

Week 8 - Feature Selection:

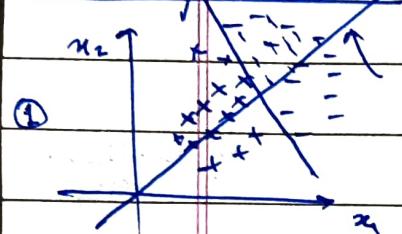
PAGE No.

DATE

5.1.) Motivation & creating new features:

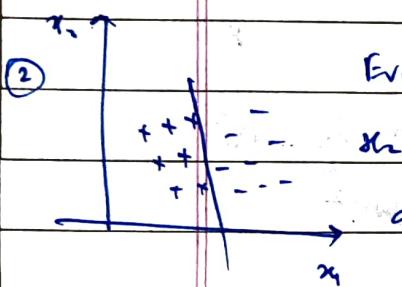
- Articulate why right features are important & wrong are harmful
- Be able to generate new ones from existing ones
- Be aware of common features for common data types
- Be able to apply feature selection methods.

Motivating Examples:



Any decision boundary of x_2/m parallel will misclassify.

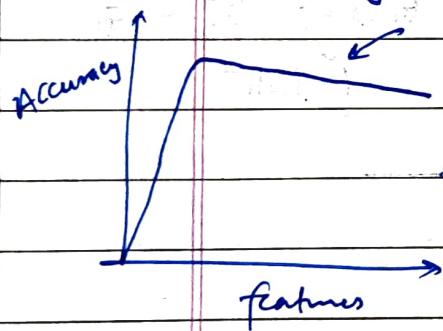
But a projection line will separate them.



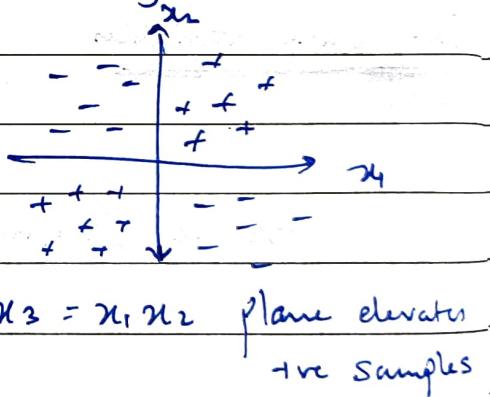
Even though x_1 can classify on its own

x_2 affects shift of the decision boundary & acts as a confounder.

So, we essentially have:-



Generating new features:-



Transforming variables:

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

- Transforming big range \leftrightarrow small range etc.

Objective is to transform to an easier distribution:

$\hat{x} \uparrow$ log for example

objective:

\hat{x} should have an easier dist.

in the analytical sense.

Let's 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,

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

Derivation :-

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

Method of transformations:

$$Y = g(x) \quad g \text{ is differentiable} \quad | \quad g \text{ is strictly increasing.}$$

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

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

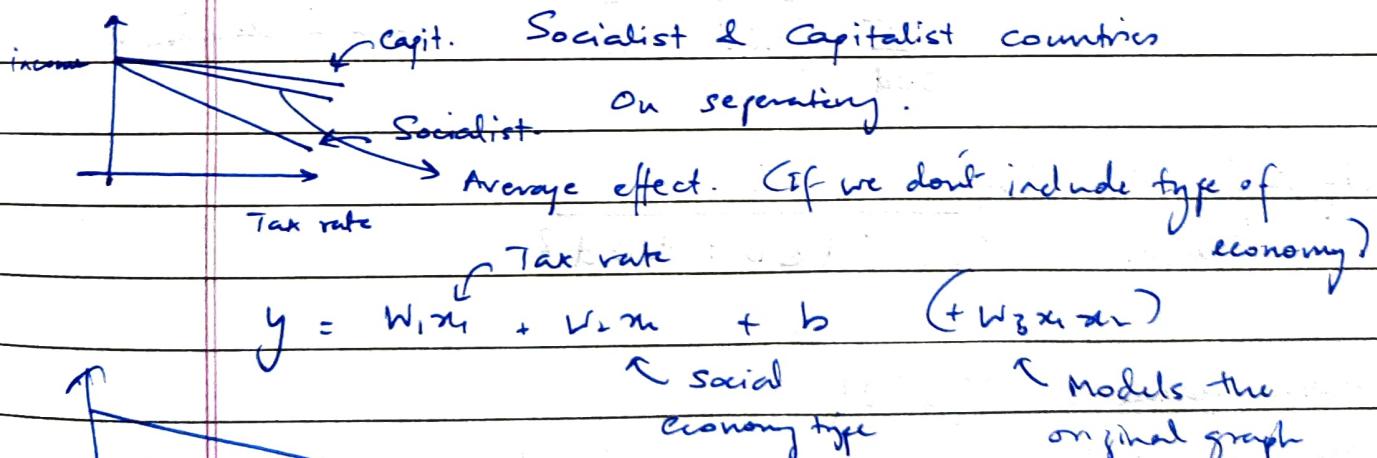
Dummy Variables:

