

Linear Classifiers

f(x) = theta^T x - theta\_0 = sum\_{i=1}^n theta\_i x\_i + theta\_0, y = { 1 if f(x) >= 0, -1 if f(x) < 0
Decision boundary H = {x in R^d: theta^T x + theta\_0 = 0}
Perceptron algorithm
Input: (x\_1, y\_1), ..., (x\_n, y\_n) in R^d x {+1, -1}
While some y\_i != sign(theta^T x\_i + theta\_0)
 pick some misclassified (x\_i, y\_i)
 theta = theta + y\_i x\_i (update)
return theta

Since theta = sum\_{i=1}^n y\_i x\_i, theta\_0 = sum\_{i=1}^n y\_i x\_i^T x + sum\_{i=1}^n y\_i theta\_0
data in any inner product space
View as stochastic gradient descent, where J(theta) = sum\_{i=1}^n (y\_i (theta^T x\_i - 1))^2
Converges only if data linearly separable, time to converge depends on margin, solution depends on the starting point

Support Vector Machines

Optimize linear classification by choosing the classifier minimizing ||theta||
Find separating hyperplane theta^T x + theta\_0 = 0 by V\_S -> min ||theta||
max\_{i=1,2} min ||theta||^T x\_i - theta\_0 >= 1 (quadratic program)
Support vectors are the points that satisfy constraints (compressed)
theta = sum\_{i=1}^n y\_i x\_i (x\_i^T x\_i = 1), ||theta|| = sum\_{i=1}^n y\_i x\_i^T x\_i (just inner products)
Hard Margin SVM
min ||theta||^2 s.t. y\_i theta^T x\_i >= 1 (i=1, ..., n)
Soft Margin SVM
Relax inequalities y\_i theta^T x\_i >= 1 by introducing slack variables epsilon\_i
min ||theta||^2 + C sum\_{i=1}^n epsilon\_i s.t. y\_i theta^T x\_i >= 1 - epsilon\_i, epsilon\_i >= 0 (epsilon\_i = (1 - y\_i theta^T x\_i)\_+)
Regularization: small C -> small ||theta||, large C -> small misclassification

Feature Selection

Linear classifier: f(x) = theta^T x, quadratic: f(x) = ||x - c||^2 - r^2
Quadratic is simply the linear classifier with features phi(x)
Want to balance complexity of feature with degree of fit
Kernels
We can work in inner product space, so (x\_i, x\_j) or phi(x\_i)^T phi(x\_j)
solve optimization problem with K(x\_i, x\_j) by defining phi for theta = sum\_{i=1}^n y\_i phi(x\_i)
since we can then compute classification by theta^T phi(x) = sum\_{i=1}^n y\_i phi(x\_i)^T phi(x)
Degree-m polynomial kernel given by K\_m(x, x') = (1 + x^T x')^m
Radial basis function kernel given by K\_rbf(x, x') = exp(-gamma ||x - x'||^2)
Scaling parameter gamma affects smoothness, small -> long margin -> few points influence
Kernels add modularity to classifier training since same optimization procedure can be used (K(x\_i, x\_j) needed)
Degree-m polynomial on R^d is inner product with (m+1) features

Decision Theory

The Prediction Problem: Given a training set (x\_1, y\_1), ..., (x\_n, y\_n) choose a function f: X -> Y s.t. that for subset S of X, f(x) is a good predictor of y
Loss function L: Y x Y -> R where L(y, y') is cost of predicting y as y'
Can define identical or asymmetric loss functions
Assume that (x\_i, y\_i) and (x\_j, y\_j) are chosen iid from some pair X x Y
A good prediction -> small expected loss, choose f to minimize R(f) = E[L(f(x), y)]
Risk is misclassification probability: R(f) = E[L(f(x), y)] = E[1 - I(f(x) = y)] = P(f(x) != y)
Two cases: R(f) = E[L(f(x), y)] = E[|f(x) - y|] = E[|f(x) - y|] = E[|f(x) - y|]
Bayes Decision Rule: f\*(x) = argmin\_{y in Y} P(y = y | x)
Risk is minimized when f is Bayes optimal (Bayes) rule: R(f) = E[L(f(x), y)]
Error risk R(f) - R(f\*) = E[L(f(x), y)] - E[L(f\*(x), y)] = E[L(f(x), y) - L(f\*(x), y)]
Error risk of decision rule can be quantified in terms of a certain distance from f\*

