

Week 8 - Feature Selection:

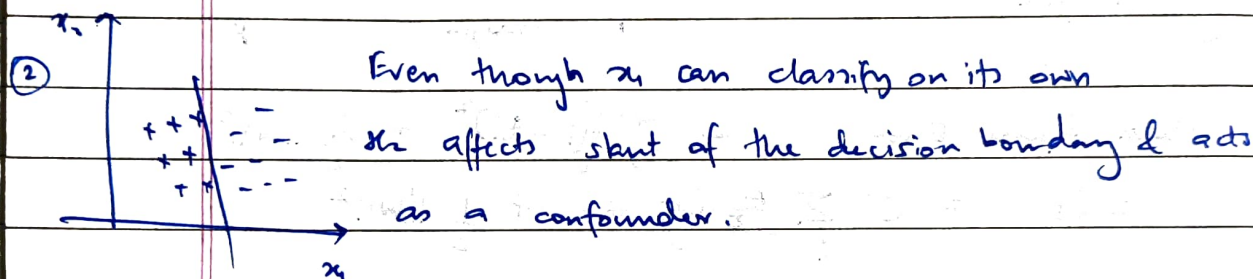
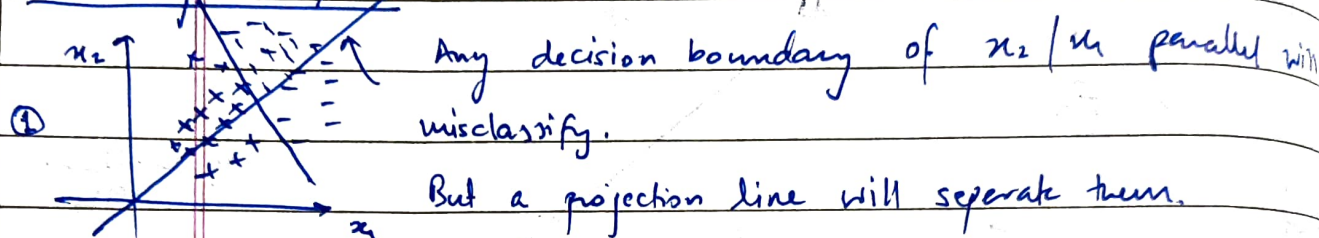
PAGE No.

DATE

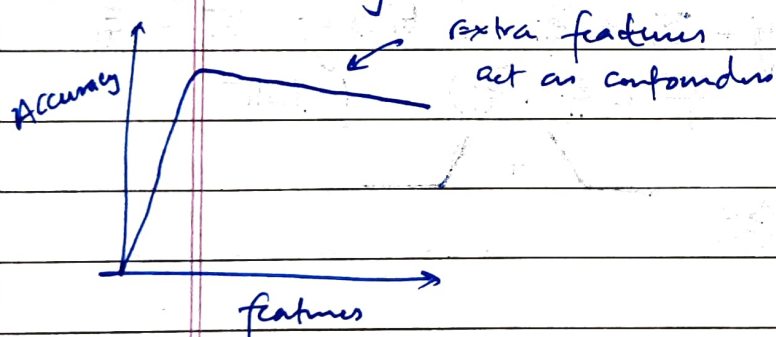
5.1.) Motivation & creating new features:

- Articulate why right features are important & wrong are harmful
- Be able to generate new ones from existing ones
- Be aware of common features for common data types
- Be able to apply feature selection methods.

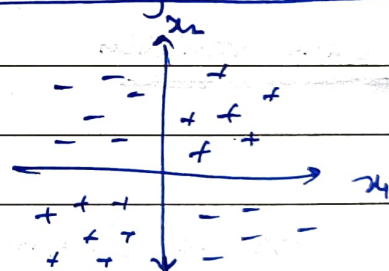
Motivating Examples:



So, we essentially have:-



Generating new features:-

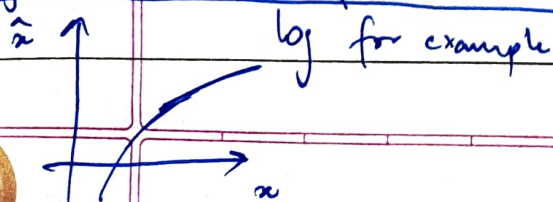


Transforming variables:

Power: $\hat{x} = x^p$ | log: $\hat{x} = \log(x+c)$ | exp: $\hat{x} = e^{ax+b}$

- Transforming big range \longleftrightarrow small range etc.

Objective is to transform to an easier distribution:



objective:

\hat{x} should have an easier dist. in the analytical sense.