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

Apple, orange, Banana, Grapes $\Rightarrow w_1 x_{1i} + w_2 x_{2i} + w_3 x_{3i}$

+ b
~ 4th

- Interaction variables



~ social

economy type

~ Models the

original graph

with different slopes.

If combination is important \Rightarrow Interaction term needed.

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

5.2) Features for images:

Common Image Features:

- Pixel level statistics

- Gray scale histogram
- Color histograms



2D array
+ color dimension

PAGE No.

DATE

- Texture based stats.

- Fourier descriptors
- GLCM

Gray level co-occurrence Matrix

- Hybrid

- Shape based.

$$\text{Hu invariant moments} \Rightarrow Y_{ij} = \sum \sum (x - \bar{x})(y - \bar{y}) I(x, y)$$

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.

- TF-IDF: Term frequency - Inverse Document frequency

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

Term frequency
over a document

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

Final feature = $tf \times idf$

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

5.4) Feature reduction 1 of 2:

PAGE No.	/ / /
DATE	/ / /

features on pre-trained deep neural networks:-

filter based methods

(Filter / Wrapper / Embedded.)

x_{11} x_{1000}

: ...

x_i

1000, 1

$x_{4000, 10}$

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:

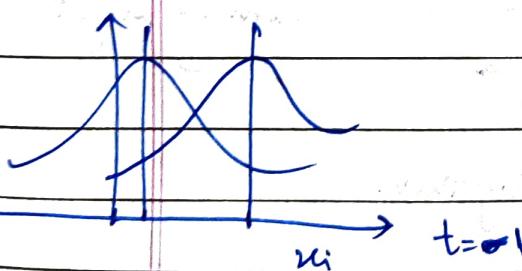
o Regression: Correlation | o Classification: t-test | o AIC & BIC

How the features related to target output.

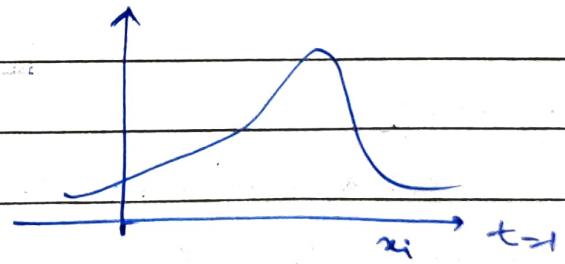
1) x_i correlation \Rightarrow Highly correlated with t may be good subset for prediction.

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

2) Classification: Two class classification.



$$t \in \{0, 1\}$$



Relative to own variances: Assume they are distributed as Gaussian

t-test formula:

$$\frac{y_1 - y_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_0^2}{n_0}}}$$

width of Gaussian

important

AIC & BIC: Akaike Info Criteria & Bayesian Info Criteria

5.5.) Forward Selection & Backward Elimination.

$$x_1 \leftrightarrow x_{100}$$

PAGE No.	/ /
DATE	/ /

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

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

for x_i in R

measure marginal utility of x_i in ML model

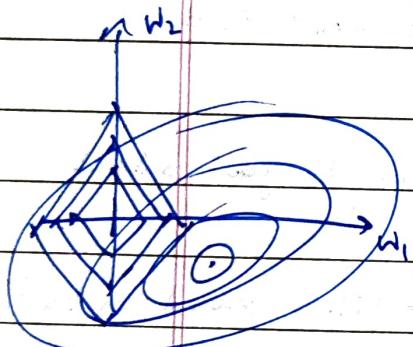
Include x_i with largest Marginal Utility in S & out of R .

