## Matrix Calculus

$\nabla(x^T A x) = (A + A^T)x, \ 2Ax$

$\nabla(x^T A y) = A^T y, \ \nabla(y^T A x) = A^T y$

$\nabla(x^T x) = 2x, \ \|x\|_2^2 = x^T x$

$\|Ax\|_2^2 = x^T A^T A x, \ \|y - Ax\|_2^2 = (y - Ax)^T(y - Ax)$

$= y^T y - y^T A x - x^T A^T y + x^T A^T A x$

$\nabla = -2x^T A^T y + x^T A^T A x$

**Generative Model:** $P(X,Y) = P(Y)P(X|Y)$. Estimate $P(Y), P(X|Y)$, use Bayes.

**Discriminative:** $P(X,Y) = P(X)P(Y|X)$. Estimate $P(Y|X)$, use in $\hat{y} = P(Y|X)$.

**Linear Regression:** $RSS(\beta) = \frac{1}{2}\|X\beta - Y\|^2$

$\nabla RSS = \nabla(\frac{1}{2}\beta^T X^T X\beta - y^T X\beta - \frac{1}{2}y^T y) = 0$

$X^T X\beta - X^T y = 0 \Rightarrow \beta = (X^T X)^{-1}X^T y.$ **Normal.**

$E(\hat\beta) = \beta. \ Var(\hat\beta) = (\sigma^2\beta) = \sigma^2(X^T X)^{-1}$

**Ridge Regression:** $\hat\beta = \min(\|y - X\beta\|^2 + \lambda\|\beta\|^2)$

$\hat\beta = (X^T X + \lambda I)^{-1}X^T y$. L2 norm. Gaussian

**Lasso:** $\hat\beta = \min(\|y - X\beta\|^2 + \lambda\|\beta\|)$. L1-norm. Sparse coefficients (set to 0). Laplace prior $= \exp(-\lambda|\beta|)$

**Logistic Regression:** For Gaussian class conditionals: $P(X|Y=1) = N(\mu_1, \Sigma)$.

Posterior in logodds: $P(Y=1|X) = (\exp(-\theta^T x) + 1)^{-1}$

Score log odds $\log\left(\frac{P(Y=1|X)}{P(Y=0|X)}\right) = \ldots$

$= \beta_0 + \beta^T x. \ \beta = \Sigma^{-1}(\mu_1 - \mu_0). \ \beta_0 = \frac{1}{2}(\mu_0^T\Sigma^{-1}\mu_0 - \mu_1^T\Sigma^{-1}\mu_1) + \log\frac{P(Y=1)}{P(Y=0)}. \ \log\frac{P}{1-P} = \log\frac{P(Y=1|X)}{P(Y=0|X)}$

**Linear Discriminant Analysis:** Choose class that maximizes $\delta_k(x)$ - class are distributed MVGaussian w/ common $\Sigma$.

$\hat\pi_k = Pr(Y=k) = \frac{n_k}{n}, \ \hat\mu_k = E(X|Y=k) = \frac{1}{n_k}\sum x_i.$

$\hat\Sigma = \frac{1}{n}\sum_k\sum_{i:y_i=k}(x_i-\mu_k)(x_i-\mu_k)' = Var(X|Y=k)$

**D function:** $\delta_k(x) = \hat\mu_k\hat\Sigma^{-1}x - \frac{1}{2}\hat\mu_k\hat\Sigma^{-1}\hat\mu_k + \log\hat\pi_k$

**KKT:** (1) Primal Feasibility: $f_i(x) \le 0 \ \forall i$

(2) Dual Feasibility: $\lambda \ge 0$ possible.

(3) Complementary Slackness: $\lambda_i f_i(x) = 0.$ Thus since $\lambda \ge 0, f_i(x) = 0$ or $f_i(x) \le 0, \lambda = 0.$

(4) Stationarity: $\nabla f_0(x) + \sum_i \lambda_i\nabla f_i(x) = 0$

**kNN Algorithm:** Given $x_q$, locate k-nn, distance function, take the vote among k. Not good in high-D; if $D > n$:

(1) Get more data, increase $n$ ... Set k by CV

(2) Reduce features, decrease D

(3) Use better distance metric. KNN optimal as $n\to\infty$, $k\to\infty$, $\frac{k}{n}\to 0$. Error $\le 2 \cdot$Bayes.

