

stretching/reducing parameters with proportionality box-cox transformation

leverage points have proportionality box-cox transformation

outliers: $\hat{y}_i = \text{original value} + \text{Box-Cox transformation}$

heteroscedasticity unequal variance: $\hat{y}_i = \text{original value} + \text{Box-Cox transformation}$

is a logarithmic transformation. Stretch out smaller range

to enlarge variability by shrinks larger range to reduce its variability

$(\hat{y}_i - \bar{y}) = (y_i - \bar{y})/\lambda$, $y_i - \bar{y}$ can be close to normal

stretches/reduces effects from time series data

transform data before using time series data in regression model

1.3 Principal Component Analysis (PCA): PCA transforms data

to remove correlation within the data. @ make coordinates

less important. Concentrate on \hat{x}_i principal component

reduce effects of randomness & likely to have higher signal to noise ratio

9.4 PCA: Eigenvectors & Eigenvalues

\hat{x}_i is eigenvalue of A s.t. $\det(A - \lambda I) = 0$

Given λ , solve $\hat{x}_i = \lambda \hat{x}_i$ to find corresponding eigenvector

Given each λ value, \hat{x}_i is the jth factor of data point

when we solve $\hat{x}_i^T \hat{x}_j = \lambda_j \hat{x}_i^T \hat{x}_i$, $\hat{x}_i \in \mathbb{R}^n$, \hat{x}_i no. of factors

$X^T X$ (covariance matrix) s.t. $\det(X^T X - \lambda I) = 0$

given each λ_j value, \hat{x}_i is the jth factor of data point

scale s.t. $\hat{x}_i^T \hat{x}_i = \mu_j = 0$. Then find all eigenvalues of

$X^T X$ (covariance matrix) s.t. $\det(X^T X - \lambda I) = 0$

$\lambda = [\lambda_1 \lambda_2 \dots]$ when \hat{x}_i is the eigenvector of $X^T X$

Principal Components note each \hat{x}_i has order $n \times 1$

while $n = m$, no. of factor columns in X

$X^T \hat{x}_i = 1^{\text{st}} \text{ component } \hat{x}_i$ order is $m \times 1$

new factor value for i-th data point: $t_{ik} = \sum_{j=1}^m x_{ij} y_{jk}$

PCA finds new L factor \hat{x}_i when $L < n$

regression finds coefficient b_0, b_1, \dots, b_L

$y_i = b_0 + \sum_{k=1}^L b_k t_{ik} = b_0 + \sum_{k=1}^L b_k [\sum_{j=1}^m x_{ij} y_{jk}]$

$V = \text{matrix of eigenvectors (sorted by eigenvalue)}$

PCA finds new L factor \hat{x}_i when $L < n$

regression finds coefficient b_0, b_1, \dots, b_L

$y_i = b_0 + \sum_{k=1}^L b_k t_{ik} = b_0 + \sum_{k=1}^L b_k [\sum_{j=1}^m x_{ij} y_{jk}]$

$V = \text{matrix of eigenvectors (sorted by eigenvalue)}$

PCA finds new L factor \hat{x}_i when $L < n$

regression finds coefficient b_0, b_1, \dots, b_L

the subsets of data points that follow all of the branches leading to the leaf node

potential branch (low improvement benefit)

too few data points (leaf node contains > 5% of original data)

Random Forests ① Building n points from original data points

choose best feature will either fit to branch or

partition tree (true random forest) or common partition

Classification Tree \rightarrow use random forest (common partition)

forests are black boxier (models tree are less explainable might give us better predictions)

logistic regression: standard linear regression $y_i = a_0 + a_1 x_{ij} + a_2 x_{ij}^2$

log regression: P : probabilities of event you want to observe

$\log \frac{P}{1-P} = a_0 + a_1 x_{ij} + a_2 x_{ij}^2$ if $a_2 > 0 \Rightarrow P = 0$

$P = \frac{1+e^{-(a_0+a_1 x_{ij}+a_2 x_{ij}^2)}}{1+e^{-(a_0+a_1 x_{ij}+a_2 x_{ij}^2)}}$

$\Delta x = +\infty \Rightarrow P = 1$

$\Delta x = -\infty \Rightarrow P = 0$

Explanatory variable x_{ij} has order $n \times 1$

are constant coefficients

optimization model for regression

minimize $\sum_{i=1}^n (y_i - (a_0 + \sum_{j=1}^m a_j x_{ij}))^2$

a_j are variables

Three things: Answer "yes" to probability P is at least some number

over was "no". E.g. if $P > 0.5$, then "yes"

