

EE353 AIDS

01 August 2024 23:37

EE 353 * scalar vector operation

$a \in \mathbb{R}$

$\bar{x} \in \mathbb{R}^{N \times 1}$

$a\bar{x} = \begin{bmatrix} ax_1 \\ ax_2 \\ ax_3 \end{bmatrix}$

$\bar{x} + a = \begin{bmatrix} x_1 + a \\ x_2 + a \\ x_3 + a \end{bmatrix}$

$\bar{x} + a \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + a \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$

$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \odot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \end{bmatrix}_{N \times 1}$

$\bar{x} \cdot \bar{y} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \\ \vdots \\ x_N y_N \end{bmatrix} \quad \bar{x} \cdot \bar{y} = \langle \bar{x}, \bar{y} \rangle$

Subspaces

vector space

$\left[\begin{array}{l} \forall v_1, v_2 \in V \quad \text{if } c_1 v_1, c_2 v_2 \in V \\ \quad \text{if } c_1 v_1 + c_2 v_2 \in V \end{array} \right]$

Subspace of $V = W$

$\forall w_1, w_2 \in W \subset V$

$\begin{array}{l} \text{① } w_1 + w_2 \in W \\ \text{② } c w_1 \in W \end{array}$

$c_1, c_2 \in \mathbb{R}$ all vector of form $c_1 w_1 + c_2 w_2$ form subspce

e.g. a plane in \mathbb{R}^3 space V defined by a linear comb of w_1, w_2 is W

$$X \in \mathbb{R}^{M \times N}$$

$$Y \in \mathbb{R}^{M \times N}$$

$$Z = X + Y \in \mathbb{R}^{M \times N}$$

$$Z_{ij} = \sum_{k=0}^m x_{ik} y_{kj}$$

$$Z_{ij} = x_{ij} + y_{ij}$$

$$Z = X \odot Y \in \mathbb{R}^{M \times N}$$

$$Z_{ij} = x_{ij} y_{ij}$$

transpose, dot product

$$X^T \in \mathbb{R}^{N \times M}$$

$$X^T X = I_{M \times M} = \begin{bmatrix} 1 & & & \\ & 0 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix}$$

$\in \mathbb{R}^{M \times M} \rightarrow \mathbb{R}^{M \times M}$

↓ (m times)

If X is invertible \Rightarrow full rank \equiv # independent rows or columns of Matrix.

Pseudo Inverses

$$M \neq N \rightarrow \mathbb{R}^{M \times M}$$

$$X^+_{N \times M} = (X^H X)^{-1} X^H \in \mathbb{R}^{N \times M}$$

conjugate transpose

$$X^+ X = (X^H X)^{-1} (X^H X) = I_{N \times N}$$

Eigen decomposition

$$\bar{A}_{N \times N} \underbrace{\bar{V}_i}_{\mathbb{R}^{N \times N}} = \bar{D}_i \bar{V}_i$$

$i \in \{1, \dots, N\}$

eigen value eigen vector unit vector

$$\bar{A} = \bar{Q} \bar{\Lambda} \bar{Q}^{-1}$$

$$\bar{V}_i = \begin{bmatrix} v_1 & \dots & v_N \end{bmatrix}$$

$$\bar{D}_i = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{bmatrix}$$

$$v_i^T v_i = \delta_{ii}$$

$$\delta_{ij} = 0 \quad \forall i \neq j$$

$$\delta_{ii} = 1 \quad \forall i = j$$

$$A^{-1} = (\bar{\theta} \bar{\pi} \bar{\theta}^{-1})^{-1} \quad (AB)^{-1} = \bar{B}^T A^{-1}$$

PAGE NO.: 1
DATE: 1/1/2024

Implication of $\Delta i = 0 \Rightarrow$ rank deficiency

Tensors

$$\alpha \in \mathbb{R} \rightarrow \bar{\alpha} \in \mathbb{R}^N ; \bar{\gamma} \in \mathbb{R}^{M \times N}$$

$$T \in \mathbb{R}^{M \times N \times P \times \dots}$$

$$\text{transpose}(T, [3, 1, 2, 0])$$

$\uparrow \quad \downarrow \quad \leftarrow$
 $[0, 2, 1, 3]$

$$T \in \mathbb{R}^{M \times N \times P} \rightarrow \text{transpose}(T) = T' \in \mathbb{R}^{P \times N \times M}$$