**Mahalanobis Distance:** $D_m(x,y) = (x-y)^T\Sigma^{-1}(x-y)$

**Transformation Invariance:** Augment dataset w/ slight transformations of training data

---

$J(\theta) = \sum_i(-y_i(\theta^T x_i))$

$\nabla J_i(\theta) = -y_i x_i \text{ if } (y_i x_i < 0)$

**Algorithm:** $\theta = \theta - y_i x_i$ if misclassified.

**Kernels:** Replace inner product of feature vectors $\phi(x)$ w/ Kernel $K(x_i, x_j)$. Examples: Polynomial $K(x,z) = (1 + x^Tz)^m$, RBF $K(x,z) = \exp(-\frac{1}{2}\|x-z\|^2)$

**MLE:** max w/o form prior $P(x_1, x_2, \ldots x_n|\theta)$

**MAP:** $P(x_1, \ldots x_n|\theta)\pi(\theta)$

$\pi =$ Gaussian $\to$ **Ridge**

$\pi =$ Laplace $= \exp(-\lambda|\beta|) \to$ **Lasso**

**Spectral Theorem:** Symm, real, $A \in R^{n\times n}$, northonorm eigenvec $(v_1, \ldots v_n)$, $n$ eigenvalues $(\lambda_1, \ldots \lambda_n)$. Can write $AU = U\Lambda, A = U\Lambda U'$, $U = [v_1, v_2, \ldots v_n], \Lambda = diag[\lambda_1, \ldots \lambda_n]$

**Gaussian MLE:** $\mu = \frac{1}{n}\sum_i x_i.$

$\Sigma = \frac{1}{n}\sum_i(x_i - \hat\mu)(x_i-\hat\mu)', \ \sigma^2 = \frac{1}{n}\sum_i(x_i-\mu)^2$

**Subset selection:** $2^p$ diff subset of features. Greedy search via CV error. Forward/backward stepwise.

**MLE for Logistic:** $\ell(\beta) = \log P(y_1, \ldots y_n|x_1, \ldots x_n, \beta)$

$= \sum_i y_i\log\mu_i(\beta) + (1-y_i)(1 - \mu_i(\beta)), \ \mu_i(\beta) = P(Y=1|X=x_i, \beta)$

$= (1 + \exp(-\beta^T x_i))^{-1}. \ \nabla_\beta\mu_i(\beta) = \mu_i(\beta)(1 - \mu_i(\beta))x_i.$

**First Deriv:** $\nabla_\beta\ell(\beta) = \sum_i\left(\frac{y_i}{\mu_i(\beta)} - \frac{1-y_i}{1-\mu_i(\beta)}\right)\nabla_\beta\mu_i(\beta) = \sum_i(y_i - \mu_i(\beta))x_i.$ **Second:** $\nabla_\beta^2\ell(\beta) = -\sum_i\mu_i(\beta)(1 - \mu_i(\beta))x_i.$

Set $\nabla_\beta\ell(\beta) = 0 \Rightarrow \sum_i y_i x_i = \sum_i\mu_i(\beta)x_i.$ **No closed form.**

**Gradient ascent:** $\beta^{t+1} = \beta^t + \eta(\nabla_\beta\ell(\beta^{(t)})) = \beta^t + \eta\sum_i(y_i - \mu_i(\beta^t)).$ $O(np)$ per iter.

**Stochastic Gradient:** Use single value, $O(p)$/iter

**Linear Reg:** $\nabla_\beta RSS(\beta) = X^T X\beta - X^T y = \sum_i x_i(x_i'\beta - y_i) \Rightarrow \beta^{(t+1)} = \beta^t + \eta x_{i,t}(y_{i,t} - x_{i,t}^T\beta^t).$

**Logistic Reg:** $\beta^{t+1} = \beta^t + \eta(y_{i,t} - \mu_{i,t}(\beta^t))x_{i,t}.$

**Lagrange Multipliers:** min/max $f_0(x)$ s.t. $f_i(x) \ge 0$

Replace w/ $L(x,\lambda) = f_0(x) + \lambda f_i(x)$. Take PD's w.r.t. $x_i$'s, $\lambda$, set to 0, solve.

