

①② Classifiers

n observations, each w/ d features
 Perceptron Alg: slow but for linearly separable
 risk func: $R(w) = \frac{1}{n} \sum_{i=1}^n -y_i x_i \cdot w$

③ Grad Descent

Grad Descent: $w \leftarrow w + \epsilon \sum_{i=1}^n y_i x_i$
 Stochastic Grad Desc: $w \leftarrow w + \epsilon y_i x_i$

④ Support Vector Mach (SVM); Features

Hard Margin SVM: find w, α to min $\|w\|^2$
 constraints $y_i(w \cdot x_i + \alpha) \geq 1$

Soft Margin SVM: slack var to violate
 find w, α and ξ to min $\|w\|^2 + C \sum \xi_i$
 st $y_i(x_i \cdot w + \alpha) \geq 1 - \xi_i, \xi_i \geq 0$
 big C : keep slack var small, overfitting, sensitive to outliers
 Features: make nonlinear features into higher dim space, separable

⑤ ML Abstractions & Numerical Optimization

- 1) Unconstrained Optimization: min cost of func $f(w)$
 - 2) Constrained: $f(x, y) = 0$, find w min $f(w)$
 - 3) LP: max $c \cdot w$ st $A \cdot w \leq b$, n linear constraints $A: n \times d; b: n \times 1$
 - 4) QP: find w : min $f(w) = w^T Q w + c^T w$ st $A \cdot w \leq b$
- Positive Definite (PD): $w^T Q w > 0$ for all $w \neq 0$

⑥ Design Theory

Risk: expected loss over x, y $E[L(w; x, y)]$
 Bayes Classifier: min $R(w)$ given $L(z, y) = 0$ for $z=y$
 optimal Risk $R^*(w) = \int \min_y L(y, y) f(x=x|Y=y) P(Y=y) dx$

⑦ Gaussian Discriminant Analysis

Assume Normal $X \sim N(\mu, \sigma^2) = \frac{1}{(\sqrt{2\pi})^d} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
 QDA: diff σ, μ
 To recover posterior probability $S = \frac{1}{\mu e^T}$
 $P(Y=C|X) = \frac{f(x|Y=C) \pi_c}{f(x|Y=C) \pi_c + f(x|Y=D) \pi_D} = S(a_c(x) - a_D(x))$

LDA: same σ , diff μ , approx Bayes Decision rule
 max $\frac{y_c \cdot x}{\sigma^2} - \frac{\|u\|^2}{2\sigma^2} + \ln \pi_c, P(Y=C|X=x) = S(w \cdot x + b)$

Decision Boundary: $a_c(x) - a_D(x) = 0$


MLE: estimate model params

$L(\mu, \sigma; X_1, \dots, X_n) = f(x_1) f(x_2) \dots f(x_n)$
 log likelihood $\ell = \max_{\mu, \sigma} \sum_{i=1}^n \left(-\frac{\|x_i - \mu\|^2}{2\sigma^2} - d \ln \sqrt{2\pi} - d \ln \sigma \right)$
 $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - \mu\|^2$

⑧ Eigenvectors & Multivariate Normal Distr

if A is invertible v is eigenvector of A^{-1} w/ eigenvalue $\frac{1}{\lambda}$
 if v is an eigenvector of A w/ eigenvalue λ
 v is eigenvector of A^k w/ eigenvalue λ^k

⑨⑩ Regression, LS, Linear & Logistic Regression

converges

 Ridge: linear reg func + squared loss + L2 reg
 Lasso: linear reg func + squared loss + L1 reg
 Sparse sol
 Logistic: logistic reg func + weighted sum + mean cost
 LS: linear reg func + squared loss + mean cost

⑪ Regression, Newton's Method, ROC

LS: w that min $\|Xw - y\|^2$
 $w = (X^T X)^{-1} X^T y$
 Logistic Regression: find w min
 $J = \sum_{i=1}^n L(x_i; w, y_i) = - \sum_{i=1}^n (y_i \ln s(x_i; w) + (1 - y_i) \ln (1 - s(x_i; w)))$
 $\nabla_w J = - \sum \left(\frac{y_i}{s} v_i - \frac{1-y_i}{1-s} v_i \right) = -X^T (y - s(Xw))$
 Gradient descent: $w \leftarrow w + \epsilon X^T (y - s(Xw))$
 Stochastic Gradient descent: $w \leftarrow w + \epsilon (y_i - s(x_i; w)) x_i$
 Newton's Method: $w \leftarrow w + \epsilon (X^T S X)^{-1} X^T (y - s)$