In this we have to define the order of dimension of array

$$\text{transpose}(T, [0, 2, 1]) \rightarrow \begin{bmatrix} 0 & 1 & 2 \\ 1 & 0 & 2 \\ 0 & 2 & 1 \\ 1 & 2 & 0 \\ 2 & 1 & 0 \\ 2 & 0 & 1 \end{bmatrix}$$

possibilities

Functions &

→ continuity of function:

$$\text{If } \lim_{\Delta x \rightarrow 0} f(x - \Delta x) = f(x + \Delta x) \Rightarrow \text{continuous @ } x$$

→ smooth function → derivative is continuous @ x

Lipschitz continuity: $|f(x_1) - f(x_2)| \leq k|x_1 - x_2|$

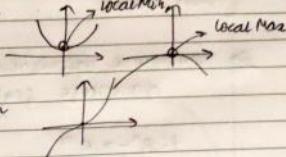
$K \in \mathbb{R}^+$

(Lipschitz Continuity)

1/8/24
PAGE NO.: / / DATE: / /

→ derivatives of functions
 $f'(x) = \frac{df(x)}{dx} = \lim_{\Delta x \rightarrow 0} \left(\frac{f(x+\Delta x) - f(x)}{\Delta x} \right)$

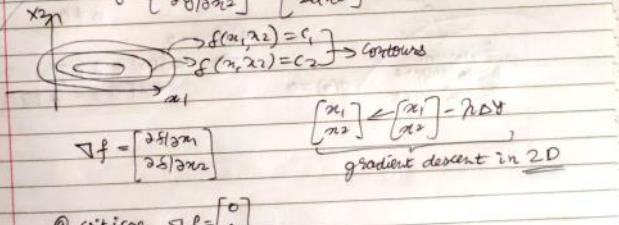
→ critical point of a function
 $f'(x) = 0$
 ↗ Maxima
 ↗ Minima
 ↗ Inflection



$f''(x) = 0 \rightarrow f''(x) > 0 \Rightarrow \text{Minima}$
 $f''(x) < 0 \Rightarrow \text{Maxima}$
 $f''(x) = 0 \Rightarrow \text{Inflection}$

Multivariate function
 $y = f(x_1, x_2)$
 eg. $y = ax_1^2 + bx_2^2$

$\nabla y = \begin{bmatrix} \frac{\partial y}{\partial x_1} \\ \frac{\partial y}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2ax_1 \\ 2bx_2 \end{bmatrix}$



$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} - \beta \Delta y$
 gradient descent in 2D

@ critical $\nabla f = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$
 eg. $\nabla y = \begin{bmatrix} 2ax_1 \\ 2bx_2 \end{bmatrix} = 0$ for @ critical pt.
 if $a > 0$, b < 0
 $(0,0)$ is min for x_2
 min for x_1
 → saddle point

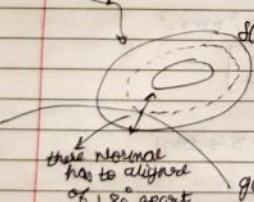
1/8/24
PAGE NO.: / / DATE: / /

$H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 2a & 0 \\ 0 & 2b \end{bmatrix}$
 (Hessian matrix)

constrained optimization using Lagrange multiplier
 maximize $f(x)$; subject to $g(x) = 0$

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $x \in \mathbb{R}^n \quad f(x) \in \mathbb{R}^m$ $J_{ij} = \frac{\partial f_i}{\partial x_j}$ Jacobian Matrix	$f: \mathbb{R}^n \rightarrow \mathbb{R}$ $x \in \mathbb{R}^n \quad f(x) \in \mathbb{R}$ $(H_f)_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}$ $H(f(x)) = J(\Delta^2 f(x))^T$
---	---

$f(x) = c$
 $L(x) = f(x) + \lambda g(x)$
 $\lambda \neq 0$
 $\nabla L(x) = 0$
 $\Rightarrow \nabla f(x) + \lambda \nabla g(x) = 0$
 $\nabla f(x) = -\lambda \nabla g(x)$



$\sum m_i + m/n - m = \min n + m$
 $\sum m_i - \min n \times n = m$

EE 353 \Rightarrow Eigen vector are orthonormal for symmetric Matrix

Probability
random variable
 $P(X=x)$

Discrete probability ; pmf [Probability mass function]
Cont' probability ; pdf [Probability density fn]

→ ① Bernoulli dist.

