shortest distances seen so far.

Query Time: $O(n \log k)$

n points, d dims, k shortest distances

Preprocess training points to get sublinear query time

- 2-5 dimensions: Voronoi diagrams

- median: $k-d$ trees

- large: exhaustive w dim red.

Voronoi Diagram

each point has cell which contains space closest to that point.

Diagram: set of k Voronoi cells.

Size (# of vertices) $\in O(n^{d/2})$

2D: $O(n \log n)$ to compute V.D. and trapezoidal map for pt location.

$O(\log n)$ query time.

3D: use binary space partition tree for pt location.

really only supports 2 nearest neighbors

$k-d$ Trees

- decision tree

- choose splitting feature w greatest width

→ maximize distance from leftmost value to rightmost

- splitting value: median point for feature i

each subtree represents particular box in space.

Query Algorithm maintains:

- nearest neighbor found so far - distance d

- binary heap of unexplored subtrees, keyed by distance d

Alg:

1. add heap containing root node w key zero

2. while Q not empty and $(1+\epsilon) \cdot \text{minkey}(Q) < \text{key}(Q)$

a. $B \leftarrow \text{remove min}(Q)$

b. $w \leftarrow B$'s sample point

c. $r \leftarrow \min_{i \in \{1, \dots, d\}} \text{dist}(w, x_i)$

d. $B' \leftarrow$ child nodes of B

e. if $C(1+\epsilon) \cdot \text{dist}(g, B) < \text{key}(B')$, insert $(B', \text{dist}(g, B'))$ into Q

return pt that determined r .

Approximate nearest neighbor speeds up alg.

$\frac{1}{n} E(X) = E(\frac{X}{n})$

$\text{Var}(\frac{X}{n}) = \frac{1}{n^2} \text{Var}(X)$

$\text{Cov}(X, Y) = \text{Cov}(X, Y) \cdot \sqrt{\text{Var}(X) \cdot \text{Var}(Y)}$

calculating distance in p polynomial space takes d^p points in \mathbb{R}^d

features get raised to p polynomial power, number features $n \cdot p$.

Valid Kernel:

$K(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$

$K = [K(x_1, x_1) \dots K(x_1, x_n); \dots; K(x_n, x_1) \dots K(x_n, x_n)]$ PSD

if K is a valid kernel, $\exists \lambda \in \mathbb{R}$

$\frac{d}{dx} \|Ax - x\|_2^2 = -2x^T(Ax - x)$

o left multiply by orthogonal matrix preserves length of vector

cov matrix: $\frac{1}{n} X^T X$ Convex Concave

nonlinear
2 node types:

- 1) internal nodes test feature values & branch accordingly
- 2) leaf nodes specify class
- if all sample pts in a node are not of same class, we need to split
- choose best splitting feature j and splitting value θ

$$S_L = \{i: X_{ij} < \theta\}, S_R = \{i: X_{ij} \geq \theta\}$$

Try all features/all splits in a feature (i, j) in cost \mathcal{C}

$$\text{Goal: minimize } J(CS_L) + J(CS_R) \text{ or } |S_L| J(CS_L) + |S_R| J(CS_R)$$

Cost: Mean entropy
let $P(CY=C) = P_C$

Surprise: being class $C = -\log_2 P_C$
Entropy:

$$H(C) = -\sum P_C \log_2 P_C$$

P_C = proportion of points in S that are in class C .
expected number of bits of info we need to transmit to identify class of a sample point chosen uniformly at random
Hatter = $|S_L| H(CS_L) + |S_R| H(CS_R)$

Choose split that minimizes Hatter
 $H(C) - \text{Hatter} = \text{minimize Hatter}$

Key: cost function is concave

Training: quantitative feature: $O(N \cdot d)$
splits, d features, n points in node

Total runtime: $O(N \cdot d \cdot \text{depth})$

Decision trees can use multivariate splits, i.e. consider more than one feature to split on

better classifier at cost of worse interpretability or speed

$$\text{Regression: } J(CS) = \frac{1}{|S|} \sum_{i \in S} (Y_i - \bar{Y})^2$$

\bar{Y} = mean (Y_i) for points in S .
Can stop building tree early:

- most split doesn't reduce entropy enough
- most of node's points have same class
- node contains few sample pts
- cell's edges too tiny
- depth too great
- use cross validation

Pruning:

- greedily remove each split whose removal improves cross validation performance

Decision trees have high variance

- take average of output of different learning alg
- some learning alg on many training sets
- bagging: same learning alg on many random subsamples of one training set
- Random Forests: randomized decision trees on random subsamples
- use learners with low bias - averaging reduces variance.