⑫ Bias-Variance Decomposition

Bias: error due to inability of hypothesis h to fit g perfectly, wrong model
 Variance: error from fitting random noise in data,
 - underfitting = too much bias
 - overfitting = too much variance
 - training error tests bias not variance, test error for both
 - variance $\rightarrow 0$ as distributions $n \rightarrow \infty$
 - if h fit g well, distributions bias $\rightarrow 0, n \rightarrow \infty$
 - good feature reduces bias
 - adding feature increases variance
 - variance decreases as $\frac{1}{n}$ (sample pts), increases as d (features)

⑬ Ridge & Lasso Regression

Ridge: w to min $\|Xw - y\|^2 + \lambda \|w\|^2$
 reg reduces overfitting by reducing variance
 Maximum a posteriori (MAP): using likelihood, but max posterior
 posterior $f(w|x, y) = \frac{f(y|x, w) \times \text{prior } f(w)}{f(y|x)} = \frac{L(w) f(w)}{f(y|x)}$
 Maximize log posterior = $\ln L(w) + \ln f(w) - \text{const}$
 \Rightarrow Minimize $\|Xw - y\|^2 + \lambda \|w\|^2$
 Lasso: $\min \|Xw - y\|^2 + \lambda \|w\|_1$
 Feature Subset Selection: find poorly predictive features

Heuristic 1: Forward Stepwise Selection $O(d^2)$
 1) Start w/ null model (0 features)
 2) add best feature until valid error increasing
 Heuristic 2: Backward stepwise Selection $O(d^2)$
 1) start w/ d features
 2) remove gives best reduction in valid error

⑭ Decision Trees classification & regression

Tree w/ 2 node types:
 - internal nodes to test 1 feature value & branch
 - leaf nodes specify class $h(x)$
 Entropy: $H(S) = - \sum p_i \log_2 p_i$ $p_i = \frac{|S_i|}{|S|}$
 \hookrightarrow same class: 0, half C, half D: 1, n pts, d diff classes: $\log_2 n$
 Weighted Avg Entropy: $H_{\text{after}} = \frac{|S_L| H(S_L) + |S_R| H(S_R)}{|S_L| + |S_R|}$
 Information gain: $H(S) - H_{\text{after}}$ max info gain

⑮ Decision Tree, Ensemble Learning, Random Forest

Multivariate Split: find non-axis aligned splits or generate randomly
 Decision Tree Regression: creates piecewise constant regression func
 Stopping Early: most points same, depth too great, prone
 Pruning: greedily remove each split that improves validation, more reliable
 Ensemble Learning: decrease variance
 Bagging: same learning alg on random subsets of one training set
 Random Forests: randomized decision trees on random subsamples

⑯ Kernel Trick

Optimize $w = X^T a = \sum_{i=1}^n a_i x_i, (x^T x + \lambda I) w = X^T y$
 For solution $a, (X X^T + \lambda I) a = y, (x^T x - \lambda I) x^T a = x^T y$
 Dual Form Ridge Regression: min $\|X X^T a - y\|^2 + \lambda \|X^T a\|^2$
 Training: Solve $(X X^T + \lambda I) a = y$
 Testing: Regression $h(z) = W^T z = a^T X z = \sum_{i=1}^n a_i (x_i^T z)$
 Kernel func: $k(x, z) = x^T z$
 Kernel matrix: $K = X X^T, K_{ij} = k(x_i, x_j)$
 Polynomial Kernel: degree p is $k(x, z) = (x^T z + 1)^p$
 Thm: $(x^T z + 1)^p = \sum_{i=0}^p \binom{p}{i} x^T z^i$
 Gaussian Kernel: \exists st $k(x, z) = \exp(-\frac{\|x - z\|^2}{2\sigma^2})$

⑰ Neural Networks classification, regression

use a logistic function between linear combinations
 Network w/ 1 Hidden Layer
 Input: $x_1, \dots, x_n; x_{n+1} = 1$ Hidden: $h_1, \dots, h_m; h_{m+1} = 1$
 Output: z_1, \dots, z_k
 Layer 1 weights: $m \times (n+1)$ matrix V, V^T row i
 Layer 2 weights: $k \times (m+1)$ matrix W, W^T row i
 Backpropagation
 - dynamic programming algorithm for computing gradients we need for grad descent