$x \in \{0, 1\}$ $E(x) = \sum_x x P(X=x) = \bar{x}$
 $\bar{x} = P(X=1)$ $E((x-\bar{x})^2) = \bar{x}(1-\bar{x})$

$\text{Bern}(x|\bar{x}) = \bar{x}(1-\bar{x})^{1-x} = \begin{cases} \bar{x} & ; x=1 \\ 1-\bar{x} & ; x=0 \end{cases}$

② Binomial distⁿ
N times Binary

$y \in \{0, 1, \dots, N\}$
next space
 $\text{Bin}(y|\bar{x}, N) = \binom{N}{y} \bar{x}^y (1-\bar{x})^{N-y}$
 $\rightarrow (\text{out of } N, y \text{ occurs})$

5/18/24
PAGE NO. / / DATE / /

$E(f(y)) = \sum_{y=0}^N \text{Bern}(y|\bar{x}, N) y$

$E(f(y)) = \sum_y f(y) \cdot P(y)$

Entropy(y) = $E[-\log p(y)]$

\Rightarrow for Bernoulli distⁿ. Entropy(y) = $-\frac{1}{\bar{x}} \log \bar{x} - \left(1 - \frac{1}{\bar{x}}\right) \log \left(1 - \bar{x}\right)$
 $= -\log \bar{x}$

$-\log \bar{x} - 1 + \log(1-\bar{x}) - 1 = 0$
 $\log \left(\frac{1-\bar{x}}{\bar{x}}\right) = 2$