Backward Elimination works in opp. dirn:

Issues: If x_i & x_j correlated, we may not be able to find that

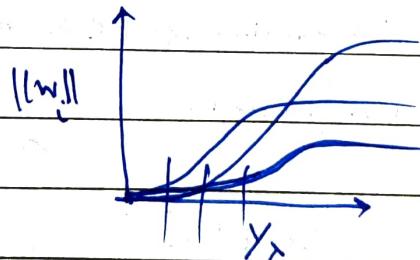
Backward elimination - correlated variables may already be in S

LASSO & elastic net:



Contours will graze to one of the corners.

$$L(w) = E(w) + \lambda \|w\|_1 + \lambda_2 \|w\|_2$$



Elastic Net
summing weights

Elastic Net:

keeps correlated variables kept in or out together.

$$\text{Assume: } x_1 = x_2$$

$$w_1 x_1 + w_2 x_2 + w_3 x_3 \quad \text{we want } w_1 \leq w_2$$

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

but L1 norm is same. How to diff.? \Rightarrow L2 Norm

L2 penalty large for $(w_1+\alpha)$ & $(w_2-\alpha)$ case.

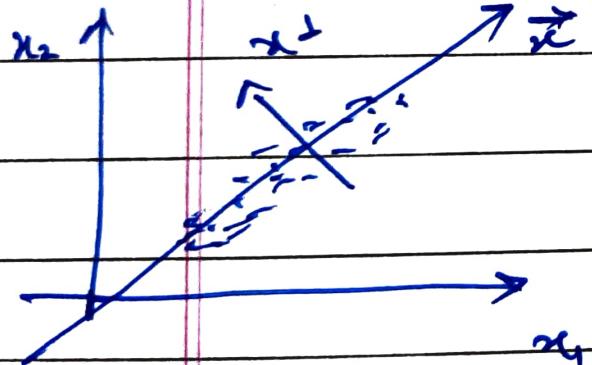
L1 SVM: (Cor L2 SVM is original SVM)

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

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.

(7-1-1) Introduction to Neural Networks (12:21)

• Objectives

- Adding layers ↑ modelling power
- Mathematical expression - Backpropagation algorithm
- List ways to improve backprop.

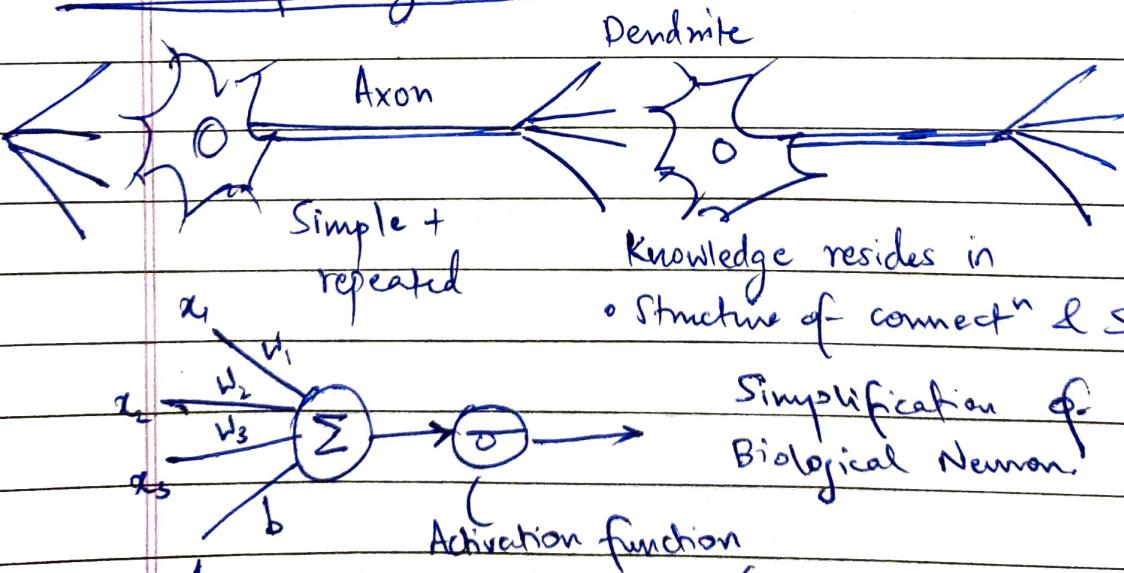
Increasing non linearity in models:

1. Linear models (Regressor / Classifier) $L(\sigma(Wx))$
2. Support vector machines → Fixed features $L(\sum_i \alpha_i K(x, x_i))$
3. Neural Networks → Trainable features $L(\sigma(W_2 \sigma(W_1 x)))$
[or multiple layers (But not too deep)]
4. Deep Neural Networks $L(\sigma(W_n \dots (\sigma(W_1 x))))$
[σ is a non linear function, not necessarily sigmoid]

(7-1-2) Layered functional representation of NN :

W are matrices for neural networks

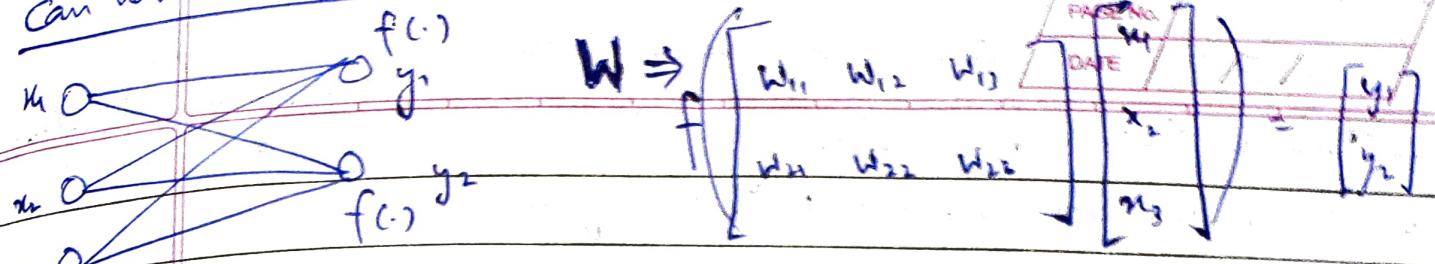
Transform feature vectors into other vectors.

Update based on $\delta L / \delta W_{ij}$ (chain rule of derivatives)Structure of Biological Neuron:

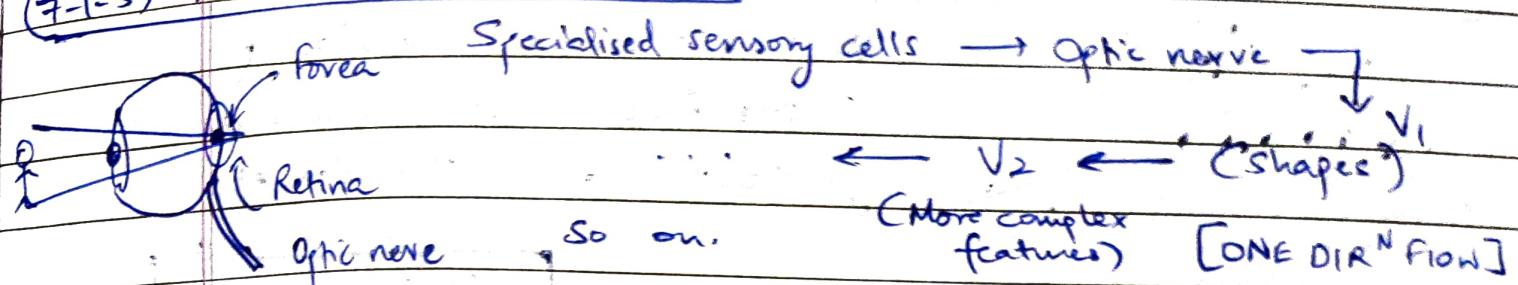
$$w_2(w, x) = (w_2 w_1)x = w_3 x$$

Which is why every layer but the last has some non linear function ' σ '

