

KDE \rightarrow small $\alpha \rightarrow$ lots of bumps and less smooth

Tabular
Nested XML

Excel SQL

CSV/TSV: what if data has commas/tabs?

JSON: Each record can have diff fields

Records can contain records (nested)

Primary key: Unique ID for each entry

Foreign key: cols that ref primary keys in other tables

Ordinal: has orders but no sense of magnitude/intervals

Nominal: no specific ordering

Quantitative: histograms, box plots, rug plots, smoothed

interpolation (KDE) \rightarrow look for spread, shape,

modes, outliers, unreasonable values

nominal/ordinal: bar plots \rightarrow look for skew,

freq and rare categories, or invalid categories

consider grouping categories and repeating analysis

Inferential plot: draw conclusion beyond data

Gaussian kernel:

$$K_{\alpha}(r) = \frac{1}{\sqrt{2\pi}\alpha^2} \exp\left(-\frac{r^2}{2\alpha^2}\right)$$

query data

$$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K_{\alpha}(x - x_i)$$

Big n (many rows): aggregation & smoothing:

avoid over-plotting or use transparency

Big p (many columns): use additional cols

to adjust shape, size, color; combine cols

XML: elem must have open and close tag

unless it is empty $\langle \text{tagname}/\rangle$: must

be properly nested, tag names case-sensitive,

no space allowed btwn \langle and tag name, tag

names must begin w/ letter & contain alphanumeric,

attributes must be in quotes, use < for <

and > for >, must have one root node

that contains all other nodes

str methods

re.findall(pat, str) \rightarrow series.str.findall

re.search(pat, str)

replace, split, contains,
len, \in string

Rest API: Get, Post, Put, delete

SQL: schema: desc of cols, types and constraints

- instance: data satisfying the schema

- attribute (col), Tuple (record, row)

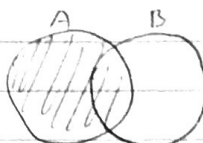
- DELETE from.

- UPDATE table SET gpa = 1.0 + gpa

- Sorting: Order by gpa DESC, name ASC

- Need to use WHERE rating IS NOT NULL
to find null vals or IS NULL

- Can use boolean logic in where clause



SELECT col FROM Table A
LEFT JOIN Table B B
ON A.key = B.key



SELECT col FROM Table A
LEFT JOIN Table B B
ON A.key = B.key
WHERE B.key IS NULL



WHERE A.key IS NULL
OR B.key IS NULL



SELECT col from Table A
FULL OUTER JOIN Table B B
ON A.key = B.key



SELECT col FROM Table A
INNER JOIN Table B B
ON A.key = B.key

loss func: $L(\theta, y) = (y - \theta)^2$ L^2 loss

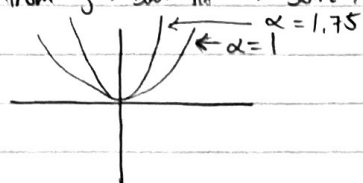
\uparrow predicted \uparrow observed

$L(\theta, y) = |y - \theta|$ L^1 loss, abs loss

$$L_{\alpha}(\theta, y) = \begin{cases} \frac{1}{2}(y - \theta)^2 & |y - \theta| < \alpha \\ \alpha(|y - \theta| - \frac{\alpha}{2}) & \text{otherwise} \end{cases}$$

$\theta = y \rightarrow$ good fit \rightarrow no loss

θ far from $y \rightarrow$ bad fit \rightarrow some loss



$$\log_b M \cdot N = \log_b M + \log_b N \quad \log_b b^k = k$$

$$\log_b \left(\frac{M}{N}\right) = \log_b M - \log_b N \quad b^{\log_b(k)} = k$$

$$\log_b Mk = k \cdot \log_b M \quad b > 1 \text{ \& } M, N > 0$$

$$\log_b(1) = 0 \quad \log_b b = 1$$

HTTP Status

100s Informational
200s Success
300s Redirection or conditional action
400s Client Error
500s Internal Server Error or Broken Request

Average Loss

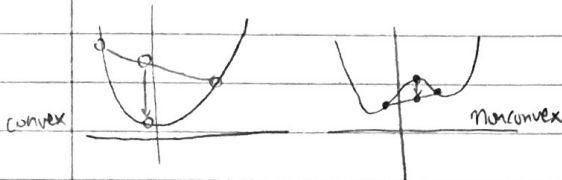
$$L(\theta, D) = \frac{1}{n} \sum_{i=1}^n L(\theta, y_i)$$