AUC: probability that the model estimates a random "yes" point higher than a random "no" point

If AUC = 0.5 means just guessing

Model classification: True if refer to whether the model is correct or not. Or P or N refer to whether the model is correct or not.

say the point is in the category

Logistic Regression Variable a_0, a_1, \dots, a_m

Maximize $\prod_{i=1}^n p(x_i) \prod_{i=1}^n (1-p(x_i))$

where $p(x_i) = \frac{1}{1+e^{-(a_0+\sum_{j=1}^m a_j x_{ij})}}$

1st derivative

Time Series Models

ARIMA Variable μ, ρ_1, ρ_2

constraint: more

Exponential Smoothing minimize $\sum_{t=1}^n (x_t - \hat{x}_t)^2$ where

Variables μ, ρ_1, ρ_2

Constraints: $0 \leq \rho_1 \leq 1$

$0 \leq \rho_2 \leq 1$

$0 \leq \rho \leq 1$

Convex Quadratic program

$f(x)$ is a convex quadratic function

minimize $f(x)$ or maximize $-f(x)$

constraint set X is defined by linear

equations or inequalities

Convex (if minimization)

Linear Program / Convex Optimization Program

Obj function $f(x)$ is concave (if maximization)

Convex (if minimization)

constraint set X is a convex set

equation to solve but not as quickly as

Easy to solve even for large matrices

Convex Quadratic program

$f(x)$ is a convex quadratic function

minimize $f(x)$ or maximize $-f(x)$

constraint set X is defined by linear

equations or inequalities

Convex (if minimization)

Linear Program / Convex Optimization Program

Obj function $f(x)$ is concave (if maximization)

Convex (if minimization)

constraint set X is a convex set

equation to solve, but slow can take longer

Convex Non-Convex Program

Optimization Problem is not convex

Upward to find optimal solutions

General Non-Convex Program

Optimization Algorithm: Find a good direction

to move from the current solution, and

determine how far to go in that direction

Survival Model / Cox Proportional Hazard Model

Given (for each data point) predictor variables x_1, x_2, \dots, x_n

Event time y

$h(t) = h_0(t) e^{(p_1 x_1 + \dots + p_n x_n)}$

$p_j = \text{coefficient for } x_j$

$p_j = \text{prob. each predictor variable multiplies factor } j$

Baseline (if all predictor variables are zero) probability of event happening at time t

Prob of event happening at time t given specific values of predictors (k coefficients)

Variables are zero, probability of event happening at time t

Prob of event happening at time t given specific values of predictors (k coefficients)

$p_j = \text{prob. each predictor variable multiplies factor } j$

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Prob of event happening at time t given specific values of predictors (k coefficients)

Ridge (global)

$$y = w_0 x_0 + w_1 x_1 + \dots + w_n x_n + b$$

introduce small amount of bias to reduce variability.

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M (y_i - \sum_{j=0}^P w_j x_{ij})^2 + \lambda \sum_{j=0}^P w_j^2$$

$\lambda \uparrow \Rightarrow$ slope gets close to zero (i.e. horizontal)

When sample sizes are relatively small, then ridge regression can improve prediction made from new data (i.e. reduce variance) by making the predictor less sensitive to training data.

Coefficients shrink toward zero to reduce variance in estimate i.e. regularization.

Lasso (global)

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M (y_i - \sum_{j=0}^P w_j x_{ij})^2 + \lambda \sum_{j=0}^P |w_j|$$

Lasso can exclude variables.

Data needs to be scaled for Lasso as constraints can go to zero.

Elastic Net (global) & quality of model

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M (y_i - \sum_{j=0}^P w_j x_{ij})^2 + \lambda \sum_{j=0}^P |w_j| + (1-\lambda) \sum_{j=0}^P w_j^2$$

Scale data If two features are highly correlated, Lasso will prob just one & leave the other out. Ridge will include both, & because both coefficients shrink, they get closer to each other.

Design of Experiments.

Multi Armed Bandits \rightarrow As we get more sure of the best answer, we're more likely to use it.
Exploration - Exploitation Trade off.

Missing Data

Categorical \rightarrow add categories

Numerical \rightarrow Create new column w/ binary yes/no (i.e. missing)

Create interaction terms w/ binary variable

Split 2 models \rightarrow one w/ no missing & one w/ missing.