Risk in Regression

Risk is expected squared error R(f) = E[(f(x) - y)^2] = E[E[(f(x) - y)^2 | x]]
Minimize the conditional expectation of the loss E[L(f(x), y) | x] by f\*(x) = E[y | x]
Bias-Variance Decomposition
R(f) = E[E[(f(x) - y)^2 | x]] = E[E[(f(x) - E[y | x])^2 | x]] + E[E[(E[y | x] - y)^2 | x]]
Minimizing R(f) = E[E[(f(x) - y)^2 | x]] = E[E[(f(x) - E[y | x])^2 | x]] + E[E[(E[y | x] - y)^2 | x]]
Using random training data to choose f, we would like E[R(f)] to be small
E[R(f)] = E[E[(f(x) - y)^2 | x]] = E[E[(f(x) - E[y | x])^2 | x]] + E[E[(E[y | x] - y)^2 | x]]
Three approaches to choosing classifiers
1. Estimate a generative model by P(x) and P(y|x), use Bayes theorem
P(y = +1 | x) = P(y = +1) P(x | y = +1) / (P(y = +1) P(x | y = +1) + P(y = -1) P(x | y = -1))
2. Estimate a discriminative model by P(y|x) and use Bayes rule as actual P(y|x)
Use as plug-in estimator (but only require accuracy around P(y = +1 | x))
3. Cross a classifier (heuristic) based on optimization of criterion

Generative and Discriminative Models

Recall P(y = +1) = E[P(x | y = +1)] = E[E[x | y = +1]] = E[E[x | y = +1]] = E[E[x | y = +1]]
For Gaussian class conditional densities P(x | y = +1) and P(x | y = -1) (with the same variance), the posterior probability is logistic: P(y = +1 | x) = 1 / (1 + exp(-x^T (mu\_+ - mu\_-)))
Suppose class conditional distributions are Gaussians:
P(x | y = +1) = N(x | mu\_+, Sigma)
P(x | y = -1) = N(x | mu\_-, Sigma)
Consider the 2-dimensional case x = [x\_1, x\_2]^T, mu\_+ = [mu\_+1, mu\_+2]^T, mu\_- = [mu\_-1, mu\_-2]^T, Sigma = Sigma
(apply Bayes rule) P(y = +1 | x) = (P(x | y = +1) P(y = +1)) / (P(x | y = +1) P(y = +1) + P(x | y = -1) P(y = -1))
Simplify and take log: log P(y = +1 | x) = log(P(x | y = +1) P(y = +1)) - log(P(x | y = +1) P(y = +1) + P(x | y = -1) P(y = -1))
exp(-theta^T x - theta\_0) = exp(-theta^T x - theta\_0) = exp(-theta^T x - theta\_0) = exp(-theta^T x - theta\_0)
Inv-dimensions: theta = Sigma^{-1} (mu\_+ - mu\_-)
The discriminant function minimizes Bayesian data overlap - log(P(+1) / P(-1))
log b = -1/2 (x - mu\_+)^T Sigma^{-1} (x - mu\_+) - 1/2 (x - mu\_-)^T Sigma^{-1} (x - mu\_-) + 1/2 ln |Sigma| - 1/2 ln |Sigma| + log pi
delta\_k(x) = mu\_k^T Sigma^{-1} x - 1/2 mu\_k^T Sigma^{-1} mu\_k + log pi\_k (predict class k that maximizes delta\_k(x))
Assume class conditional distributions are Gaussian, common covariance
MLE estimate theta\_hat = P(y = +1) mu\_+ - P(y = -1) mu\_- / (P(y = +1) + P(y = -1))
theta\_hat = E[x | y = +1] - E[x | y = -1] / (P(y = +1) + P(y = -1))