⑱ Neurobio, Variations on NN

ReLU, ramp, hinge fn: $r'(x) = \begin{cases} 1 & x \geq 0 \\ 0 & x < 0 \end{cases}$

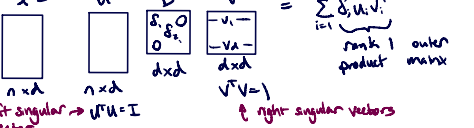
⑲ Better NN Training, CNN

Heuristics for Faster Training
 1) Stochastic Gradient Descent: faster than batch on large, redundant datasets
 2) Epoch: presents every training pt once
 3) Normalizing: center mean at zero, scale so variance = 1
 4) Different learning rate for each layer of weights
 Heuristic to Avoid Overfitting
 1) Fewer hidden units
 2) Weight decay: L_2 regularization
 Convolutional Neural Networks
 1) Local Connectivity: hidden layer to small patch of prev layer
 2) Shared weight: hidden units share mask, filter, kernel
 Convolution: same linear transformation applied to diff parts of image by shifting

⑳ Unsupervised Learning & Principal Components

Unsupervised Learning: sample points but no labels, discover structure
 PCA: find k directions to capture variance
 find w max $\text{Var}(\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_k) = \frac{1}{n} \sum_{i=1}^n (x_i \cdot w)^2 = \frac{1}{n} \sum_{i=1}^n \frac{X_{i,j}^2 w_j^2}{w_j^2} = \frac{1}{n} \sum_{j=1}^d X_{i,j}^2 w_j^2$

㉑ Singular Value Decomposition, Clustering

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

 $\delta_1, \dots, \delta_n$ are nonnegative singular values of X
 v_i is an eigenvector of $X^T X$ w/ eigenvalue δ_i^2
 u_i : of $U D$ gives principle coordinates of sample point x_i
 k -means clustering (Lloyd's Algorithm) partition n points into k disjoint clusters, assign x_i a cluster label $y_i \in \{1, \dots, k\}$
 min y in $\sum_{i=1}^n \|x_i - \mu_j\|^2$
 Hierarchical Clustering
 Creates a tree, every subtree is a cluster
 Agglomerative Clustering: each point a cluster and combine
 Hierarchical Clustering: all pts in one cluster, split

Dist func for clusters A, B

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

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

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

Centroid Linkage: $d(A, B) = d(\mu_A, \mu_B)$ μ mean of S

22 High Dimensions, Random Projection, Pseudoinv

P: sampled from univariate normal dist w/ mean 0, variance 1

$P_i \sim N(0, 1)$ $P_i^2 \sim \chi^2(1)$ $E[P_i^2] = 1$ $Var(P_i^2) = 2$

$E[\|P_i P\|^2] = d E[P_i^2] = d$ $Var(\|P_i P\|^2) = d Var(P_i^2) = 2d$ $SD(\|P_i P\|^2) = \sqrt{2d}$

Random Projection: Project onto random subspace, could preserve dist better

Choose small ϵ , small δ , random subspace $S \subset \mathbb{R}^d$ dimension k , $k \geq \left[\frac{2 \ln(\frac{2}{\delta})}{\epsilon^2} \right]$

For pt q , let \tilde{q} be orthogonal proj of q onto S , multiplied by \sqrt{d}

Pseudo inverse D^+ by transposing D and replacing zeros w/ reciprocal

$DD^+D = D$ $D^+DD^+ = D^+$ $D^+D^+D = D$

X is $n \times d$ matrix $X = UDV^T$ $rank(D) = rank(X)$

Moore-Penrose pseudoinverse X is $X^+ = V D^+ U^T$

- $XX^+ = UDV^+V^+U^+ = U(DD^+)U^+$ symmetric, PSD
- $X^+X = V D^+U^+UDV^+ = V(D^+D)V^+$ symmetric, PSD
- same rank: $D, D^+, DD^+, D^+D, X, X^+, XX^+, X^+X$
- $XX^+X = X$
- $X^+X^+X = X^+$

Solution to normal eq $X^T X w = X^T y$ is $w = X^+ y$

23 Learning Theory

Constrain hypotheses we allow learner to consider

Range Space (Set System): pair (P, H) , P set of all possible test/training

H : hypotheses class, set of hypotheses

- Power set classifiers: P set of k numbers, H power set of P , 2^k subsets of P
- Linear classifier: $P = \mathbb{R}^d$, H set of all halfspaces form $\{x: w \cdot x \geq -\alpha\}$

Hoeffding's inequality: prob of bad estimate, how likely number drawn from binomial dist far from mean

