

Linear Classifiers

Decision Boundary: H = {x in R^d: f(x) > 0}
f(x) = theta^T x - theta_0 = sum_{i=1}^d theta_i x_i - theta_0 > 0
Decision Boundary: H = {x in R^d: f(x) > 0}
Normal to H: For any x in H, the vector theta is normal to H at x.

Support Vector Machines
Optimize linear classification by choosing classifier maximizing margin.
Find a separating hyperplane and scale theta by margin.

Kernel Methods
Kernel trick: map data into a higher dimensional space where it is linearly separable.
Linear classification in the kernel space: f(x) = theta^T phi(x) + theta_0.

Convex Optimization and SVM
Convex problems are those in which it is possible to draw lines, and stay within the feasible region.
Support Vector Machines (SVM): Maximize the margin between two classes.

Linear Regression
Least Squares: Minimize the sum of squared residuals.
Normal Equations: (X^T X) beta = X^T y.
Ridge Regression: Add a regularization term to the normal equations.

Bayesian Linear Regression
Bayesian View of Linear Regression: Model the weights as random variables.
Posterior distribution: p(beta | X, y) proportional to exp(-lambda/2 (beta - mu)^T Sigma (beta - mu)) exp(-1/2 (y - X beta)^T Sigma^-1 (y - X beta)).

Discriminative Models
Logistic Regression: Model the probability of class membership.
Softmax Regression: Generalization of logistic regression for multiple classes.

Generative and Discriminative Models
Generative Models: Model the joint distribution of inputs and outputs.
Discriminative Models: Model the conditional distribution of the output given the input.

Parameter Estimation Methods

Maximum Likelihood
Choose parameters to maximize the likelihood of the observed data.
Log-likelihood: log L(theta) = sum_{i=1}^n log p(x_i | theta).

Bayesian Estimation
Model parameters as random variables with a prior distribution.
Posterior distribution: p(theta | X, y) proportional to p(theta) p(y | X, theta).

Multivariate Normal Distribution
Recall the multivariate Gaussian density function.
Covariance Matrix: Sigma = Cov(X, X).

Diagonal Covariance Matrices
Simplify the covariance matrix by assuming independence between features.
Lasso Regression: Penalize the L1 norm of the weights.

Properties of Multivariate Gaussians
Probability density function: p(x) = (2*pi)^(-d/2) |Sigma|^-1/2 exp(-1/2 (x - mu)^T Sigma^-1 (x - mu)).

Affine Transformations
Linear transformations of Gaussian distributions.
Change of variables: p(y) = p(x) |J|.

Linear Discriminant Analysis
Linear Discriminant Analysis (LDA): Find a linear combination of features that best separates the classes.

Regression

Linear Regression
Least Squares: Minimize the sum of squared residuals.
Ridge Regression: Add a regularization term to the normal equations.

Bayesian Linear Regression
Bayesian View of Linear Regression: Model the weights as random variables.
Posterior distribution: p(beta | X, y) proportional to exp(-lambda/2 (beta - mu)^T Sigma (beta - mu)) exp(-1/2 (y - X beta)^T Sigma^-1 (y - X beta)).

Regularization and Regression
Bias-variance trade-off: Balancing model fit and generalization.
Lasso Regression: Penalize the L1 norm of the weights.

Subset Selection
Lasso Regression: Penalize the L1 norm of the weights.
Ridge Regression: Penalize the L2 norm of the weights.

Bayesian View of Linear Regression
Bayesian View of Linear Regression: Model the weights as random variables.
Posterior distribution: p(beta | X, y) proportional to exp(-lambda/2 (beta - mu)^T Sigma (beta - mu)) exp(-1/2 (y - X beta)^T Sigma^-1 (y - X beta)).

Logistic Regression
Logistic Regression: Model the probability of class membership.
Softmax Regression: Generalization of logistic regression for multiple classes.

Gradient Ascent
Gradient Ascent: Optimize a function by moving in the direction of the gradient.
Newton's Method: Use the Hessian matrix to approximate the minimum.

Newton's Method
Newton's Method: Use the Hessian matrix to approximate the minimum.
Hessian Matrix: H = -E[second derivatives].

MLE
Maximum Likelihood Estimation (MLE): Estimate parameters by maximizing the likelihood of the observed data.