Parameter Estimation Methods

We are interested in estimating Bernoulli and Gaussian parameters - distribution
Bernoulli Random Variable (biased coin)
Method of Moments
Choose p such that distribution has same expectation as average of data
Setting n to be large, define x\_1, ..., x\_n in {0, 1} and choose p = 1/n sum\_{i=1}^n x\_i
Maximum Likelihood
Count number of +1 outcomes and write L(p) = p^{sum x\_i} (1-p)^{n - sum x\_i} -> log-likelihood
log L(p) = sum x\_i log p + (n - sum x\_i) log(1-p) -> d/dp L(p) = 0 -> p\_hat = sum x\_i / n
Regularized Maximum Likelihood
Incorporate prior information that p close to some p\_0
log L(p) + log pi(p) = log L(p) + log pi(p) = log L(p) + log pi(p) -> d/dp L(p) = 0
Bayesian estimation
Model p as a random variable and belief about p captured by prior over possible values
Can start by modeling as uniform distribution if no a priori preference
Update belief using Bayes Theorem: P(p | x) = P(x | p) P(p) / P(x)
Bayesian approach gives us distribution and uncertainty (via Bayes rule)
Assumes that parameters randomly chosen with fixed, known distribution
Something is just fit with known distribution
Bayesian estimation is just a computation of conditional probability distribution
Maximum a posteriori (MAP) estimate is just mode of the distribution
Loss: likelihood if uniform prior, penalized with other priors

Gaussian Random Variable

Model p as a random variable and belief about p captured by prior over possible values
Can start by modeling as uniform distribution if no a priori preference
Update belief using Bayes Theorem: P(p | x) = P(x | p) P(p) / P(x)
Bayesian approach gives us distribution and uncertainty (via Bayes rule)
Assumes that parameters randomly chosen with fixed, known distribution
Something is just fit with known distribution
Bayesian estimation is just a computation of conditional probability distribution
Maximum a posteriori (MAP) estimate is just mode of the distribution
Loss: likelihood if uniform prior, penalized with other priors