Lets say $\hat{x} = f(x)$ assuming they are bijective
 $\Rightarrow f^{-1}(\hat{x}) = x$

PAGE No.	
DATE	/ /

If $p(x)$ and $q(\hat{x})$ are the distributions of the two,

$$q(\hat{x}) = p(f^{-1}(\hat{x})) \left| \frac{d f^{-1}(\hat{x})}{d \hat{x}} \right| \Rightarrow q(\hat{x}) = \frac{p(f^{-1}(\hat{x}))}{\left| \frac{d f(\hat{x})}{d \hat{x}} \right|}$$

Derivation :-

$$P(f(x) \leq \hat{x}) = P(x \leq f^{-1}(\hat{x})) = F(f^{-1}(\hat{x}))$$

Method of transformations :

$Y = g(x)$ g is differentiable | g is strictly increasing.

$$\text{Now, } F_Y(y) = P(Y \leq y) = P(g(x) \leq y) = P(x \leq g^{-1}(y)) = F_X(g^{-1}(y))$$

$$\Rightarrow \frac{d}{dy} F_Y(y) = f_Y(y) = \frac{d}{dy} F_X(g^{-1}(y)) = \frac{f_X(x_1)}{g'(x_1)} [x_1 = g^{-1}(y)]$$

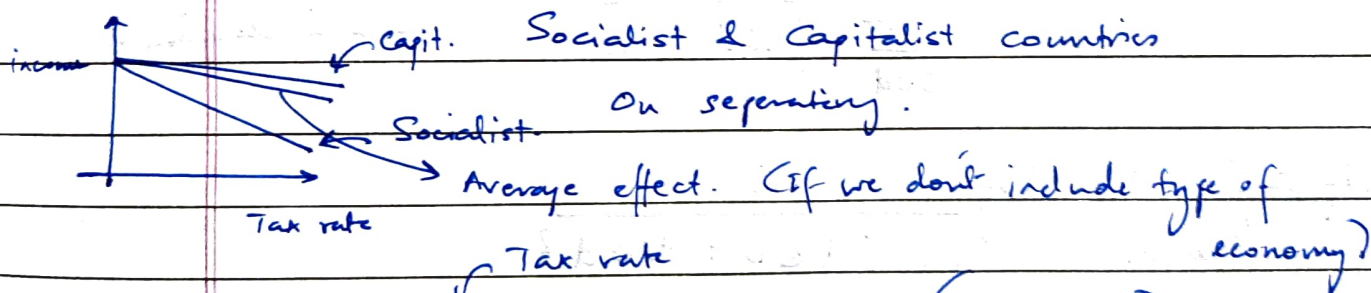
Dummy Variables:

- Convert discrete variables with N levels to $N-1$ binary variables.

Apple, Orange, Banana, Grapes $\Rightarrow W_1 x_1 + W_2 x_2 + W_3 x_3$

+ b (Intercept)

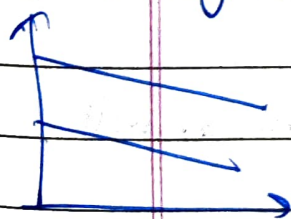
- Interaction variables



$$y = W_1 x_1 + W_2 x_2 + b (+ W_3 x_1 x_2)$$

↑ Social economy type

↑ Models the original graph with different slopes.



If combination is important \Rightarrow Interaction term needed.

$x_1 \leftrightarrow x_2 \leftrightarrow x_1 x_2$ are correlated and can have some effects.

5.2) Features for images:



2D array
+ color
dimension

PAGE No.

DATE

Common Image Features:

- Pixel level statistics

- Gray scale histogram
- Color histograms

- Texture based stats.

- Fourier descriptors
- GLEM

Gray level Co-occurrence Matrix

Hybrid

- Shape based.

• Hu invariant moments

$$\eta_{ij} = \sum \sum (x - \bar{x})(y - \bar{y})^2 I(x, y)$$

5.3) Features for audio & text.

- Meaningful features - power in different frequency bands

Signal \rightarrow Fourier transform \rightarrow Useful data.

Time windows also important. Window \rightarrow Fourier Transform

MFCC is a default feature selection ① Windows of time (overlapping)

② DFT of window - Power Spectral Density

③ Filter bank - log energy

④ DCT (Discrete Cosine Transform)

Common Features to Extract from text:

"The cat was chasing the rat" or "I got lucky in the test today"

"I was not very lucky today"

- Histogram of words from a dictionary (feature vector)

Problem: Too much importance to all words.

$\begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$ are not encoding

- TF-IDF: Term frequency - Inverse Document Frequency

$$tf(c, d) = \frac{f(c, d)}{\sum_{t \in d} f(t, d)}$$

Term frequency over a document

$$idf(t, d) = \log \frac{|D|}{1 + \{d \in D : t \in d\}}$$

↑ corpus

Entire feature = $tf \times idf$

(Also, pretrained deep neural networks help for features as well.)

5.4) Feature reduction 1 of 2:

Features on pre-trained deep neural networks:-

PAGE No.	
DATE	/ /

filter based methods / (Filter / Wrapper / Embedded.)

x_1 x_{1000}

\vdots \vdots

$x_{1000,1}$

$1000,1$

2^{100} subsets so too much

Filter based: For each feature decide keep vs discard \rightarrow Train ML with kept features.

Wrapper: Generate subset, train & validate ML model on subset.

Embedded: - LASSO regularization, Subset selection & ML integrated.