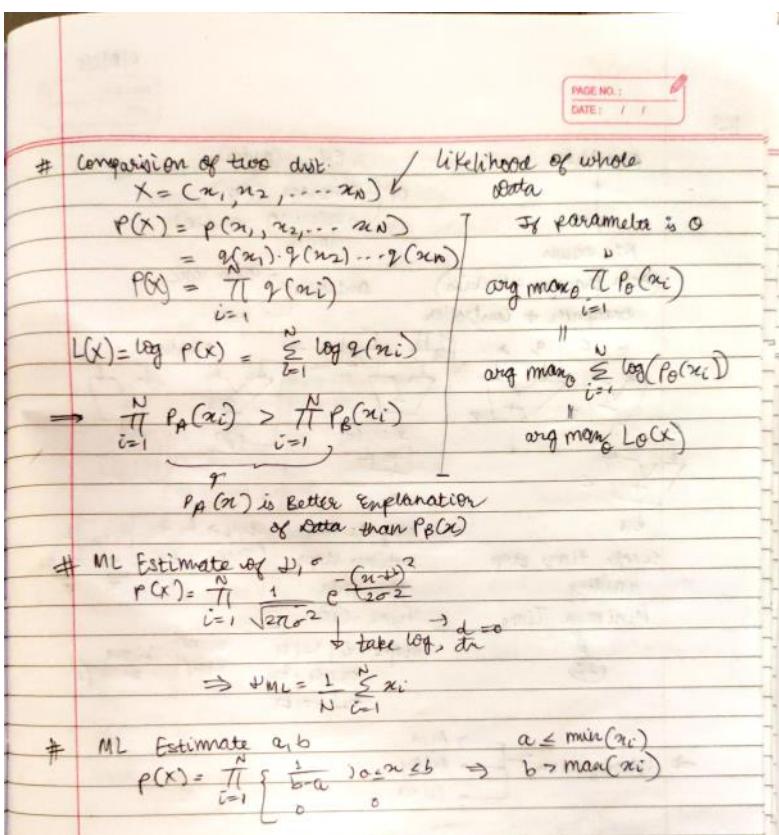
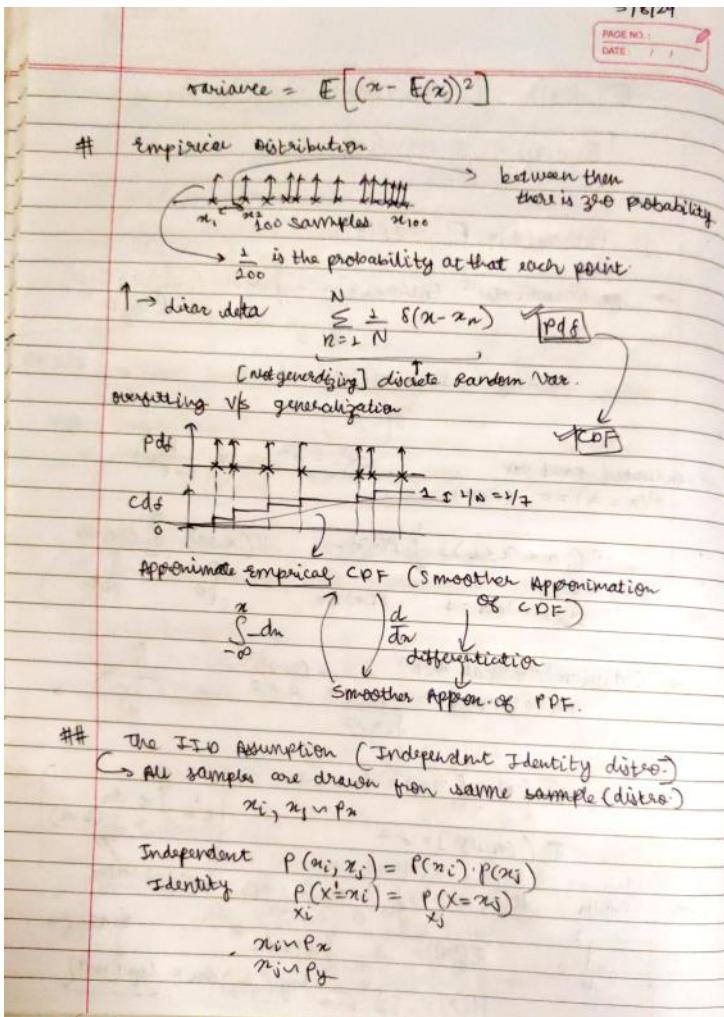
Continuous Rand Var
 $P(x=a) = 0$ (PDF)
 $P(a < x \leq b) = \int_a^b f(x) dx$
 $\int_{-\infty}^{\infty} f(x) dx = 1$
 $f(x) \geq 0$
 $\int_a^b f(x) dx \geq 0$

$N(x'; \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-x')^2}{2\sigma^2}}$

$E(x) = \int_{-\infty}^{\infty} x \cdot p(x) dx = \mu$
 $E[(x-\mu)^2] = \sigma^2$

$\int_a^b x \cdot \frac{1}{b-a} dx = \frac{1}{b-a} \int_a^b x dx = \frac{1}{b-a} (\mu - a)$

Uniform $U(a; a+b) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 & \text{otherwise} \end{cases}$
 $E(g(x)) = \frac{1}{b-a} \int_a^b g(x) dx$
 $H(x) = -\int_a^b \frac{1}{b-a} \log \left(\frac{1}{b-a}\right) dx = \log(b-a)$



Independent:

- **Mathematical Definition:** Two random variables X_1 and X_2 are independent if the joint probability of them occurring together equals the product of their individual probabilities:

$$P(X_1 = x_1, X_2 = x_2) = P(X_1 = x_1) \cdot P(X_2 = x_2)$$

For more than two variables, this extends to:

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n) = \prod_{i=1}^n P(X_i = x_i)$$

This means the occurrence of one event doesn't affect the probability of the others.

Identically Distributed:

- **Mathematical Definition:** Random variables X_1, X_2, \dots, X_n are identically distributed if they all have the same probability distribution function (PDF) or probability mass function (PMF):

$$F_{X_1}(x) = F_{X_2}(x) = \dots = F_{X_n}(x) = F_X(x)$$

where $F_X(x)$ is the cumulative distribution function (CDF) common to all X_i . This implies that each random variable has the same mean, variance, and other distributional characteristics.