Multivariate Normal Distribution

Review the Gaussian density function p(x) = 1 / (sqrt(2\*pi))^d exp(-x^T x / (2\*sigma^2))
M, R, sigma^2 in R, normalization: arg of exp: -x^T x / (2\*sigma^2) so (mu, cov, max at x = mu)
Likelihood scale
p(x) = 1 / (sqrt(2\*pi))^d exp(-1/2 (x - mu)^T Sigma^{-1} (x - mu)) (x in R^d)
MAP estimate: maximize log-likelihood: arg max\_x log p(x) = arg max\_x -1/2 (x - mu)^T Sigma^{-1} (x - mu)
Location: MLE estimate: Sigma^{-1} x = Sigma^{-1} mu -> mu\_hat = Sigma x
Consider two independent Gaussian random variables X and Y:
p(x, y) = p(x) p(y) = 1 / (sqrt(2\*pi))^d exp(-x^T x / (2\*sigma\_x^2)) \* 1 / (sqrt(2\*pi))^d exp(-y^T y / (2\*sigma\_y^2))
Covariance Matrices
Symmetric: Sigma = Sigma^T, Non-negative diagonal entries: Sigma\_{ii} >= 0
Positive semidefinite: for all v in R^d, v^T Sigma v >= 0
Diagonal Covariance Matrices
Consider X = N(mu, Sigma) where Sigma is diagonal, so components of x uncorrelated
Look at super-level sets of p(x): E = {x in R^d: (x - mu)^T Sigma^{-1} (x - mu) <= r^2}
(x - mu)^T Sigma^{-1} (x - mu) <= r^2 <= sum\_{i=1}^d (x\_i - mu\_i)^2 / sigma\_i^2 <= r^2
In this case corresponds to axis-aligned ellipsoids in R^d
If all terms but x\_i = mu\_i: (x\_i - mu\_i)^2 / sigma\_i^2 <= r^2 -> |x\_i - mu\_i| <= r sigma\_i
The volume of the ellipsoid E is proportional to r^d sigma\_1 sigma\_2 ... sigma\_d
shows that normalization factor, for fixed r, P(E) is same for all Sigma
Non-Diagonal Covariance and Principalization
Start with a diagonal Sigma and consider transformed vector y = Ax + G(x - mu)
where G is matrix corresponding to rotation about origin [orthogonal]
Mean mu unchanged by E[y = mu + G(x - mu)] = mu + G(mu - mu) = mu
Covariance is transformed Sigma\_y = E[(y - mu)(y - mu)^T] = E[G(x - mu)(x - mu)^T G^T] = G Sigma G^T
Eigenvalues and Eigenvectors
For a square matrix A, lambda is an eigenvalue and x is an eigenvector if Ax = lambda x
These eigenvalues are root of characteristic polynomial det(A - lambda I)
Spectral Theorem
For a symmetric real matrix A in R^n, we can find n orthonormal eigenvectors of A (x\_1, ..., x\_n) and eigenvalues (lambda\_1, ..., lambda\_n) are real
Can then write AU = UA, A = U Lambda U^T where U = [v\_1, v\_2, ..., v\_n], Lambda = Diag([lambda\_1, lambda\_2, ..., lambda\_n])
Covariance matrix Sigma symmetric so Sigma = U Lambda U^T and positive semidefinite
v\_1, v\_2, ..., v\_n: v\_i^T v\_j = delta\_{ij}, v\_i^T U = e\_i^T, U^T v\_i = e\_i
Regular matrix can be written Sigma = U Lambda U^T
x = U [z\_1, z\_2, ..., z\_n]^T, z\_i = U^T x = U^T (mu + G(x - mu)) = U^T mu + U^T G(x - mu)
U^T x = U^T mu + U^T G(x - mu)
U^T x - U^T mu = U^T G(x - mu)
where the covariance is the diagonal matrix Lambda
Level sets E = {x in R^d: (x - mu)^T Sigma^{-1} (x - mu) <= r^2} corresponds to orthogonal aligned ellipsoids in R^d
(x - mu)^T Sigma^{-1} (x - mu) = (U^T x - U^T mu)^T Lambda^{-1} (U^T x - U^T mu) = sum\_{i=1}^d (z\_i - mu\_i)^2 / lambda\_i
so it's (x - mu)^T Sigma^{-1} (x - mu) <= r^2 <= sum\_{i=1}^d (z\_i - mu\_i)^2 / lambda\_i <= r^2
Value of ellipsoid E proportional to product of lambda\_i = lambda\_1 lambda\_2 ... lambda\_n = |Lambda| = |U^T Sigma U| = |Sigma|
Properties of Multivariate Gaussians
Probability: P(x in E) = E[exp(-x^T Sigma^{-1} x / (2\*sigma^2))] = E[exp(-x^T Sigma^{-1} x / (2\*sigma^2))]
Covariance: Cov(x) = E[(x - mu)(x - mu)^T] = E[(U^T x - U^T mu)(U^T x - U^T mu)^T] = U^T E[(x - mu)(x - mu)^T] U = U^T Sigma U
Sum of independent Gaussians: x = x\_1 + x\_2, y = mu\_1 + mu\_2, Sigma = Sigma\_1 + Sigma\_2
Mean: E[x] = mu\_1 + mu\_2, Cov(x) = E[(x - mu)(x - mu)^T] = E[(x\_1 - mu\_1 + x\_2 - mu\_2)(x\_1 - mu\_1 + x\_2 - mu\_2)^T] = E[x\_1 x\_1^T] + E[x\_2 x\_2^T] + E[x\_1 x\_2^T] + E[x\_2 x\_1^T] = Sigma\_1 + Sigma\_2
Marginals
Write d-dimensional Gaussian as x = [x\_1, x\_2]^T, mu = [mu\_1, mu\_2]^T, Sigma = [Sigma\_11, Sigma\_12; Sigma\_21, Sigma\_22] where
y in R^m and z in R^n, then y = N(mu\_y, Sigma\_y), z = N(mu\_z, Sigma\_z)
Mean and covariance of y evident from definition, show (basis) by writing determinant as integral (in part 2), complete squares -> N(mu\_y, Sigma\_y)
Affine Transformations
Given d-dimensional Gaussian X ~ N(mu, Sigma), matrix A in R^n x d, b in R^n -> y = Ax + b ~ N(mu\_y, Sigma\_y)
mu\_y = A mu + b, Sigma\_y = E[(y - mu\_y)(y - mu\_y)^T] = E[(Ax + b - A mu - b)(Ax + b - A mu - b)^T] = A E[(x - mu)(x - mu)^T] A^T = A Sigma A^T
Given d-dimensional normal X ~ N(mu, Sigma) Y composed in R^n, Z ~ N(mu\_z, Sigma\_z)
Y ~ N(mu\_y, Sigma\_y) Y composed in R^n, Z ~ N(mu\_z, Sigma\_z)