Correlation-based elimination:

Blocks around diagonal - Feature reduction [Correlation based feature reduction]

Utility based subset selection:

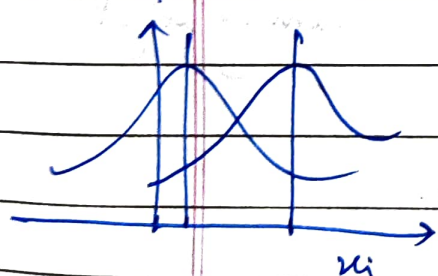
• Regression: Correlation | • Classific: t-test | • AIC & BIC

How the features related to target output.

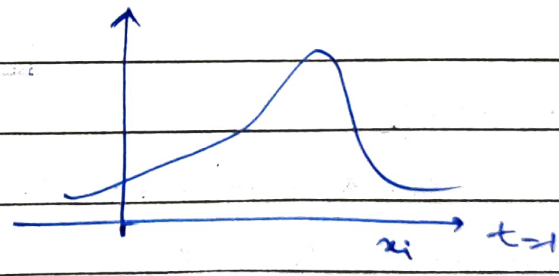
1) x_i correlation \Rightarrow Highly correlated with t may be good subset for prediction.

• Doesn't take in acc. interaction of x_i & x_j • Filter based method

2) Classif: Two class classification.



$t \in \{1, 13\}$



Relative to own variances: Assume they are distributed as Gaussian

t-test formula:
$$\frac{y_i - y_{-1}}{\sqrt{\frac{\sigma_i^2}{n_i} + \frac{\sigma_{-1}^2}{n}}}$$
 width of Gaussian important

AIC & BIC: Akai Info. Criteria & Bayesian Info Criteria

5.5.) Forward Selection & Backward Elimination.

PAGE No.

DATE

$$x_1 \longleftrightarrow x_{100}$$

Subset $S \leftarrow \emptyset$ initially $R = (x_1, \dots, x_{100})$

for $i = 1 \rightarrow n(100)$ we measure the marginal utility for x_j in R

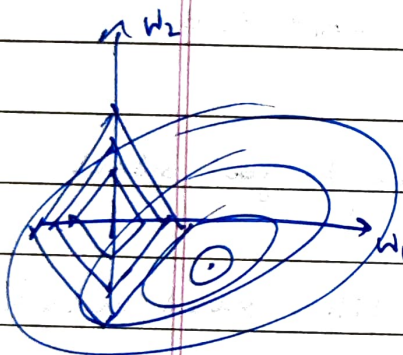
measure marginal utility of x_j in ML model

Include x_i with largest Marginal Utility in S & out of R .

Backward Elimination works in opp. dirⁿ:

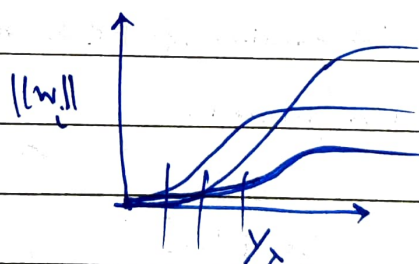
Issues: If x_i & x_j correlated, we may not be able to find that Backward elimination - correlated variables may already be in S

LASSO & elastic net:



Contours will graze to one of the corners.

$$2C(w) = E(w) + \lambda \|w\|_1 + \frac{1}{2} \|w\|_2^2$$



Elastic Net

shrinking weights

Elastic Net:

keeps correlated variables kept in or out together.

Assume: $x_1 = x_2$

$$w_1 x_1 + w_2 x_2 + w_3 x_3$$

we want $w_1 \leq w_2$

but $(w_1 + \alpha) x_1 + (w_2 - \alpha) x_2 + w_3 x_3$ is also the same

but L_1 norm is same, how to diff.? $\Rightarrow L_2$ Norm

L_2 penalty large for $(w_1 + \alpha)$ & $(w_2 - \alpha)$ case.

L1 SVM: (Cor L2 SVM is original SVM)

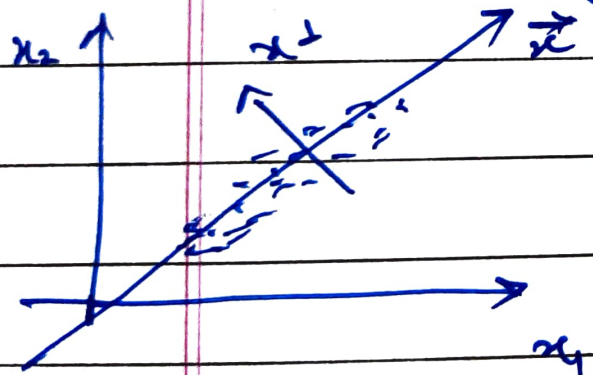
$$\frac{1}{2} \|w\|_2^2 + c \sum_i [1 - t_i y_i]_+$$

PAGE No.	
DATE	/ /

$\lambda \|w\|_1 + c \sum_i [1 - t_i y_i]_+$ for some values of λ you have certain w to be 0.

Principal Component Analysis:

$(N \rightarrow d)$



Captures essence of x_1 & x_2 .
Eigen vectors corresponding to higher
Eigen values. Each captures contribution
from all underlying contribution.