L^1 loss \rightarrow less sensitive, not smooth

L^2 loss \rightarrow sensitive, but smooth

Huber \rightarrow less sensitive but extra param

second deriv has to be positive to be a min



$$\text{minimize } L\theta = (1 - \log(1 + \exp(\theta)))^2$$

$$\frac{\partial}{\partial \theta} L(\theta) = \frac{\partial}{\partial \theta} (1 - \log(1 + \exp(\theta)))^2$$

$$= 2(1 - \log(1 + \exp(\theta))) \frac{\partial}{\partial \theta} (1 - \log(1 + \exp(\theta)))$$

$$= 2(1 - \log(1 + \exp(\theta))) (-1) \frac{\partial}{\partial \theta} \log(1 + \exp(\theta))$$

$$= 2(1 - \log(1 + \exp(\theta))) \frac{-1}{1 + \exp(\theta)} \frac{\partial}{\partial \theta} (1 + \exp(\theta))$$

$$= 2(1 - \log(1 + \exp(\theta))) \frac{-1}{1 + \exp(\theta)} \exp(\theta)$$

Derivatives

$$f(x) = a^x \rightarrow f'(x) = \ln a \cdot a^x$$

$$y = f(x)g(x) \rightarrow f'(x)g(x) + f(x)g'(x)$$

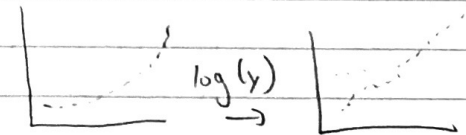
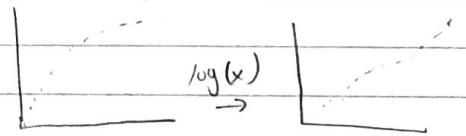
$$\frac{d}{dx} \left[\frac{f(x)}{g(x)} \right] = \frac{f'(x)g(x) - f(x)g'(x)}{g(x)^2}$$

$$\frac{dy}{dx} = \frac{dy}{du} \cdot \frac{du}{dx} = f'(u) \cdot g'(x)$$

$$= f'(g(x)) \cdot g'(x)$$

$$\text{Bayes: } P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

$$= \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)}$$



$$tf(a) + (1-t)f(b) > f(ta + (1-t)b) \quad \forall a, \forall b, t \in [0, 1]$$

$$\text{Gradient Descent Algorithm } L(\theta, y) = \sum_{i=1}^n \theta^2 \cdot x_i^2 - \log(y_i)$$

$$\theta^{(t+1)} \leftarrow \theta^t - \alpha L'(\theta^t) \quad L'(\theta) = \sum_{i=1}^n 2\theta x_i^2$$

Init val θ^0 (zero, random guess)

for t until convergence update

α is the learning rate and not const ($\frac{1}{t}$)

might stop once $\nabla L(\theta^t)$ is small

Data Collection sampling: SRS,

Stratified: SRS within partitions

cluster: divide pop into groups then take SRS of groups

Non Probability Samples: administrative, voluntary, convenience

EPA - structure, granularity, scope, temporality, faithfulness

Pandas - .loc uses the value in the index col

and not the actual row id \rightarrow use .iloc instead

df['w'].value_counts()

sum(), min(), max(), count(), median(), var(), mean(), apply()

df.plot.hist() df.plot.scatter(x='w', y='h')

&, !, ~ df.loc[df['a'] > 10, ['a', 'c']]

rows satisfying boolean and only specific cols

SELECT FROM WHERE GROUP BY HAVING ORDER BY

[DISTINCT] for unique vals

ASC DESC

$$P(B|A) = \frac{P(A \text{ and } B)}{P(A)}$$

Sum values is ascending by default

Average Cross Entropy Loss

$$L(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(\sigma(\phi(x_i)^T \theta)) + (1-y_i) \log(1 - \sigma(\phi(x_i)^T \theta)))$$

If $y_i=1 \Rightarrow -\log(\sigma(\phi(x_i)^T \theta))$

$y_i=0 \Rightarrow -\log(1 - \sigma(\phi(x_i)^T \theta))$

log loss \Rightarrow log of the predicted prob for the true class

Stochastic Gradient Descent

$$\nabla_{\theta} L(\theta) = \frac{1}{|B|} \sum_{i \in B} (\sigma(\phi(x_i)^T \theta) - y_i) \phi(x_i)$$

Both size \uparrow random samples \rightarrow

$$\theta^{(t+1)} \leftarrow \theta^{(t)} - \rho^{(t)} \left(\frac{1}{|B|} \sum_{i \in B} \nabla_{\theta} L_i(\theta) \right)_{\theta=\theta^{(t)}}$$

Linear models

$f_{\theta}(x) = \theta_0 + \theta_1 + \theta_2 x^2$ (still linear in the parameters θ)