Regression

Linear Regression
Recall quadratic loss function L(y, y') = (y - y')^2, risk R(theta) = E[(f(x) - y)^2]
Goal: minimize L(theta) = E[(f(x) - y)^2] = E[(theta^T x - y)^2]
One-variable: L(theta) = E[(theta\_0 + theta\_1 x - y)^2] = E[(theta\_0 + theta\_1 x - y)^2]
Given X in R^d, Y in R, linear regression model: f(x) = theta\_0 + theta\_1 x + theta\_2 x^2 + ... + theta\_d x^d
Consider class of linear prediction rules and minimize empirical risk
Model the process of generating y\_i as linear function of x\_i plus Gaussian noise. Compute MLE for linear coefficients.
Both equations should arise at normal equations: the choice of beta corresponds to projection onto linear subspace

Empirical Risk Minimization
Risk R(theta) = E[(f(x) - y)^2] is the expected squared error while
Empirical risk R\_hat(theta) = sum\_{i=1}^n (f(x\_i) - y\_i)^2 is sample average of squared error
Choose linear prediction rule f in F that minimizes empirical risk by beta = arg min\_{beta in F} sum\_{i=1}^n (f(x\_i) - y\_i)^2 = arg min\_{beta in F} sum\_{i=1}^n (f(x\_i) - y\_i)^2
We can consider off-set term beta\_0 and minimize component to x
Solve beta\_0 = arg min\_{beta\_0} sum\_{i=1}^n (beta\_0 + theta^T x\_i - y\_i)^2 = arg min\_{beta\_0} sum\_{i=1}^n (beta\_0 - y\_i + theta^T x\_i)^2
beta\_0 = (sum\_{i=1}^n (y\_i - theta^T x\_i)) / n = sum\_{i=1}^n (y\_i - theta^T x\_i) / n
beta\_0 = (sum\_{i=1}^n y\_i - theta^T sum\_{i=1}^n x\_i) / n = sum\_{i=1}^n y\_i / n - theta^T sum\_{i=1}^n x\_i / n
beta\_0 = y\_bar - theta^T x\_bar
beta\_0 = y\_bar - theta^T x\_bar
View as finding linear combination of columns of X: X = [x\_1, x\_2, ..., x\_n]
to minimize Euclidean distance to y in R^n by y - X beta = 0
Optimal approximation q in space spanned by columns of X has error ||y - q|| orthogonal to column space
(y - q)^T X = 0 <=> X^T (y - X beta) = 0 <=> X^T y = X^T X beta
Risk of empirical risk minimizer q close to minimal if
X comes from compact set S in R^d, y has bits not too heavy (e.g. Gaussian)
||beta|| is not too large <=> ||y - q|| <= ||y|| + ||q||

Linear Model with Additive Gaussian Noise
Model conditional distribution of y | x as P(y | x) = N(y | theta^T x, sigma^2) or y = theta^T x + epsilon
Using MLE, L(theta) = sum\_{i=1}^n (y\_i - theta^T x\_i)^2 / (2\*sigma^2) -> d/d theta L(theta) = 0 -> X^T (y - X theta) = 0
we get the least squares solution
Show p unbiased E[y\_i | x\_i] = X beta, Cov(y) = sigma^2 I
E[y] = E[X theta] = X theta, Cov(y) = sigma^2 I
Cov(y) = E[(y - E[y])(y - E[y])^T] = E[(X theta - X theta)(X theta - X theta)^T] = 0
If data generated by linear model with Gaussian noise, we can compute dist. of parameter estimates (beta\_hat - mu) | (beta\_hat - mu) ~ N(0, sigma^2 (X^T X)^{-1})
approximate confidence sets for param beta\_hat = {beta: ||beta\_hat - beta|| <= c} which is normal and design error for non-zero values of param