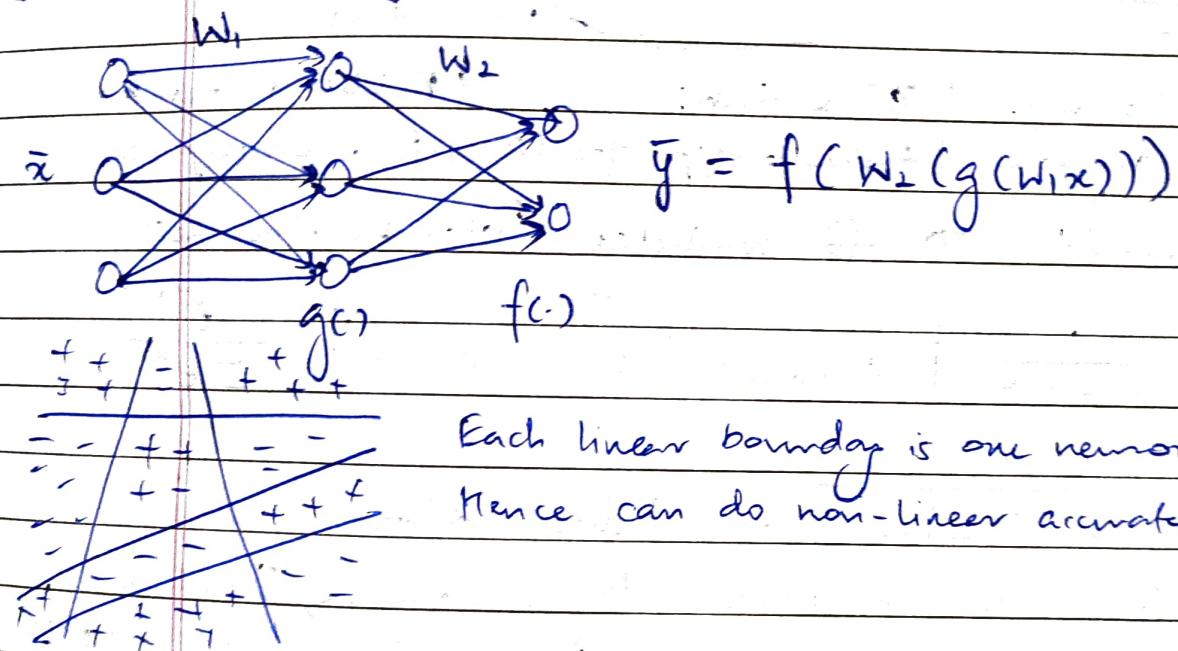
Can have multiple outputs as well:



(7-1-3) Mammalian Visual Cortex:

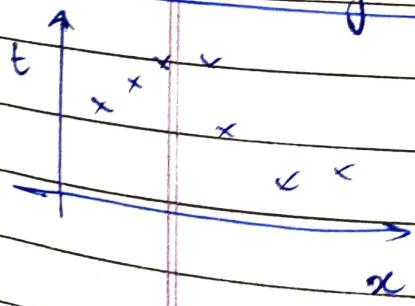


(7-1-4) Introducing Hidden layer in NN :



Each linear boundary is one neuron.
Hence can do non-linear accurate classification

What hidden layers are doing:-



Basic $f(x) = c$

Lin. $f(x) = W_1x + W_0$

Non lin. $f(x) = W_1\phi_1x + W_2\phi_2x + b$

SVM $f(x) = \sum W_i K(x, x_i)$

NN $f_{NN} = \text{Summation of activation function with weights \& biases.}$

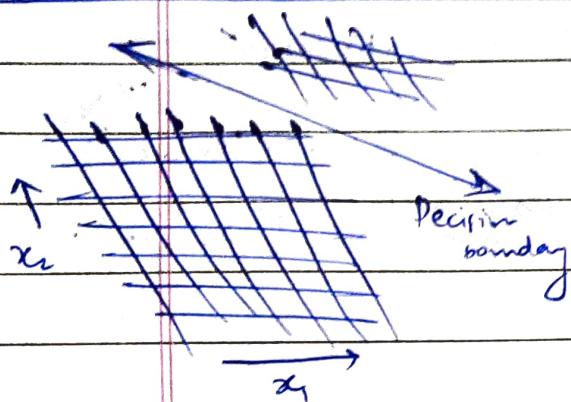
Repeated combination of canonical function with W \& b \leftarrow layers.