Nearest Neighbor Classification

Key idea: Store all training examples $\langle x_i, f(x_i) \rangle$
k-MN - given query instance x_q , take a vote among its k nearest neighbors (if discrete target), take mean of k-MN f values (real valued)
Non-parametric, number of parameters grows with N
Training is fast; learn complex target functions easily, don't lose intuition
but slow at query time, requires lots of storage, easily fooled by irrelevant attr
Behavior at the limit: $\lim_{k \rightarrow \infty} \text{error} = \text{error of opt}/N$ $\lim_{k \rightarrow \infty} \text{error} = \text{error of opt}$
Curse of Dimensionality
Easy to mislead MN in high dimensions (point on hyper-grid, hypercube...)
① cost more data (if $D \ll N$, make $D \ll N$)
② reduce D to get better features (compress pixels, use bag-of-words model), compare histograms using Inner/Outer Product (k-NN)
③ Use a better distance metric

Distance Metrics
Minkowski distance $D_{(p)}(x, y) = (\sum_{i=1}^n |x_i - y_i|^p)^{1/p}$
L1 norm (Manhattan / City-block distance) $D_{(1)}(x, y) = \sum_{i=1}^n |x_i - y_i|$
L2 norm (Euclidean distance) $D_{(2)}(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
0-norm (Number of nonzero elements) $D_{(0)}(x, y) = \sum_{i=1}^n \mathbb{1}_{\{x_i \neq y_i\}}$
Hamming distance for binary strings, or edit/Levenshtein distance
Mahalanobis distance $D_M(x, y) = \sqrt{(x - y)^T \Sigma^{-1} (x - y)}$
Decorrelate and renormalize features, Euclidean special case with $\Sigma = I$

Similarities
Inner $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$, symmetric, output $\in \mathbb{R}$
Cosine $\cos(x, y) = \frac{\langle x, y \rangle}{\|x\| \|y\|}$, symmetric, output $\in [-1, 1]$, normalized inner product
Corr $\text{corr}(x, y) = \frac{\langle x - \bar{x}, y - \bar{y} \rangle}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$, symmetric, output $\in [-1, 1]$, invariant to scaling

Criteria
given an instance space X , a distance metric is function $D: X \times X \rightarrow \mathbb{R}$
such that for any $x, y, z \in X$:
 $D(x, x) = 0$ (point and itself zero) $D(x, y) \geq 0$ ($x \neq y$)
 $D(x, y) = D(y, x)$ (symmetric) $D(x, z) \leq D(x, y) + D(y, z)$

Problems
Easy to mislead L2 MN in template matching, normalized cross-correlation
Euclidean distance for digit recognition sensitive to thickness, rotation

Strategies
① Augment the dataset with transform copies of the digit
② Build invariance into the feature vector (orientation histograms)
③ Build invariance into classification strategy (change dist, conv MN)
Clustering \rightarrow set k by cross-validation, form prototypic, remove noisy data

Reducing Computational Cost

k-d tree
Binary tree data structure for organizing set of points in k-dim space
Internal node associated with axis-aligned hyper-plane, splits into subtrees
Dimensions with high variance chosen first
Position of splitting hyperplane taken as median of projected points
Query for MN may require backtracking
Indexing local features: approximate MN search
Best-Bin First (BBF), variant of k-d trees, uses priority queue
to examine most promising branches first
Locality-sensitive Hashing (LSH): randomized hashing technique to map
similar points to same bin with high probability
we want a high degree of collisions and hence each hash function
must satisfy $\Pr[h(x) = h(x')] = \frac{1}{\pi}$

Decision Trees

Utilize labels y (unless k-d trees) in performing classification
Internal nodes test value of features x_j and branch accordingly
Leaf nodes specify a class $h(x)$
Allow for great feature selection, interpretable results, decision
regions look similar to L2MN/linear classifiers

Learning/Tree Construction
Choosing Best Attribute: 1-step greedy (local) heuristic, does not account for "progress"
Split on entropy $H(x) = -\sum_{i=1}^n p_i \log p_i$
Information gain = entropy(parent) - average entropy(children)
Mutual information is a better measure (A vs. absolute error)
Non-binary features: Multiple discrete values - multiply split (normalise into gain)
one vs. all split (binary), group into two disjoint subsets
Real-valued features - consider threshold split at each stage

Unknown Attributes
If node n tests A , assign most common value of A among other examples
Assign most common value of A among other examples γ same target t
Assign prob p_i to each value v_i of A , assign frac p_i to each descendant