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

Dichotomy: $X \cap h, h \in H$, func that assigns training point to $\{0, 1\}$ or not \subset

Shatter func: $\Pi_H(n) = \max_{S \subseteq \{1, \dots, n\}} |\{h(S) : h \in H\}|$ most dichotomies of point set size n

Vapnik-Chervonenkis dimension: $VC(H) = \max \{n : \Pi_H(n) = 2^n\}$

VC dimension is upper bound on exponent of polynomial

24 Boosting; Nearest Neighbor Classification

AdaBoost: "adaptive Boosting" ensemble method for classification

- trains multiple learners on weighted sample points

- diff weights for each learner

- increases weights of misclassified training pts

- bigger vote to accurate learners

- Metalearner: linear combination of learners, $M(x) = \sum_{t=1}^T \beta_t f_t(x)$

$\min_{\beta_t} \frac{1}{n} \sum_{i=1}^n L(M(x_i), y_i)$, $M(x) = \sum_{t=1}^T \beta_t f_t(x)$, $L(p, r) = e^{-pr} = \begin{cases} e^{-p} & r=1 \\ e^p & r=-1 \end{cases}$

optimal β_t after $\frac{d}{d\beta_t} Risk = 0$, $\beta_t = \frac{1}{2} \ln \left(\frac{1 - err_t}{err_t} \right)$

Boosting benefits \uparrow : fast, no hyperparameter, subset selection,

Nearest Neighbor Classification

Idea: query point q find k sample pts nearest q distance metric

Regression: return avg label of k pts

Classification: return class w mos votes from k pts

25 NN Algos, Voronoi Diagrams, k-D trees

Exhaustive k-NN Alg: query pt q

- Scan through all pts and find squared dist to q

- keep max-heap w/ k smallest dist found so far

Voronoi Diagrams: X point set

Voronoi cell: $w \in X$ is $Var w = \{x \in \mathbb{R}^d : \|x - w\| \leq \|x - v\| \forall v \in X\}$

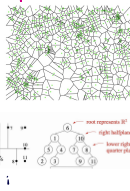
Point location: for query point $g \in \mathbb{R}^d$, find point $w \in X$ where $g \in Var w$

\rightarrow good for 1-NN in 2 or 3 dimensions

k-D Trees: Decision Trees for NN search

- Choose splitting feature: w/ greatest width

\rightarrow splitting value: median point for feature



Misc Algebra

L2 Norm (Euclidean): $\|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T x}$

L1 Norm: $\|x\|_1 = \sum_{i=1}^n |x_i|$

L ∞ Norm: $\|x\|_\infty = \max_{1 \leq i \leq n} |x_i|$

Frobenius Norm: $\|A\|_F = \sqrt{\sum_{i=1}^n \sum_{j=1}^n A_{ij}^2} = \sqrt{Tr(A^T A)}$

Orthogonal: $U^T U_j = 0, i \neq j$ and $\|u_i\|_2 = 1$

Trace: $Tr A = \sum_{i=1}^n A_{ii}$: sum of diagonal elements

Cauchy-Schwartz Inequality: $|z^T y| \leq \|y\|_2 \cdot \|z\|_2$

Fundamental Theorem of Linear Algebra

Range of a matrix is the orthogonal complement of the nullspace of its transpose

$R(A)^+ = N(A^T)$

Spectral Theorem (Symmetric eigenvalue decomposition (SEVD))

$A = \sum_{i=1}^n \lambda_i u_i u_i^T = U \Lambda U^T$, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$

$Tr \sum_{i=1}^n \lambda_i u_i u_i^T = \sum_{i=1}^n \lambda_i Tr(u_i u_i^T) = \sum_{i=1}^n \lambda_i \|u_i\|_2^2 = \sum_{i=1}^n \lambda_i$

Total Variance: $Tr \Sigma = Tr(U \Lambda U^T) = Tr(\Lambda U^T U) = Tr \Lambda = \lambda_1 + \dots + \lambda_n$

rank: linearly independent columns $X: \mathbb{R}^{m \times n}$

Rank-Nullity Theorem: $n - \dim(\text{nullspace}(X)) = \text{rank}(X)$

$\dim(\text{rowspace}(X)) = \dim(\text{columnspace}(X^T)) = \text{rank}(X^T) = \text{rank}(X)$