Bagging:

- Given n -point training sample, generate random subsample of size n' by sampling w replacement
- build learner, then repeat on n' samples

Random Forests better approximate posteriors
at each tree node, take random sample of n features, choose best split from n features

Disadvantage: (use interpretability / inference)

Kernels don't have to explicitly compute feature space n

Suppose $w = X^T a = \sum_{i=1}^n a_i X_i$ for some $a \in \mathbb{R}^n$

optimize n dual weights a .

Kernel Ridge Regression:

- center X and y so means are zero.
 $(X^T X + \lambda I) w = X^T y$

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

$$\min_a \|X X^T a - y\|^2 + \lambda \|X^T a\|^2$$

Regression fn:

$$h(x) = w^T x = a^T X x = \sum_{i=1}^n a_i (X_i^T x)$$

$K(x, z) = x^T z$ is kernel fn

$$K = X X^T = n \times n \text{ Kernel matrix, } K_{ij} = K(X_i, X_j)$$

Dual n/2 reg alg:

$$V_{ij}, K_{ij} = K(X_i, X_j)$$

Solve $(K + \lambda I) a = y$ for a .

for each test pt z

$$h(z) = \sum_{i=1}^n a_i K(X_i, z)$$

Polynomial kernel of degree p

$$K(x, z) = (x^T z + 1)^p = \phi(x)^T \phi(z), \phi(x) \text{ contains every monomial in } x \text{ of degree } 0 \dots p$$

W/it we can compute $\phi(x)^T \phi(z)$ in $O(d)$ time instead of $O(d^2)$

- don't compute $\phi(x)$, $\phi(z)$ directly & evaluate $w = \phi(x)^T a$, do everything in terms of a .

$$\phi(X_i)^T w = [\phi(X_i) w]_i = [\phi(X_i) \phi(X)^T a]_i = (K a)_i$$

Gaussian Kernel: $K(x, z) = \exp(-\frac{\|x-z\|^2}{2\sigma^2})$, $O(d)$

think of kernel k as measure of how similar points are

$$h(z) = \sum_{i=1}^n a_i K(X_i, z) \text{ is linear comb of Gaussians centered at sample pts}$$

\rightarrow each sample pt votes for classification of surrounding area

Popular!

- Oscillates less than polynomials - wide σ has less oscillation
- $K(x, z)$ interpreted as similarity measure
- sample points vote for values of z , closer points get higher vote.

Standard Kernel $x \cdot z$ assigns more weight to vectors that point in same direction. Direction vs Magnitude

σ trades off bias vs variance
larger $\sigma \rightarrow$ wider Gaussian, smoother \rightarrow more bias, less variance.

Kernel Nuts
Logistic fn: $\sigma(x) = \frac{1}{1 + e^{-x}}$

$$\text{Training: } S \&D / B \&D$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n L(h(X_i), Y_i)$$

\rightarrow average loss per prediction

- start w as random weights to differentiate and get better local minima.

Backprop:

$$\frac{\partial}{\partial z} L(z, \hat{z}), \frac{\partial}{\partial z} L(z, \hat{z}) = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial z} + \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial z}$$

NN Variations

- Regression: usually linear and with - and sigmoid fn

Classification - for 2,3 classes, use softmax fn: probability of z_j is $\frac{e^{z_j}}{\sum_{i=1}^n e^{z_i}}$

Used Saturation/Vanishing gradient
when unit output is close to 0 or 1 for most training points, $S'(x) \approx 0$.