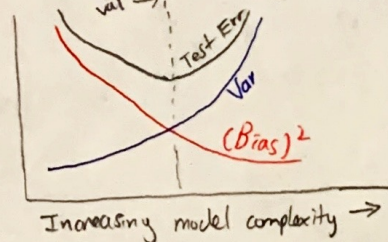
Bias: expected deviation between true and predicted value

Variance \Rightarrow observation \Rightarrow variability of random noise

\Rightarrow model \Rightarrow variability in the predicted value across diff training sets \leftarrow noise \quad bias $^2 \rightarrow$

$$E[(y - f_{\theta}(x))^2] = E[(y - h(x))^2] + (h(x) - E[f_{\hat{\theta}}(x)])^2 +$$

$$E[(E[f_{\hat{\theta}}(x)] - f_{\hat{\theta}}(x))^2] \quad \text{Model Var}$$



Least sq

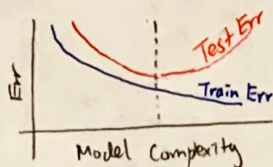
$$\hat{y} = f_{\theta}(x) = \sum_{j=1}^d \theta_j \phi_j(x)$$

linear in the params

features

$$\hat{\theta} = (\Phi^T \Phi)^{-1} \Phi^T Y$$

$$0 = \Phi^T (Y - \Phi \theta)$$



Data warehouse

Extracted from remote sources
Transformed to std schemas
Loaded into the relational data sys

Extract and Load - data in a single sys historical snapshot, isolates analytics

Transform - clean and prep data for analytics in a unified representation
 \rightarrow difficult b/c different schemas, encoding and granularities

Expectations

$$E[X] = \sum_{x \in X} x P(x)$$

$$E[ax + Y + b] = aE[X] + E[Y] + b$$

$$E[XY] = E[X]E[Y] \text{ if } X \text{ and } Y \text{ indep.}$$

$$\text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$$\text{Var}[ax + b] = a^2 \text{Var}[X] + 0$$

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] \text{ if } X, Y \text{ indep.}$$

$$\text{SD}[X] = \sqrt{\text{Var}[X]} \quad \text{SD}[ax + b] = |a| \text{SD}[X]$$

Bernoulli

$$\rightarrow E[X] = p \quad \text{Var}[X] = p(1-p) \quad \frac{\sigma^2}{n}$$

Variance of sample mean decreases at rate $\frac{1}{n}$

$$\text{Standard Error } \text{SD}(\bar{X}) = \frac{\sigma}{\sqrt{n}}$$

Regularization

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n \text{Loss}(y_i, f_{\theta}(x_i)) + \lambda R(\theta)$$

$$R_{\text{ridge}}(\theta) = \sum_{i=1}^d \theta_i^2 \quad R_{\text{Lasso}}(\theta) = \sum_{i=1}^d |\theta_i|$$

distrib weights across related features

analytical soln

small but nonzero weights

encourage sparsity by setting weights to 0

used to select informative features

No analytical solution

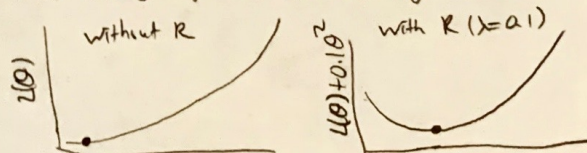
Logistic Regression

$$\hat{P}_{\theta}(y=1|x) = \sigma(\phi(x)^T \theta) = \frac{1}{1 + \exp(-\phi(x)^T \theta)}$$

linear model

use cross entropy loss instead b/c non-convex

If linearly separable \Rightarrow needs regularization



Feature Eng

Quantitative \rightarrow log, normalize

Categorical \rightarrow one-hot-encode

Missing vals

\rightarrow predict vals

\rightarrow binary field for missing

Categorical

\rightarrow True, False, missing

Bag of words

\rightarrow text as long rect of word count

N-gram

\rightarrow The book was well-written

Standardization

\rightarrow each dimension has the same scale

Online Analytics Processing (OLAP)

↳ constructing complex SQL queries

↳ showing views that summarize data across important dimensions

cross tabulation (pivot table)

Cube Operator

↳ generalizes cross-tabulation to higher dimensions.

Slicing \Rightarrow select val for a dimension

dicing \Rightarrow select a range of vals in multiple dimension

Rollup \Rightarrow Aggregate along a dimension

Drill-down \Rightarrow de-aggregating along a dimension

Data Lake

↳ store copy of all data in its natural form

↳ schema on read

↳ lots of dirty data

↳ hard to know what data contain and where data came from