

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  $\mu_c, \mu_x$

Perceptron

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

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

$R(w) = \sum_{i \in V} -y_i x_i \cdot w, V = \{i: y_i x_i \cdot w < 0\}$   
Goal: minimize risk  $= \sum_{i \in V} LC(x_i, y_i)$

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  $w, \alpha$  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

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  $\xi_i$  iff it violates the margin

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

large maximize margin  $\Rightarrow$  keep margin large  $\Rightarrow$  any point in slab/margin considered support vector

overfitting  $\Rightarrow$  overfitting  $\Rightarrow$  very sensitive  $\Rightarrow$  more sensitive

outliers less sensitive  $\Rightarrow$  more robust

boundary more "flat"  $\Rightarrow$  more robust

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

- 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

Optimization Problem

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

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.

$R(c) = E[L(c(X), y)]$

Bayes Decision Rule: - assumes no loss

for correct prediction

$R(c) = \int L(c(x), y) P(y) P(x) dx$

$R(c) = \int L(c(x), y) P(y) P(x) dx$

Bayes decision rule

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

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

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

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

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

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

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

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

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

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  $P(X=x|Y=c)P(Y=c)$

$Q_c(x) = R_c(\frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma_c^2}))$

$= \frac{1}{\sqrt{2\pi}\sigma_c} \exp(-\frac{(x-\mu_c)^2}{2\sigma_c^2})$

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  $\sigma^2$

$Q_c(x) - Q_p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

$= \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x-\mu_c)^2}{2\sigma^2})$

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 and radii are the reciprocals of

eigenvalues, in matrix  $A$  maps

spheres to ellipsoids

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

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

$\Sigma$ : covar matrix

$\Sigma^{-1}$ : precision matrix

$P(x) = n(C(x), \mu) \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))$

$g(x)$  is signed distance from



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:  
 $\mu_C^T x - \frac{1}{2} \mu_C^T \Sigma^{-1} \mu_C + \ln \Pi_C$

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  $\mu^T$  from each row of  $X$   
 $\rightarrow \bar{X}$

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

decomposing  $X'$ : applying  $Z = \bar{X} V$ ,  $\text{Var}(C) = V \Delta V^T$   
 $\rightarrow$  transform sample points to eigens coordinate system  
 $\rightarrow \text{Var}(Z) = \Delta$

sphering  $X'$ : applying transform  $W = \bar{X} \text{Var}(C)^{-1/2}$

whitening  $X'$ : centering + sphering,  $X \rightarrow W$

$W$  has covariance matrix  $I$

whitening: put 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:  $XTXw = XT y$   $= \sum_{i=1}^n (y_i \cdot w + d - y_i)^2$

Let  $X' = (XTX)^{-1} XT$

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

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

$X \in \mathbb{R}^{n \times d}$  at most  $d+1$  dim

$y \in n$  d's space

Minimizing  $\|y - y'\|^2$  finds  $\hat{y}$  nearest  $y$  in subspace

= orthogonal projection.

Advantages:

- easy to compute; just solve a linear system

- unique, stable soln.

Disadvantages:

- sensitive to outliers

- fails if  $XTX$  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 trustful 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 - iteratively

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.

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

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

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

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

$\Omega$  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 from

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 negative rate

True negative 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 dist

$y$ -values are sum of unknown non-zero

fn + random noise

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

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

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

Max likelihood explains when 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

adding a feature increases

variance

noise in test set affects only

$\text{Var}(C)$

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

$\hookrightarrow (XTX + \lambda I)^{-1} w = XT y$

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. (Odds)

Heuristic 2: Backward stepwise selection

- start w all d features, remove feature

whose removal gives best reduction in validation

error (Odds)

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.

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

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

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

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

Posterior:  $PCY = p(y|x)$

$\Theta_{MLE} = \arg \max_{\Theta} PCX(\Theta)$

$= \arg \max_{\Theta} \prod_{i=1}^n PCX_i(\Theta)$

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

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

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

Covariance matrix must be PSD

Gaussian prior = L2 reg.

COV =  $\frac{\text{cov}}{\sigma_x \sigma_y}$  Laplace prior = L1 reg

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  $\Sigma$

- 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  $\lambda I$  to covariance matrix

causes isocountours to be more spherical

$SC_0 = 1/2$  is decision boundary

No covariance terms: isocountour

is axis aligned

- space between contours is

$\sqrt{\lambda_i}$  of  $\Sigma$

- higher SD  $\Rightarrow$  larger gaps

between significant

isocountours.

- eigenvectors form

orthonormal basis

- tell of something is

convex by whether

its Hessian is positive

definite.

$A^T A = |A|^2$