**Optimization:** primal $= p^* = \min_x \max_{\lambda\ge 0} L(x,\lambda)$, $g(\lambda) = \min_x L(x,\lambda)$, dual $= d^* = \max_{\lambda\ge 0}\min_x L(x,\lambda)$

**SVM hard margin dual:** $\max_\alpha \sum_i\alpha_i - \frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j x_i^T x_j, \ 0 \le \alpha_i$

**SVM soft margin dual:** $\max_\alpha \sum_i\alpha_i - \frac{1}{2}\sum_{i,j}\alpha_i\alpha_j y_i y_j x_i^T x_j, \ 0 \le \alpha_i \le \frac{C}{n}$

corresponding primal: $\min_\theta \frac{1}{2}\|\theta\|^2 + C\sum_i\xi_i : 1 - y_i\theta^T x_i - \xi_i \le 0, \xi \ge 0$

**k-d tree:** Binary tree for organizing points in k-dim space. Each node associated w/ axis-aligned hyperplane splitting into two subtrees (choose dim w/ highest var first). Position chosen @ median to balance tree. Best bin first uses priority queue to branch. **Locality sensitive hashing:** use hash to map similar points to same bin. Project data $n^T x$ w/ bits.

---

**Softmax:** ...

**Decision Theory:** Loss $\ell(\hat y, y) = $ cost of ...

**Risk:** Expected loss $R(f, y) = \mathbb{E}[\ell(f(x), y)]$

**Bayes:** ... $f^*(x) = \ldots$

**Regression Risk:** $E[(f(x) - y)^2|X]$; Optimal Risk: $R^*$

**Excess Risk:** $R(f) - R^* = E[(f(X) - f^*(X))^2] \ge 0$

**Bias-Variance:** $R(f) = E[(f(X) - E(y|X))^2] + E[(E(y|X) - Y)^2]$

**Gaussian Generative:** $P(X|Y=k) = N(\mu, \Sigma) \Rightarrow$

**Logistic Discriminative:** $P(Y=1|X) = (1 + \exp(-\theta^T\ldots))^{-1}$

**Two class case:** $\theta_1 = \frac{\mu_0 - \mu_1}{\sigma^2}, \theta_0 = \frac{\mu_1^2 - \mu_0^2}{2\sigma^2} - \log\frac{P(Y=0)}{P(Y=1)}$

**MVGaussian:** $(2\pi)^{-n/2}|\Sigma|^{-1/2}\exp(-\frac{1}{2}(x-\mu)^T\Sigma^{-1}(x-\mu))$

$X + Y \sim N(\mu_x + \mu_y, \Sigma_x + \Sigma_y); \ if \ A = \ldots, Y = N(A\mu, A\Sigma A^T)$

**Covariance Matrix:** $\Sigma = E(X - \mu)(X - \mu)^T$ Symm, PSD.

$Ax = \lambda x, \lambda = $ eigenval, $x = $ eigenvec. Solve $\det(A - \lambda I) = 0$ to get $\lambda$'s, solve $(A - \lambda I)[\ ,\ ] = $ for diagonal covariance. $((\sigma(x_i, x_j) = 0, i\ne j)$

**Super level set:** $\{x \in R^n : (x-\mu)^T\Sigma^{-1}(x-\mu) \le r^2\}$ ... $r^2$ correspond to eigenvalue. For non-diagonal, $\Sigma_d$ is rotated version of diagonal $\Sigma_x : \Sigma_x = Q\Sigma_d Q'$. Volume of ellipsoids proportional to $\prod_i\sigma_i^2 = \sqrt{|\Sigma|}$

**Newton's Method:** iterative reweighted least squares

$\nabla f(x_t) = x_{t+1} = x_t - f''(x)^{-1}\ldots$

**Newton-Raphson for Logistic:**

$\ell(\beta) = \sum_i y_i\log\mu_i(\beta) + (1-y_i)\log(\ldots)$

$\nabla_\beta\ell(\beta) = \sum_i(y_i - \mu_i(\beta))x_i, x = X(y - \mu)$

$\nabla_\beta^2\ell(\beta) = \sum_i -\mu_i(\beta)(1-\mu_i(\beta))x_i x' = -X diag(\mu(1-\mu))X.$

$\beta^{t+1} = \beta^t - [\nabla^2\ell(\beta)]^{-1}\nabla\ell(\beta^t) = \beta^t + [X' diag(\mu(1-\mu))X]^{-1}X'\ldots$

**Using SVM kernel:** classify point $x$ by $\theta^T\phi(x) = \sum_j\alpha_j y_j\langle\phi(x_j), \phi(x)\rangle = \sum_j\alpha_j y_j K(x_j, x)$ where $\alpha_j \ne 0$ only for support vectors.