Regularization in Regression
Bias-variance tradeoff: E[R(theta\_hat)] = R(theta) + E[(theta\_hat - theta)^T (theta\_hat - theta)] = E[(theta\_hat - theta)^T (theta\_hat - theta)]
Use random training points so R(theta\_hat) is varying
Increase in bias decreases variance and overall complexity
Linear considers number of points in S and size of coefficients
Subset Selection
Ideas consider 2^d subsets of variables, fit linear model, and decide which is best
Ideas which model complexity to use with cross validation
We can use greedy computational shortcuts to get approximately best solution
Forward stepwise selection - iterate through all fittings using 2 param, include in next iteration (only if better)
Backward stepwise selection - start with all and remove least contributing

Shrinkage Methods

Shrinkage Methods
Encourage linear predictor coefficients to be small using penalty
Ridge Regression
Penalty using penalty term beta\_hat = arg min (sum\_{i=1}^n (y\_i - x\_i^T beta)^2 + lambda sum\_{j=1}^d beta\_j^2) or
forbid using constraint beta\_hat = arg min (sum\_{i=1}^n (y\_i - x\_i^T beta)^2 s.t. ||beta||\_2 <= B
Do not penalize beta\_0, beta\_0 = (sum\_{i=1}^n y\_i) / n + lambda sum\_{j=1}^d beta\_j^2 / (2\*n)
Even if X singular, penalty term prevents an increase ||X^T X + lambda I||^{-1}
Solution also depends on scaling of covariances (become in standardize)
Lasso
Use L1 norm beta\_hat = arg min (sum\_{i=1}^n (y\_i - x\_i^T beta)^2 + lambda sum\_{j=1}^d |beta\_j|), or with
constraint ||beta||\_1 <= B
Ridge regression leads to reduced coefficients, Lasso pushes to zero

Bayesian View of Linear Regression

Bayesian View of Linear Regression
Given linear model P(y | x) = N(y | theta^T x, sigma^2)
Model beta as RV (beta = N(mu, Sigma)) then compute posterior distribution
P(theta | x, y) = P(y | theta, x) P(theta) / P(y) = exp(-y^T (X theta - mu) / (2\*sigma^2) - 1/2 (theta - mu)^T Sigma^{-1} (theta - mu)) / P(y)
MAP estimates beta as mode of distribution
Ridge regression is MAP 7 Bayesian prior, Lasso is MAP Laplace prior

Laplace Regression

Laplace Regression
P(y | x) = 1 / (sqrt(2\*pi)) exp(-y^2 / (2\*sigma^2))
Example beta by MLE
log likelihood: L(theta) = log P(y | x, theta) = log P(y | theta, x) P(theta) = log P(y | theta, x) P(theta)
MLE yields beta\_hat = (sum\_{i=1}^n y\_i x\_i) / (sum\_{i=1}^n x\_i^2 + 1)
First derivative: d/d beta L(theta) = sum\_{i=1}^n (y\_i x\_i - beta x\_i^2) / (sum\_{i=1}^n x\_i^2 + 1) = 0 -> beta\_hat = (sum\_{i=1}^n y\_i x\_i) / (sum\_{i=1}^n x\_i^2 + 1)
second derivative: d^2/d beta^2 L(theta) = -sum\_{i=1}^n x\_i^2 / (sum\_{i=1}^n x\_i^2 + 1)^2 < 0
MLE yields beta\_hat = (sum\_{i=1}^n y\_i x\_i) / (sum\_{i=1}^n x\_i^2 + 1)

Gradient Ascent

Gradient Ascent
beta^{(t+1)} = beta^{(t)} + eta \* grad L(beta^{(t)}) = beta^{(t)} + eta \* sum\_{i=1}^n (y\_i - x\_i^T beta^{(t)}) x\_i
Each gradient calculation takes O(n\*p) time
Stochastic Gradient Descent/Ascent
Instead compute approximation: grad L(theta) = (y - x^T theta) x
which takes O(p) time and update using the equation
beta^{(t+1)} = beta^{(t)} + eta \* (y - x^T theta^{(t)}) x

Minimizing RLS in Linear Regression can also be done by gradient of statistical gradient descent with convex regularization