IID in Practice:

- When we say a set of random variables $\{X_1, X_2, \dots, X_n\}$ is IID, it means:

1. **Independence:** X_i and X_j for any $i \neq j$ satisfy:

$$P(X_i \leq x_i \text{ and } X_j \leq x_j) = P(X_i \leq x_i) \cdot P(X_j \leq x_j)$$

2. **Identical Distribution:** Each X_i shares the same distribution function $F_X(x)$.

Why It Matters:

- **Central Limit Theorem (CLT):** The concept of IID is central to the CLT, which states that the sum (or average) of a large number of IID random variables will be approximately normally distributed, regardless of the original distribution.

If X_1, X_2, \dots, X_n are IID with mean μ and variance σ^2 , then as $n \rightarrow \infty$:

$$\frac{1}{n} \sum_{i=1}^n X_i \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right)$$

- **Law of Large Numbers (LLN):** For IID random variables, the sample mean converges to the population mean as the sample size increases.

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu$$

Mathematical Explanation

Suppose you have a dataset consisting of n independent and identically distributed (IID) observations X_1, X_2, \dots, X_n from some probability distribution with probability density function (PDF) $f(X; \theta)$, where θ represents the parameters of the distribution.

Likelihood Function

The likelihood function for the entire dataset is the joint probability (for discrete data) or joint density (for continuous data) of all observations. For IID data, because each X_i is independent and identically distributed, the likelihood $L(\theta)$ can be expressed as:

$$L(\theta) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n | \theta)$$

Since the data points are independent, this joint probability (or density) factorizes into the product of the individual probabilities (or densities):

$$L(\theta) = \prod_{i=1}^n P(X_i = x_i | \theta) = \prod_{i=1}^n f(x_i; \theta)$$



Key Insight

The factorization of the likelihood function is what allows us to simplify many complex problems in statistics and machine learning. When data points are IID, the complexity of calculating the joint likelihood of the entire dataset reduces to calculating the product of individual likelihoods.

This factorization is crucial in:

- **Maximum Likelihood Estimation (MLE):** Where you find the parameter θ that maximizes $L(\theta)$.
- **Bayesian Inference:** Where you update the prior distribution based on the likelihood of observed data.

Likelihood function

Likelihood Function

Given a dataset $X = \{x_1, x_2, \dots, x_n\}$ and a statistical model parameterized by θ , the likelihood function $L(\theta | X)$ represents the probability of observing the data X given the parameter θ . For a set of independent observations from a distribution with parameter θ , the likelihood function is:

$$L(\theta | X) = \prod_{i=1}^n f(x_i | \theta)$$

where $f(x_i | \theta)$ is the probability density function (PDF) or probability mass function (PMF) of the distribution.

Maximizing the Likelihood Function

- **Objective:** We seek to find the value of θ that maximizes $L(\theta | X)$. This value of θ is called the **Maximum Likelihood Estimate (MLE)**.
- **Mathematical Expression:**

$$\hat{\theta} = \arg \max_{\theta} L(\theta | X)$$

Here, $\hat{\theta}$ is the estimate of the parameter θ that maximizes the likelihood function.

Summary

- $\arg \max_{\theta} L(X)$ denotes the parameter value θ that maximizes the likelihood function given the data X .
- Finding this value involves differentiating the likelihood function (or its logarithm), setting the result to zero, and solving for θ .
- The result is known as the Maximum Likelihood Estimate (MLE) of the parameter θ .



Sufficient statistics are a key concept in statistical theory that help to simplify the process of parameter estimation. Here's a detailed explanation:

Definition of Sufficient Statistic

A statistic $T(X)$ is said to be **sufficient** for a parameter θ if it captures all the information about θ contained in the sample data X . In other words, once you know the value of the sufficient statistic, the data provides no additional information about the parameter θ .

Formally, a statistic $T(X)$ is sufficient for θ if the conditional probability distribution of the sample data X given $T(X)$ does not depend on θ .

Formal Definition: Factorization Theorem

A common way to establish that a statistic $T(X)$ is sufficient for θ is by using the **Factorization Theorem**. According to this theorem, a statistic $T(X)$ is sufficient for θ if the joint probability density function (PDF) or probability mass function (PMF) of the data X can be factorized into two parts:

$$f(x | \theta) = g(T(x) | \theta) \cdot h(x)$$

where:

- $g(T(x) | \theta)$ is a function of the statistic $T(x)$ and the parameter θ ,
- $h(x)$ is a function of the sample data x but does not depend on θ .

where $f_2(x_i | \theta_2)$ is the PDF or PMF for Family 2 with parameter θ_2 .