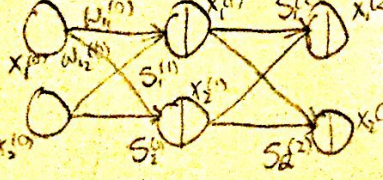**Quad Kernel Features:** $[1, x_1^2, x_2^2, \sqrt 2 x_1, \sqrt 2 x_2, \ldots]$

**Decision Tree Alg:** GrowTree(S):
- if $(y = 0$ for all $(x,y) \in S)$ return new leaf(0)
- if $(y = 1$ for all $(x,y) \in S)$ return new leaf(1)
- choose best attribute $x_j$
  - $S_0$: all $(x,y) \in S$ w/ $x_j = 0$;
  - $S_1$: all $(x,y) \in S$ w/ $x_j = 1$
- return new node $(x_j, GrowTree(S_0), GrowTree(S_1))$

**ChooseBestAttribute(S):** L = ...
- choose j to minimize $J_j$:
  - $S_0$: all $(x,y)$ w/ $x_j = 0$, $S_1$: all $(x,y)$ w/ $x_j = 1$
  - $y_0 = $ Most common $y$ in $S_0$, $y_1 = $ Most common $y$ in $S_1$
  - $J_0 = $ \#$\{(x,y) \in S_0 : y \ne y_0\}$, $J_1 = \#\{(x,y) \in S_1 : y \ne y_1\}$
  - $J_j = J_0 + J_1$
- return ...

Surprise: $S(V=v) = -\log(V=v)$
$P(V=1) = \infty$ surprise; $P(V=1) = 0$ surprise
Entropy = average surprise
$H = \Sigma - p_i \log p_i$. For group w/
all one class, entropy $= 0$.
For group w/ 50% split, entropy
$H = 0.5 \log \frac{1}{2} - \frac{1}{2}\log\frac{1}{2} = 1$.
Information Gain: entropy (parent)
$- [$average entropy (children)$]$.
Want to make split that
maximizes info gain.
Choose Best Attribute based on IG.
Classification tree: $x \to \{0, 1\}$
Regression tree: $x \to y \in R$
Instead of minimizing entropy,
minimize sum of variance at each
split: $\sum_{x_i \in R_1}(y_i - \hat{y}_{R_1})^2 + \sum_{x_i \in R_2}(y_i - \hat{y}_{R_2})^2$

To prevent overfitting:
1. Impose stopping criteria:
- Max tree depth
- Min # points @ node
- Tree complexity penalty
- Validation error monitor
Stop splitting earlier &
make leaf node:
avg/majority of
$y$'s @ leaf
2. Reduced Error Pruning:
Split into training &
validation data.
- Do until further
pruning harmful:
1. Evaluate impact on
validation set of
pruning node & tree below
2. Greedily remove nodes w/
most improvement

Ensemble Methods
Combine several models
- Bagging, Averaging
- Random Forest, Boosting
Bagging: Generate bootstrap
replicates of training set
by sampling w/ replacement.
Learn one model on each.
Combine by uniform voting.
Random Forest: Bagging
on dataset, also random
subset features for each tree
ADABoost: Adapt weights for each tree. As error grows,
"importance" of tree (adaptive weight) goes down. Error
decreases, tree importance. For each example,
boost weight if incorrectly predicted, else decrease.
$\varepsilon = $ sum of misclassifications $\alpha_t = \frac{1}{2}\ln(\frac{1-\varepsilon}{\varepsilon})$
$D_n^{(t)} = \frac{1}{Z} \times D_n^{(t-1)} \times \{^{e^{-\alpha_t} \; y_n = h}_{e^{\alpha_t} \; y_n \neq h}}$. $Z$ is normalization factor

Boosting: Build strong learner as
combo of weak learners
$L(x) = a_1 h_1(x) + \dots + a_n h_n(x)$
ADABOOST $(\Psi, D, K)$
1. $d^{(0)} = [\frac{1}{n}, \dots, \frac{1}{n}]$ // uniform importance over each example
2. for $k = 1, 2, \dots K$ do
3. $f^{(k)} = W(D, d^{(k-1)})$ // train $k$th base
4. $\hat{\varepsilon}^{(k)} = \sum_n d_n^{(k-1)}$ // misprediction
5. $\hat{\varepsilon} = $ 2nd $\Sigma$ // weighted err
6. $\alpha^{(k)} = \frac{1}{2}\log((1-\hat{\varepsilon})/\hat{\varepsilon})$ // adaptive param
7. $d_n^{(k)} = \frac{1}{Z} d_n^{(k-1)} \exp(-\alpha^{(k)} y_n f^{(k)})$
8. return $f(x) = \text{sgn}(\sum_k \alpha^{(k)} f^{(k)}(x))$ // return required classifier