Optimization Binary Variable

Variable $\rightarrow y_i = 1$ if invest in stock 0 if not

Constraint $\rightarrow y_{AMZN} + y_{GOOG} + y_{AAPL} \geq 1$

Object \rightarrow Another constraint w/ binary

New Soln = old soln + βt

Non Parametric Test

McNemar (yes/no)

Wilcoxon Signed Rank Test (numerical)

Mann Whitney (categorical)

Paired vs Unpaired

Escape Velocity \rightarrow parametric (mean)

Rank of Value \rightarrow Non parametric (median)

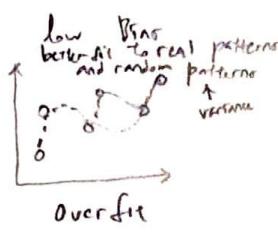
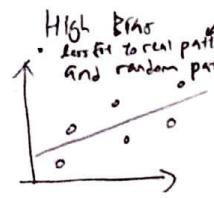
Count of Binary Outcome \rightarrow Bernoulli (which is better)

Comparing number of strata \rightarrow Fisher's test

Variable Selection

$$\sum_{i=1}^M (y_i - \hat{y}_i)^2 = \sum_{i=1}^M (y_i - \sum_{j=0}^P w_j x_{ij})^2$$

↑↑↑
Observed Estimate Cost Function



Underfit
(allow misclassification)
Arrange job fitting training set

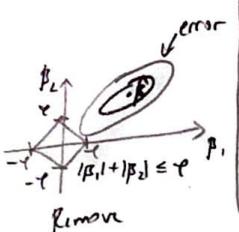
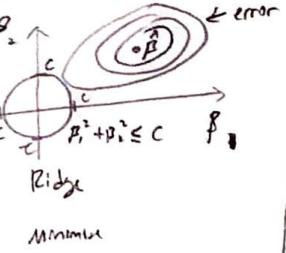
Ok job fitting testing set

Low Variability/training
SSE similar for varying distance

Good job fitting training set

Bad job fitting testing set

High Variability/
SSE very different for varying distance



to find sweet spot between bias & variance due regularization, Lasso & LASSO

Forward selection, Backward Elimination

& Stepwise Regression (greedy)

At every step, the stepwise regression fits a different model

A/B Testing.

Ad	Clicker	Users	%
A	46	1003	5%
B	97	993	10%

Factorial Design \rightarrow ANOVA to determine factor importance

Fractional Factorial Design
 \rightarrow Test each choice the same # of times
 \rightarrow Test each pair of choices the same # of times

Optimization

Variables \rightarrow decisions to be made

Constraints \rightarrow restrictions on variables

Objective Functions \rightarrow solution quality measure

Solution \rightarrow values for each variable

Diet Problem

Variables $\rightarrow x_i = \text{amt of food } i \text{ in daily diet}$

Constraints

$\sum_i a_{ij} x_i \geq m_j$ for each nutrient j

$\sum_i a_{ij} x_i \leq M_j$ for each nutrient j

$x_i \geq 0$ for each food i

Objective function

minimize $\sum_i c_i x_i$

where $a_{ij} = \text{amt of nutrient } j \text{ per unit of food } i$

$m_j = \text{min daily intake of nutrient } j$

$M_j = \text{max daily intake of nutrient } j$

$c_i = \text{per-unit cost of food } i$

Importance

mean/median (numerical)
mode (categorical)

Regression Imputation

+ perturbation \rightarrow imputed value less accurate but will give better estimate of the overall spread of the values

\rightarrow no more than 5% per factor

\rightarrow data used twice in regression will lead to overfitting

Bayesian Model
 $P(A|B) = \frac{P(B|A)P(A)}{P(B)}$

Graph to Communities \rightarrow Latent Algo

Neural Network
Input Hidden Output
recognize images, speech, writing, language

Competition \rightarrow Game Theory

(1) perfect info (know everyone's strategy)
(2) imperfect info (no info about others)

(3) zero sum (win/lose)

(4) non-zero sum (total benefit higher or lower)

\rightarrow player & opponents' decisions must be accounted for.