Mitigation!
1) Initial weight of edge into unit w fanin q random w mean zero, SD $\frac{1}{\sqrt{q}}$

2) Set target values to 0.15 & 0.85 instead of 0 or 1.

small constant to ϵ
- hack so unit can't get stuck

4) Cross Entropy Loss instead of sq. error

$$L(z, y) = -\sum_i (y_i \ln z_i + (1-y_i) \ln (1-z_i))$$

$$\frac{\partial L}{\partial z_i} = \frac{z_i - y_i}{z_i(1-z_i)}, z = \text{prediction, } y = \text{truth}$$

Used only for sigmoid & softmax outputs (classification)

Regression: use linear outputs which don't saturate

5) replace sigmoids w ReLUs: $\text{ReLU}(x) = \max\{0, x\}$

- vulnerable to exploding gradient problem because output is arbitrarily large; problem in deep / recurrent net

Heuristics for Faster Training

S&D: faster than batch on large, redundant data sets. takes less on one sample point

Normalizing Data

- center each feature so mean is zero. makes it easier for hidden units to get into good operating region of sigmoid or ReLU

- scale each feature so variance ≈ 1 makes each objective fn better conditioned so GD converges faster

Converting hidden units to replace sigmoids with tanh $\frac{e^x - e^{-x}}{e^x + e^{-x}}$

\rightarrow range from -1 to 1 $\frac{e^x - e^{-x}}{e^x + e^{-x}}$

Different learning rate for $= 2s(2x) - 1$

- each layer of weights

- earlier layers have smaller gradients, need larger learning rate

Emphasizing schemes (uncommon samples)

Second order optimization

- Hessian too large expensive to compute

- nonlinear conjugate gradients: work well for small nets & small data

- stochastic Levenberg Marquardt

- approximates diagonal Hessian

acceleration schemes: RMS Prop, Adam, AMS Grad

Avoiding Bad Local Minima

- SGD

- momentum $w \leftarrow w + \Delta w$

$\Delta w \leftarrow -\epsilon \nabla J(w) + \beta \Delta w$

Avoiding Overfitting

- NN ensemble: random initial weights & bagging

- L2 regularization (weight decay)

- Dropout

- fewer hidden units

- too few, can't learn well

- too many, overfit

Conv Nets - late layers: higher level

1. Local connectivity: hidden units in early layer connect only to small patch of units in prev layer

2. Shared Weights: groups of hidden units share same set of input weights, called mask / filter / kernel

- Center X

- compute unit eigenvectors of $X^T X$

- choose k = dim of subspace we're projecting down to.

- pick k largest eigenv

- compute $x \cdot v_i$ of training data in PC space

$$\% \text{ of variability} = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^n \lambda_i}$$

- want to retain as much variability as possible

PCA Deriv. 2.

- find direction w that maximizes sample variance of projected data

$$\text{Max Var}(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n) = \frac{1}{n} w^T X^T X w$$

Rayleigh quotient of $X^T X$, w .

If w is eigenv V_i , Rayleigh quotient = λ_i of $X^T X$

Deriv 3: minimize proj. error i.e. minimize error in direction orthogonal to principal component direction

- minimizing projection error = maximizing variance

Singular Value Decomposition

$$X = U D V^T = \sum_{i=1}^n \delta_i u_i v_i^T$$

and and and and

$$U^T U = I, V^T V = I$$

Orthogonal U_i 's are left singular vectors of X

orthogonal V_i 's are right singular vectors of X

- diagonal entries $\delta_i, \dots, \delta_n$ are nonnegative singular values of X .

V_i is eigenv of $X^T X$ w/ eigen δ_i^2

U_i is eigenv of $X X^T$, $X X^T U_i = \delta_i^2 U_i$

- we need eigenv of $X^T X$ for PCA!

- can find k greatest singular values & corresponding vectors in $O(nd \log k)$ time

- row i of $U D$ gives coord of sample point X_i in principal components space.

increase k : k -Means Clustering

Goal: Partition n points into k disjoint clusters

- assign each input point X_i a cluster label

- cluster i 's mean is $\mu_i = \frac{1}{n_i} \sum_{j \in i} X_j$

Find y that minimizes $\sum_{i=1}^n \sum_{j=1}^k \|X_j - \mu_i\|^2$

(Sum of squared distances from pt to cluster means)

NP-hard: $O(n^k)$, n points, k clusters

1) y 's are fixed; update μ_i 's

2) μ_i 's are fixed; update y_i 's

- both steps decrease objective fn unless they change nothing

Starting k -means...

Foggy: choose k random sample pts to be initial μ_i 's; go to 2.

Random partition: randomly assign each sample pt to a cluster; go to 1.

k -means++: like Foggy, but center i 's chosen w preference for points far from previous centers.