2. Comparing Likelihoods Directly

Comparing the likelihoods $L_{\theta_1}(X)$ and $L_{\theta_2}(X)$ directly can be done to assess how well each family fits the data, but this comparison can be challenging if the distributions are not nested or if they have different parameter spaces.

- **Likelihood Ratio:** If the two families are not nested or have different parameterizations, the likelihood ratio test may not be directly applicable. However, you can still compare the likelihood values of the best-fitting models from each family:

$$\text{Likelihood Ratio} = \frac{L_{\hat{\theta}_1}(X)}{L_{\hat{\theta}_2}(X)}$$

where $\hat{\theta}_1$ and $\hat{\theta}_2$ are the MLEs of θ_1 and θ_2 , respectively.

- **Normalized Comparison:** Normalize the likelihoods to compare their relative fit:

$$\text{Relative Likelihood} = \frac{L_{\hat{\theta}_1}(X)}{L_{\hat{\theta}_1}(X) + L_{\hat{\theta}_2}(X)}$$

This gives you a sense of how each model fits the data relative to the other.

5. Practical Considerations

- **Model Complexity:** Consider the complexity of each model and whether the improvement in fit justifies the additional complexity.
- **Domain Knowledge:** Use domain-specific knowledge to guide the choice of distribution family.

Bayesian Estimation

Bayesian Estimation

1. **Bayesian Estimate:** In Bayesian estimation, you incorporate prior beliefs about the parameters into the estimation process. The parameter θ is treated as a random variable with a prior distribution $p(\theta)$. The goal is to find the posterior distribution of θ given the observed data.

2. **Mathematical Formulation:** Using Bayes' Theorem, the posterior distribution $p(\theta | X)$ is:

$$p(\theta | X) = \frac{p(X | \theta) \cdot p(\theta)}{p(X)}$$

where:

- $p(X | \theta)$ is the likelihood function.
- $p(\theta)$ is the prior distribution of θ .
- $p(X)$ is the marginal likelihood or evidence, which ensures the posterior is a valid probability distribution.

3. **Bayesian Point Estimates:** Bayesian estimation often involves calculating the expected value of the parameter with respect to the posterior distribution. Common Bayesian point estimates include:

- **Posterior Mean:** The expected value of θ given the data:

$$\hat{\theta}_{\text{Bayes}} = \mathbb{E}[\theta | X] = \int_{\Theta} \theta \cdot p(\theta | X) d\theta$$

- **Posterior Median:** The value of θ such that the posterior distribution has 50% probability on either side.
- **Posterior Mode:** The value of θ that maximizes the posterior distribution $p(\theta | X)$.



A confidence interval provides a range of values within which we expect the true population parameter (such as the mean) to lie with a certain level of confidence. Here's how to construct and interpret a confidence interval for the mean, given a sample and sample mean:

Steps to Construct a Confidence Interval for the Mean

1. Calculate the Sample Mean (\bar{x})

- This is the average of the sample data. It's a point estimate of the population mean.

2. Determine the Sample Standard Deviation (s)

- This measures the variability of the sample data. If you have the sample standard deviation, use it. Otherwise, if you have the sample variance (s^2), take the square root to get s .

3. Choose the Confidence Level ($1 - \alpha$)

- Common confidence levels are 95% or 99%. This confidence level determines the probability that the confidence interval contains the true population mean.

4. Find the Critical Value

- For a known population standard deviation:** Use the Z-distribution (standard normal distribution). The critical value is denoted as $Z_{\alpha/2}$.
- For an unknown population standard deviation:** Use the t-distribution. The critical value is denoted as $t_{\alpha/2, df}$, where df (degree of freedom) is $n - 1$, and n is the sample size.