Newton's Method

Newton's Method
Use 2nd-order Taylor approximation f(x) = f(x\_0) + f'(x\_0)(x - x\_0) + 1/2 f''(x\_0)(x - x\_0)^2
Solve for f(x) = 0: x\_1 = x\_0 - f'(x\_0) / f''(x\_0)
Solve f(x) = 0 using x\_{i+1} = x\_i - f'(x\_i) / f''(x\_i)
f(x) = sum\_{i=1}^n y\_i log pi(x\_i) + (1 - y\_i) log(1 - pi(x\_i))
f'(x) = sum\_{i=1}^n (y\_i - pi(x\_i)) x\_i = X^T (y - mu)
f''(x) = -sum\_{i=1}^n pi(x\_i) (1 - pi(x\_i)) x\_i x\_i^T = -X^T diag(pi(x\_i) (1 - pi(x\_i))) X
Each iteration is less squares with x\_i x\_i^T weighted by pi(x\_i) (1 - pi(x\_i)) (heavily weighted least squares)
New MLEs minimizing log loss for beta^{(t)} = beta^{(t-1)} + eta \* grad L(beta^{(t-1)})
Define loss as L\_n(beta^{(t)}) = -log L(beta^{(t)}) = -log P(y | x, beta^{(t)})
Write log loss as log L(beta^{(t)}) = -log P(y | x, beta^{(t)}) = -log P(y | x, beta^{(t)})
One MLE corresponds to minimizing sample average 1/n sum\_{i=1}^n log(1 + exp(-y\_i theta^T x\_i))

Convex Optimization and SVMs

Recall hard margin SVM:  $\min_x \|G\|^2$  s.t.  $y_i G \cdot x_i \geq 1$

soft margin SVM:  $\min_x \|G\|^2 + C \sum_{i=1}^n (1 - y_i G \cdot x_i)$

ASGP: hard margin SVM:  $\min_x \|G\|^2$  s.t.  $y_i G \cdot x_i \geq 1$

soft margin SVM:  $\min_x \|G\|^2 + C \sum_{i=1}^n \xi_i$  s.t.  $\xi_i \geq 0, \xi_i \geq 1 - y_i G \cdot x_i$

Convex sets are those in which it is possible to draw line all in set

Consider convex optimization problem  $p^* = \min_x f_0(x)$  s.t.  $f_i(x) \leq 0$   $\forall i \in \{1, \dots, m\}$

Writing constraints as penalties  $p^* = \min_x f_0(x) + \begin{cases} 0 & \text{if all } f_i(x) \leq 0 \\ \infty & \text{otherwise} \end{cases}$

Replace constraint penalty with something smaller:

Introduce Lagrange multipliers (dual variables)  $\lambda_1, \dots, \lambda_m \geq 0$  and define

Lagrangian  $L: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$  as  $L(x, \lambda) = f_0(x) + \sum_{i=1}^m \lambda_i f_i(x)$

$\lambda_i$  cost of violating constraint  $f_i(x) \leq 0$

$L$  defines saddle point game: one player (MIN) chooses  $x$  to minimize

$L$ , the other player (MAX) chooses  $\lambda$  to maximize  $L$ . If MIN violates a

constraint,  $f_i(x) > 0$ , MAX can drive  $L$  to infinity.

Primal problem  $p^* = \min_x \max_{\lambda \geq 0} L(x, \lambda)$  (infeasible  $x \rightarrow \infty$ , feasible  $x \rightarrow \lambda_i f_i(x) = 0$ )

dual problem  $d^* = \max_{\lambda \geq 0} g(\lambda) = \max_{\lambda \geq 0} \min_x L(x, \lambda)$   $g(\lambda) := \min_x L(x, \lambda)$

This is a zero sum game where better to play second:

$p^* = \min_x \max_{\lambda \geq 0} L(x, \lambda) \geq \max_{\lambda \geq 0} \min_x L(x, \lambda) = d^*$  (weak duality)

If there is saddle point  $(x^*, \lambda^*)$  so that for all  $x$  and  $\lambda \geq 0$ ,  $L(x^*, \lambda) \leq L(x^*, \lambda^*) \leq L(x, \lambda^*)$

then  $p^* = \min_x \max_{\lambda \geq 0} L(x, \lambda) = \max_{\lambda \geq 0} \min_x L(x, \lambda) = d^*$  (strong duality)

Complementary Slackness

If  $p^* = d^*$  and we have primal solution  $x^*$  and dual solution  $\lambda^*$

then for  $i$ th constraint  $(f_i(x^*) \leq 0)$ ,  $\lambda_i^* f_i(x^*) = 0$ .