Classification vs. Regression Trees
 $x \rightarrow \{0, 1, 2\}$ vs. $x \rightarrow y \in \mathbb{R}$, minimize sum of variances
 $\sum_{i \in \{0, 1, 2\}} (y_i - \hat{y}_i)^2 + \sum_{i \in \{0, 1, 2\}} (y_i - \hat{y}_i)^2$

Overfitting
Don't wait until absolute purity (no more happen - noise, data slowness)
low data points / depth, deep trees slow and huge
Stop splitting earlier and make leaf node with (average/majority class)
Early termination based on max tree depth, min # points at node,
tree complexity penalty, validation error monitoring
Reduced-Error Pruning
Split data into training and validation set, do until harmful:
① Evaluate impact of pruning each possible node on validation
② Greedily remove node that most improves validation accuracy

Random Forests

Ensemble Methods combine several "weak learner" models by
Averaging (randomize each model), Bagging (randomize dataset given to model),
Boosting (specialize each model for subset of examples), Random Forests (4th, 3rd)
Bagging still leads to correlated trees, limit to subset of p splits
Use each result as a vote towards overall classification

Boosting

Define classifier using additive model $f(x) = w_1 f_1(x) + w_2 f_2(x) + \dots$
Example: Each data point starts with weight 1. Draw weak learner, then
reweight misclassified points by $w_t = \frac{1}{2} \exp(-y f_t(x))$
AdaBoost Algorithm
given m examples $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in \mathbb{R}^n, y_i \in \{0, 1\}$, initialize $D_0(x) = \frac{1}{m}$
for $t = 1$ to T :
① Train learner h_t with min error $\epsilon_t = \min_{h \in H} \sum_{i=1}^m h(x_i) y_i$
② Compute the proportional weight $\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t}$
③ For each example (x_i, y_i) $D_t(x_i) = \frac{1}{m} \exp(-\sum_{s=1}^t \alpha_s y_s h_s(x_i))$
Output: $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t h_t(x))$

Neural Networks

Input nodes, connected to hidden layers, then output nodes
Each node outputs a linear combination of its inputs, by activation
Two-layer network equivalent to single matching, use more complicated
learning scheme to extract better features
Use one-vs-all classification for multi-class

Learning
Randomly perturbing weights similar to reinforcement learning
but can be very inefficient, changing just one is not effective
Backpropagation uses error derivatives with respect to
the hidden activities
Numerical and analytic gradient slow to calculate, Damp
similar to cached analytic gradient
Using squared error, constraint enforce sum to one constraint
Use Softmax $y_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$ $z_i = y_i(1 - y_i)$, $c = -\sum_{i=1}^K \log y_i$

Combine ideas of logistic regression, template matching/MN,
decision focus/ensembles of weak learners, stochastic gradient,
and feature learning
Apply regularization using L1, L2, both, max norm constraint

Conditional Neural Networks
Alternating convolutional and pooling layers
Conclusion: Apply element wise condition using filters to produce
map layered output
Note: locally connected (not fully connected), and sparse
weights across each depth slice
Pooling: Filter responses at each location to extract features
Max (take max response), Average, L2 (sqrt), L2 over features

Unsupervised Learning

without labels, discover something interesting about x (data), discover
some kind of underlying structure (existing data may also be using)
Clustering (partitioning into tight clusters), Dimensionality Reduction (discover
low-D manifold, good features), Mode Seeking (discover frequent patterns)

Clustering
Applications in: image segmentation, super pixels, market segmentation...
Agglomerative: Starting with n points, keep merging closest two clusters
Similarity metric: dist $f(x, y)$
Single link: $\text{dist}(S_1, S_2) = \min_{i \in S_1, j \in S_2} \text{dist}(i, j)$ (sensitive to outliers)
Complete link: $\text{dist}(S_1, S_2) = \max_{i \in S_1, j \in S_2} \text{dist}(i, j)$ (compact, sensitive to outliers)
Average link: $\text{dist}(S_1, S_2) = \frac{\sum_{i \in S_1, j \in S_2} \text{dist}(i, j)}{|S_1| |S_2|}$ (in between, group decision)
Centroid link: $\text{dist}(S_1, S_2) = \| \bar{x}_1 - \bar{x}_2 \|^2$
Divisive: Starting with single cluster, recursively divide until each
cluster only has one object or minimal overlap

k-means
Specify there are k means and minimize distance of points to mean
Number of distinct assignments exponential in k and N
Use greedy algorithm alternating between mean/cluster assignment
① Randomly initialize cluster centroids μ_1, \dots, μ_k
② Repeat until convergence:
- For every point i , assign to closest cluster k : $c_i = \text{argmin}_k \|x_i - \mu_k\|^2$
- For every cluster k , compute mean $\mu_k = \sum_{i \in c_k} x_i / |c_k|$