Probability Models

Bernoulli
iid
 $P(X=1) = p$
 $P(X=0) = 1-p$

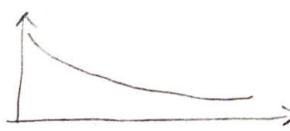
Bernoulli \rightarrow must be independent trials
Probability of getting x successes
out of n independent identically
distributed Bernoulli (p) trials
 $n \rightarrow \infty$ then Binomial Converges to Normal
pmf

$$P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

Geometric

Probability of having x Bernoulli (p)
trials/failures before first success

$$P(X=x) = (1-p)^x p$$



"number of tries between failures"

Poisson \rightarrow good model for random arrivals
 $f_x(x) = \frac{\lambda^x e^{-\lambda}}{x!}$

λ = average number of occurrences in a time period
occurrences are IID \star memoryless

mean interarrival rate is $1/\lambda$ so
 λ is the mean arrival rate per unit time

Exponential

$$f_x(x) = \lambda e^{-\lambda x}$$

If occur in Poisson (λ)
then time between successive occurrences are exponential (λ)



Weibull

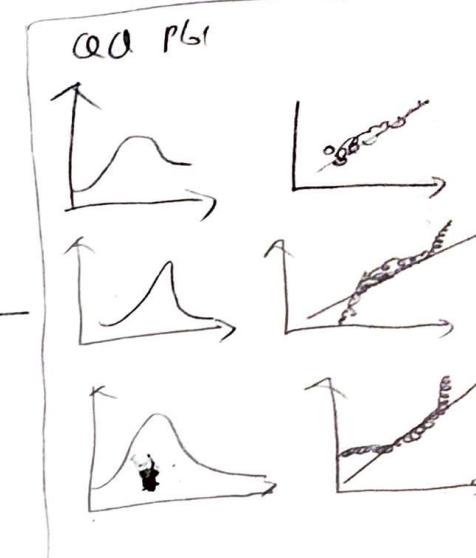
$$f_{wc}(x) = \frac{k}{\lambda} \left(\frac{x}{\lambda}\right)^{k-1} e^{-(x/\lambda)^k}$$

"time between features"

$k < 1 \rightarrow$ failure rate decreases with time
"worst things fail first"

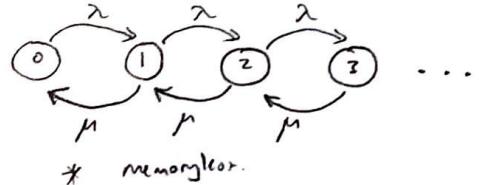
$k > 1 \rightarrow$ failure rate increases with time
"things that wear out"

$k=1 \rightarrow$ failure rate constant w/ time } exponential



Queuing

Arrival Rate = λ
Service Rate = μ
* memoryless.



Transition Equations (≥ 1 calls in the queue)

- $P(\text{next event is arrival}) = \frac{\lambda}{\lambda + \mu}$
- $P(\text{next event is finished}) = \frac{\mu}{\lambda + \mu}$

Expected fraction of time employee is busy = $\frac{\lambda}{\mu}$ } Utilization factor / traffic intensity

Expected waiting time before talking = $\frac{\lambda}{\mu(\mu - \lambda)}$

Expected number of calls waiting in queue = $\frac{\lambda^2}{\mu(\mu - \lambda)}$

Avg no. of users in system (waiting + being served) = $\frac{\lambda}{\mu - \lambda}$

Avg time a unit spends in the system = $\frac{1}{\mu - \lambda}$

$P(\text{more than } k \text{ units in system}) = \left(\frac{\lambda}{\mu}\right)^{k+1}$

Simulation

Deterministic: same inputs give same outputs (no randomness)

Stochastic: system has randomness

\hookrightarrow one random outcome might not be representative of system performance in the range of different situations that could arise

Markov Chain memoryless ie don't care previous state
only care present state
 $P_{ij} = \text{transition prob from state } i \text{ to } j$.

Steady State:

$$((\pi_1 P_1) P_2) P_3 \dots = \pi P^\infty \text{ (hard to calculate)}$$

alternative $\pi^* P = \pi^*$ when current state doesn't matter

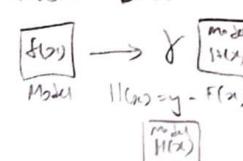
$$\text{and } \sum \pi^* = 1$$

Game Theory
Cost = \$1/gallon $\quad b = \text{total demand}$

BP \$2.50	PP \$2.00
Shell \$2.50 $\frac{1}{2} (\$1.50) = \0.75	Procter (Shell) = 0
Shell \$2.00 $\frac{1}{2} (\$1) = \0.5	Procter (Shell) = 1

Gx Proportional Hazard Model \rightarrow Probability Model

Cox's Model



Validate a simulation by comparing it to real data as much as possible. If the simulation isn't a good reflection of reality, any insights we gain from studying the simulation might not be applicable to reality

Memoryless: the next state of the process doesn't depend on any of its previous states but does depend on current state