Universal Approximation Theorem:

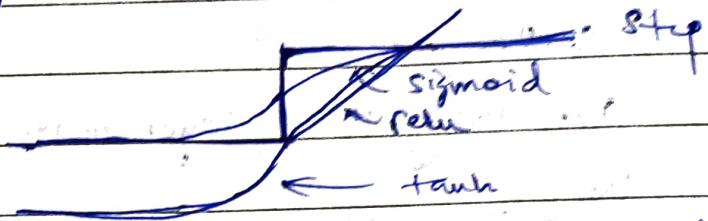
- $f \in C^k(\sigma(W, x))$ can approximate f with error ϵ in a compact interval any smooth $f(x)$
 - Provided size of W is arbitrary
 - σ is also smooth, but NOT a polynomial

PAGE NO.	/ / /
DATE	/ / /

(7-1-5) Activation function:



We need a smoother version of a step function with non zero gradients \Rightarrow Sigmoid etc.



Step: Classification alg.
Not used anymore

σ & tanh:
trainable app. of step

ReLU: Referred as
fast convergence

Softmax: Generalized σ for multiclass classification net

$\frac{\text{sgn}(x)+1}{2}$	$\frac{1}{1+e^{-x}}$	$\tanh(x)$	$\max(0, x)$	$\frac{e^{x_i}}{\sum e^{x_i}}$
-----------------------------	----------------------	------------	--------------	--------------------------------

Problems with too many layers:

- Too many parameters
- Gradient dilution

Hyperparameter: No. of hidden layers. \rightarrow Input dim \rightarrow Dimension of layers. \rightarrow Output dim fixed

The training process is gradient descent

Gradient of vectors: ∇f $\nabla = \frac{\partial L}{\partial w_{ij}}$ matrix

(7-2-1) chain rule & backpropagation:

$$f(x) = g(h(x))$$

$$f'(x) = g'(y) h'(x) \quad y = h(x), \text{ so on.}$$

$$\text{eg. } \frac{d(\sin(x^2))}{dx} = \cos(x^2) \cdot 2x$$

Algorithm:

- Make a forward pass & store partial derivatives
- During backward pass, multiply δ [Library PAGE NO. 22 DATE _____ Libraries Axis]

(7-2-2) Backpropagation for vector variables :-

$$f(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} f_1(x_1, x_2, x_3) \\ f_2(x_1, x_2, x_3) \end{bmatrix} \xrightarrow{x_j}$$

$$\mathbf{J}(f) = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \end{bmatrix},$$

$$\mathbf{J}(f) = \begin{bmatrix} \frac{\partial f_i}{\partial x_j} \end{bmatrix} \quad \begin{array}{l} i - \text{row index} \\ j - \text{col index} \end{array} \quad \begin{bmatrix} f_i \times x_j \end{bmatrix} \quad |I| \times |J|$$

Questions:

- Scale all weights & biases by factor? \Rightarrow Sign no change
NOT UNIQUE !!
- Gradients in deep neural networks?
- $\prod_i W_i$ - if $|W| < 1$ then we have vanishing ∇ \nwarrow Need to regularize
- if $|W| > 1$ then we have exploding ∇ \nearrow (depends on init)

Vanilla gradient descent:

- Iteration loop n

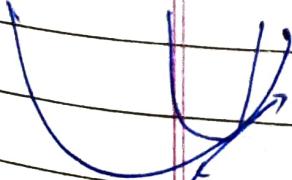
- Sample loop i

$$\nabla_{\theta} L \leftarrow \nabla_{\theta} L + \nabla_{\theta} L_i$$

$$\theta \leftarrow \theta - n \nabla_{\theta} L$$

Slow model: We have to go through all training samples before we can update even once

(7-2-3) Role of step size & learning rate :



Same gradient \Rightarrow Same step for same η .