Repeated coordinate descent on $J(\mu) = \sum_{i=1}^n \|x_i - \mu_{c_i}\|^2$
Avoid bad local minima by better initialization: choose one point among
data randomly, then others by weighted probability distribution $D(x)$
If only given distance matrix, use k-medoids algorithm
For each point x , find closest entity m and move towards
 x a little; repeat
Works well if clusters spherical, well separated, of similar volume
and have similar number of points

Vector Quantization involves finding closest clusters (often by
k-means) and creating objects to closest neighbor

Mixture Models

$P(x) = \sum_{c=1}^K \alpha_c \mathcal{N}(x; \mu_c, \Sigma_c)$ Objective function, log likelihood of data
Clustering might use K-Cluster, $\mathcal{N}(x; \mu, \Sigma)$

EM Algorithm

Initialize parameters μ, Σ manually into, repeat:
E: compute expected values of unobserved variables assuming μ, Σ
M: compute new parameters to maximize prob of data
For Gaussians,
Initialization: Choose means at random
E: For all examples x_i :
 $P(\mu, \Sigma) = \prod_{i=1}^n P(x_i; \mu, \Sigma) = \prod_{i=1}^n \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x_i - \mu)^T \Sigma^{-1} (x_i - \mu))$
M: For all components c :
 $P(\mu, \Sigma) = \prod_{i=1}^n \sum_{c=1}^K \alpha_c P(x_i; \mu_c, \Sigma_c)$, $\mu_c = \frac{\sum_{i=1}^n \alpha_c x_i}{\sum_{i=1}^n \alpha_c}$, $\Sigma_c = \frac{\sum_{i=1}^n \alpha_c (x_i - \mu_c)(x_i - \mu_c)^T}{\sum_{i=1}^n \alpha_c}$

Non-Parametric Density Estimation
Histogram method: use fixed bin widths

Kernel Density Estimation
 $P(x) = \frac{1}{n} \sum_{i=1}^n K(x - x_i)$ function of finite number of data points
• Normalized: $\int_{-\infty}^{\infty} K(x) dx = 1$, symmetric: $K(x) = K(-x)$, Exp. weight decay:
Epanechnikov kernel: $K_E(x) = \frac{3}{4} (1 - x^2) \mathbb{1}_{\{|x| \leq 1\}}$, $\lim_{h \rightarrow 0} \mathbb{1}_{\{|x| \leq 1\}}(x) = 0$
Uniform kernel: $K_U(x) = \begin{cases} \frac{1}{2} & \text{otherwise} \\ 0 & \text{otherwise} \end{cases}$
Normal kernel: $K_N(x) = \frac{1}{\sqrt{2\pi}} \exp(-\frac{1}{2} x^2)$

Mode-seeking

"Association Rules", "Frequent Itemsets", "Market Analysis"
Association Rules
Rule $X \rightarrow Y$: "support \rightarrow head (support, confidence)"
Support: $\text{supp}(X, Y) = \text{frequency}$ (i.e. $P(X, Y)$)
Confidence: $\text{conf}(X \rightarrow Y) = \frac{\text{supp}(X, Y)}{\text{supp}(X)}$ (i.e. $P(Y|X)$)
Lift $L(X \rightarrow Y) = \frac{\text{supp}(X, Y)}{\text{supp}(X) \text{supp}(Y)}$

Principle Component Analysis

Given a set of feature vectors x_i , we can represent as pattern matrix
Ex: Measurement vectors, digital images, text documents, user rating data
Problems \rightarrow too large/complexity, missing/noisy entries, lack of structure
Extract Latent Structure of matrix using factorization: $A \approx L \cdot R$
Can perform dimensionality reduction, compressing and keeping
only directions of highest variance
PCA: ① Subtract mean from every point
② (Sometimes) scale each dimension by its variance
③ Compute covariance matrix $S = x^T x$
④ Compute k largest eigenvalues of $S = V D V^T$

Singular Value Decomposition

① Subtract mean from each data point
② Decompose $X = U S V^T$, $X^T X = V S^T V^T = V D V^T$, $X X^T = U S U^T = U D U^T$
③ U, S matrix provides coefficients: $x_i = U_{i1} S_{11}^{1/2} + U_{i2} S_{22}^{1/2} + \dots$