## Neural Networks



$g(s) = \frac{1}{1 + e^{-s}}$

$s_i^{(k)} = \sum_j x_j^{(k-1)} w_{ij}$
$x_i^{(k)} = g(s_i^{(k)})$
$J(w) = \frac{1}{2}\sum_i (x_i^{(2)} - y_i)^2$

$\frac{\partial J(w)}{\partial w_{ij}^k} = \frac{\partial J(w)}{\partial x_j^{(k)}}\frac{\partial x_j^{(k)}}{\partial s_j^{(k)}}\frac{\partial s_j^{(k)}}{\partial w_{ij}}$
$\frac{\partial s_j^{(k)}}{\partial w_{ij}} = x_i; \quad \frac{\partial x_j^{(k)}}{\partial s_j^{(k)}} = g(s_j)(1-g(s_j))$
$\frac{\partial J(w)}{\partial x_j^{(out)}} = (x_j^{(out)} - y_j)$ at output

$\frac{\partial J(w)}{\partial x_j^{(k)}} = \sum_i (\frac{\partial J(w)}{\partial x_i^{(k)}}\frac{\partial x_i^{(k)}}{\partial s_i})$
$\frac{\partial J(w)}{\partial w_{ij}} = \delta_j x_i$ where
$\delta_j = \frac{\partial J(w)}{\partial x_j}\frac{\partial x_j}{\partial s_j}$

$\{(y_j - x_j)g(s_j)(1-g(s_j))$ output
$\{\sum_i \delta_i w_{ji}; g(s_j)(1-g(s_j))$
@ inner. where $i$ is node
connected in next layer

$w_{ij} = w_{ij} - \eta \frac{\partial J(w)}{\partial w_{ij}}$

## Clustering
Hierarchical: Produce tree/dendrogram
Agglomerative: bottom-up (most used). Start with points, @
each step merge closest two
Divisive: top-down.
Need to define distance function

Flat clustering: Partition data into $k$ groups, iteratively
reallocate observations to closest group to minimize loss.
K-means: Initialize cluster centroids $\mu_1, \dots, \mu_k$ randomly.
Repeat until convergence: ① Assign every point to closest
cluster $k$: $C_i = \text{argmin}_k \|x_i - \mu_k\|^2$. ② For every cluster $k$, recompute
its mean $\mu_k = (\sum_{C_i = k} x_i)/(\sum_{C_i = k} 1)$
Repeat to minimize objective function: $J(C, \mu) = \sum_i \|x_i - \mu_{C_i}\|^2$.
K-means will converge to local minima, use random restarts/better init.
K-means++: Choose one center uniformly @ random from data,
compute $D(x) = $ distance to nearest center. Choose one new data
point as center using prob. distribution $\propto D(x)^2$. Repeat to
get $k$ centers, then do $k$-means.
K-means works well if clusters spherical/well separated/similar vol/# pts.
Gaussian mixture model clustering: $P(x) = \sum_i P(c_i) P(x|c_i)$. Gaussian prior
Meanshift.

PCA: Given design matrix $X$
1. Subtract mean from each point
2. (Sometimes) scale by SD
3. Compute $S = X^TX = $ cov mat $= \Sigma$
4. Compute eigenvec. $S = VDV^T$
SVD: Decompose $X = USV^T$
$X^TX = VSSV^T = VDV^T$, $XX^T = USSU^T$



Dimensionality reduction: Taking $k$
eigenvectors amounts to reducing
dimensions of $X$ into $k$ dim w/
highest variance (project data points
onto vector). Approx matrix to
$X_i = U_{i,1} S_{11} V_1 + U_{i,2} S_{22} V_2 + \dots$

$XX^T$, $X^TX$ share same eigenvalues
TSVD $X$: Calculate $\det(XX^T - \lambda I)$
to get $\lambda$, then $(XX^T - \lambda I)/:: = 0$
for $U$ vec, $X^TX$ for $V$ vec.