- Same value + same gradient ∇f :

- Different Hessian $\nabla^2 f$ \Rightarrow Different step sizes needed

\Rightarrow We need learning rate scheduling or some function to model η

Issues with gradient descent:

- Need to find good step size η
- lots of computation before each update
- can get stuck in local minima

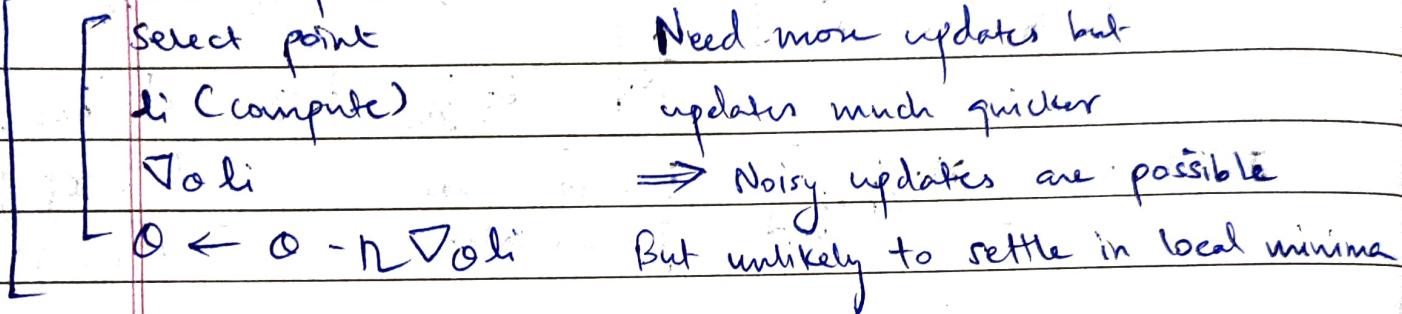
PAGE No.	/ / /
DATE	/ / /

Solutions: Stochastic GD or Batch GD

1.) Stochastic Gradient Descent:

- pick any training point randomly

Iteration loop:



2.) Batch gradient descent:

[Hybrid of Vanilla GD & Stochastic GD] Most popular. (GPU utilization)

- Batch formation loop:
 - shuffle training pt
 - divide into batches

Epoch loop

- Batch loop

 └ Batch

$\nabla \theta$ ∇_{batch}

$$\theta \leftarrow \theta - \eta \nabla_{\text{batch}}$$

key: increasing batch till GPU error.

Faster than stock

Double derivative speed up of Hessian & Jacobian:

$$\text{let } f(x) = ax^2 + bx + c$$

Assuming $a < 0$ [Inverse paraboloid]

Minima at: $-b/2a$

For any 'x' the perfect step :-

$$(-\frac{b}{2a} - x) = -\left(\frac{2ax + b}{2a}\right) = -\frac{f'(x)}{f''(x)}$$

PAGE NO.

DATE

Same logic for:

$$x \leftarrow x - \frac{1}{f''(x)} H(f(x)) \nabla(f(x))$$

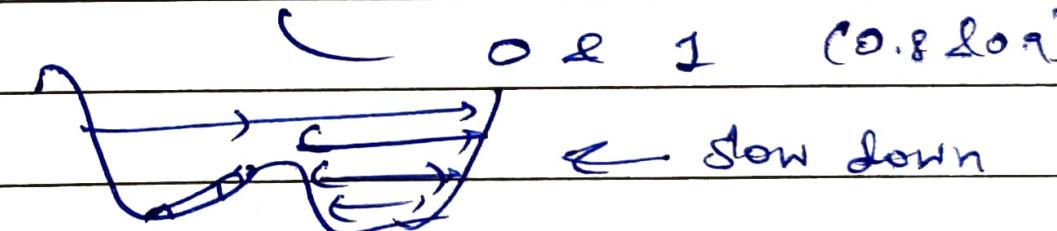
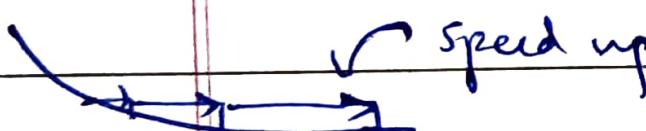
[Not used but approximations of these speed up]

→ Note that it is NOT important to find global minima.

(Since Global Min of training data \neq That of testing data)

Adding momentum instead of 2nd derivative :-

$$\theta^n \leftarrow \theta^{n-1} - \eta \nabla L + \alpha (\theta^{n-1} - \theta^{n-2}) \quad \alpha =$$



Like momentum intuition

[L2 - Weight decay] in NN