$f_i(x^*) < 0 \Rightarrow \lambda_i = 0$ ,  $\lambda_i > 0 \Rightarrow f_i(x^*) = 0$

As a result, every term in  $\sum_{i=1}^m \lambda_i^* f_i(x^*) \geq 0 \approx 0$

Karush-Kuhn-Tucker Optimality Conditions

Suppose  $f_0, f_i$  are convex and differentiable. Then  $x$  and  $\lambda$  optimal if and only if:

① Primal feasibility:  $f_i(x) \leq 0$  ② Dual feasibility:  $\lambda_i \geq 0$

③ Complementary slackness:  $\lambda_i f_i(x) = 0$  ④ Stationarity:  $\nabla f_0(x) + \sum_{i=1}^m \lambda_i \nabla f_i(x) = 0$

SVMs

Hard Margin SVMs  $\theta(w) = \min_x L(\theta; x)$

write  $\min_x \|G\|^2$  s.t.  $y_i G \cdot x_i \geq 1$  as  $L(\theta; x) = \frac{1}{2} \|G\|^2 + \sum_{i=1}^n \alpha_i (1 - y_i G \cdot x_i)$

$\nabla_x L(\theta; x) = G + (-1) \sum_{i=1}^n \alpha_i y_i x_i \rightarrow G^* = \sum_{i=1}^n \alpha_i y_i x_i$

$g(\theta) = \frac{1}{2} \theta^T \theta + \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \alpha_i y_i G \cdot x_i = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$

If  $\exists$  feasible  $\theta \Rightarrow$  strong duality  $\rightarrow \max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$

$\alpha_i > 0 \Rightarrow y_i G^* \cdot x_i = 1$  and  $y_i G^* \cdot x_i \leq 1 \Rightarrow \alpha_i = 0$

Express solution as kernel  $\langle x, x \rangle = \text{sign}(\langle \theta, \theta(x) \rangle) = \text{sign}(\sum_{i=1}^n \alpha_i \langle x, x_i \rangle \langle \theta, x_i \rangle)$

where  $\alpha$  solves dual problem  $\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$

Soft Margin SVMs  $\min_x \frac{1}{2} \|G\|^2 + C \sum_{i=1}^n (1 - y_i G \cdot x_i)$  as  $\min_x \frac{1}{2} \|G\|^2 + \frac{C}{n} \sum_{i=1}^n \xi_i$   $\xi_i \geq 0, \xi_i \geq 1 - y_i G \cdot x_i$

$L(\theta, \xi, \alpha) = \frac{1}{2} \theta^T \theta + \frac{C}{n} \sum_{i=1}^n \xi_i + \sum_{i=1}^n \alpha_i (1 - y_i G \cdot x_i - \xi_i) \rightarrow \frac{C}{n} \sum_{i=1}^n \lambda_i \xi_i$

$\nabla_x L(\theta; x) = G + (-1) \sum_{i=1}^n \alpha_i y_i x_i \rightarrow G^* = \sum_{i=1}^n \alpha_i y_i x_i$

$\nabla_{\xi} L(\theta; \xi, \alpha) = \frac{C}{n} - \alpha_i \rightarrow \alpha_i + \xi_i = \frac{C}{n}$  so  $g(\theta; \lambda) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$

Dual problem given by  $\max_{\alpha_i \geq 0} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j x_i^T x_j$

$\alpha_i \geq 0, x_i \geq 0, \alpha_i + \lambda_i = C/n$

Consequences of complementary slackness:  $\alpha_i^* (1 - y_i x_i^T G^* - \xi_i^*) = 0$

$\alpha_i^* > 0 \Rightarrow y_i x_i^T G^* = 1 - \xi_i^* \leq 1$  'support vector'  $x_i^* \xi_i^* = 0$

ie in the wrong side of half space

If  $y_i x_i^T G^* < 1$ ,  $\xi_i^* > 0$  so  $x_i^* = 0$  and  $\alpha_i^* = C/n$  so

support vectors in open halfspace have  $\alpha_i^* = C/n$