- symmetric matrix has real eigenvalues

- eigenvalues neg if concave

- Axis scaled by square roots of eigenvalues of Σ

- Covar = $\frac{XX^T}{n}$ $X \in \mathbb{R}^n$

Symmetric matrix M

Positive Definite: if $w^T M w > 0$ all $w \neq 0 \Rightarrow$ pos eigenvalues

Positive Semidefinite: if $w^T M w \geq 0$ all $w \Leftrightarrow$ nonnegative eigenvalues

Indefinite: if pos & neg eigenvalue

Invertible: no zero eigenvalue

convex: $x^2 \cup$

concave: \cap

$\begin{matrix} \nearrow & \searrow \\ x & y \end{matrix}$ $x \cdot y > 0$

$\begin{matrix} \nearrow & \searrow \\ x & y \end{matrix}$ $x \cdot y = 0$

$\begin{matrix} \nearrow & \searrow \\ x & y \end{matrix}$ $x \cdot y < 0$

$\cos \theta = \frac{x \cdot y}{\|x\| \|y\|} = \frac{x^T y}{\|x\| \|y\|}$

$E[\cos \theta] = 0$ $SD(\cos \theta) = \frac{1}{\sqrt{2}}$

Probability

Bayes Rule: $P(Y=1|X) = \frac{P(X|Y=1)P(Y=1)}{P(X)}$

$P(A, B) = \sum C P(A, B, C)$

Chain Rule: $P(A, B, C) = P(A, B|C) P(C) = P(A|B, C) P(B|C) P(C)$

A conditionally independent given C: $P(A, B|C) = P(A|C) P(B|C)$

A independent of B given C: $P(A|B, C) = P(A|C)$

A, B independent: $P(A, B) = P(A) P(B)$

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

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

Matrix Derivatives

$\frac{\partial x^T a}{\partial x} = \frac{\partial a^T x}{\partial x} = a$

$\frac{\partial a^T X b}{\partial X} = ab^T$

$\frac{\partial a^T X^T b}{\partial X} = ba^T$

$\frac{\partial a^T X a}{\partial X} = \frac{\partial a^T X^T a}{\partial X} = aa^T$

$\frac{\partial X}{\partial X_{ij}} = J^{ij}$

$\frac{\partial (X A)_{ij}}{\partial X_{mn}} = \delta_{im} (A)_{nj} = (J^{mn} A)_{ij}$

$\frac{\partial (X^T A)_{ij}}{\partial X_{mn}} = \delta_{in} (A)_{mj} = (J^{nm} A)_{ij}$

$\frac{\partial}{\partial X_{ij}} \sum_{klmn} X_{kl} X_{mn} = 2 \sum_{kl} X_{kl}$

$\frac{\partial b^T X^T X c}{\partial X} = X(bc^T + cb^T)$

$\frac{\partial (Bx + b)^T C (Dx + d)}{\partial x} = B^T C (Dx + d) + D^T C^T (Bx + b)$

$\frac{\partial (X^T B X)_{kl}}{\partial X_{ij}} = \delta_{ij} (X^T B)_{kl} + \delta_{kj} (B X)_{il}$

$\frac{\partial (X^T B X)}{\partial X_{ij}} = X^T B J^{ij} + J^{ij} B X$ $(J^{ij})_{kl} = \delta_{ik} \delta_{jl}$

$\frac{\partial x^T B x}{\partial x} = (B + B^T)x$

$\frac{\partial b^T X^T D X c}{\partial X} = D^T X b c^T + D X c b^T$

$\frac{\partial}{\partial X} (Xb + c)^T D (Xb + c) = (D + D^T)(Xb + c) b^T$

$Tr(AB) = Tr(BA)$ $PSD Q = PDP^T$

$(A^T)^T = A$ $PD^2 = D^2 P^T$

$(A+B)^T = A^T + B^T$ $= U U^T$

$(AB)^T = B^T A^T$ $A = A^4 = A^{12}$

$= (P D^2 P^T)(P D^2 P^T) = P D^4 P^T$

Eigenvalues 2×2

$\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ $m \pm \sqrt{m^2 - p}$

$\frac{a+d}{2}$ $\frac{ad-bc}{2}$ $ad-bc$

MC

- rotation does not change principal comp

- dot prod $x \cdot v_i$ to get principle comp

- variance: how much change w/ diff dataset

\rightarrow high var \rightarrow overfit

- bias: how well model predicts training

- \uparrow entropy \uparrow uncertainty

- VC of halfplane is 3

- XX^T columns are eigenvectors

- cross entropy instead of MSE for sigmoid

- pruning has better test acc than stop early