5. Compute the Margin of Error (ME)

- For large samples (or known population standard deviation):**

$$ME = Z_{\alpha/2} \times \frac{s}{\sqrt{n}}$$

- For small samples (unknown population standard deviation):**

$$ME = t_{\alpha/2, df} \times \frac{s}{\sqrt{n}}$$

6. Construct the Confidence Interval

- The confidence interval is given by:

$$\bar{x} \pm ME$$

- So the interval is:

$$[\bar{x} - ME, \bar{x} + ME]$$

- Determine the p-value:**

Use the calculated test statistic and the degrees of freedom to find the p-value from the t-distribution. This p-value tells you the probability of observing a test statistic as extreme as, or more extreme than, the one calculated, assuming the null hypothesis is true.

- Compare the p-value to the Significance Level (α):**

- Significance Level (α):** Commonly set at 0.05 (5%).

- Decision Rule:**

- If $p \leq \alpha$, reject the null hypothesis H_0 . This indicates that there is a significant difference between the means of the two groups.
- If $p > \alpha$, do not reject H_0 . This indicates that there is not enough evidence to conclude that the means are different.

LINEAR REGRESSION



Linear regression is a powerful and commonly used model due to its simplicity and flexibility. Its utility includes:

- **Simplest to code:** The implementation of linear regression is straightforward in most programming languages. Libraries such as Python's `scikit-learn` offer built-in functions to perform linear regression in just a few lines of code.
- **Simplest to understand mathematically:** Linear regression is based on simple concepts like finding the line of best fit to minimize the difference between predicted and actual values (using the least squares method). This makes it an excellent model for teaching and explaining fundamental statistical concepts.
- **Nonlinear feature extraction can be used with linear solutions:** Despite its linear nature, linear regression can still capture nonlinear relationships in the data by using polynomial or interaction terms as additional features. This allows for increased model flexibility without losing the simplicity of a linear model.

$$t = Xw + e$$

Where:

- X is the matrix of input features,
- w is the vector of weights (parameters),
- $e \sim \mathcal{N}(0, \sigma^2)$, i.e., the noise is normally distributed with mean 0 and variance σ^2 .

Alternatively, using precision $\beta = \frac{1}{\sigma^2}$, we can rewrite the noise distribution as $e \sim \mathcal{N}(0, \frac{1}{\beta})$.

Thus, the conditional probability distribution $p(t|X, w, \beta)$ is the **likelihood of observing t** , assuming the model parameters and input data are known.

Gaussian Likelihood:

Since the noise is Gaussian, the targets t are also Gaussian around the predicted values Xw :

$$p(t|X, w, \beta) = \mathcal{N}(t|Xw, \frac{1}{\beta})$$

- Likelihood in Linear Regression

Gaussian Likelihood:

Since the noise is Gaussian, the targets t are also Gaussian around the predicted values Xw :

$$p(t|X, w, \beta) = \mathcal{N}(t|Xw, \frac{1}{\beta})$$

This means that t follows a normal distribution with:

- Mean Xw (the predictions from the model),
- Variance $\frac{1}{\beta}$, or precision β .

In more detail, the likelihood function is:

$$p(t|X, w, \beta) = \prod_{i=1}^m \mathcal{N}(t_i|X_iw, \frac{1}{\beta})$$

Where each target t_i is drawn from a normal distribution centered at X_iw (the predicted value for input X_i) and with variance $\frac{1}{\beta}$.

Expanded Form:

The probability density function for a Gaussian distribution $\mathcal{N}(t_i|X_iw, \frac{1}{\beta})$ is:

$$p(t_i|X_i, w, \beta) = \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2}(t_i - X_iw)^2\right)$$

For the entire dataset, the likelihood is the product of these individual probabilities:

$$p(t|X, w, \beta) = \left(\frac{\sqrt{\beta}}{\sqrt{2\pi}}\right)^m \exp\left(-\frac{\beta}{2} \sum_{i=1}^m (t_i - X_iw)^2\right)$$

Where m is the number of data points.

Log-Likelihood:

The log-likelihood is often used to simplify optimization. Taking the logarithm of $p(t|X, w, \beta)$ gives:

$$\log p(t|X, w, \beta) = \frac{m}{2} \log \beta - \frac{m}{2} \log 2\pi - \frac{\beta}{2} \sum_{i=1}^m (t_i - X_iw)^2$$

This log-likelihood is maximized to find the optimal values of w and β .

To maximize the likelihood in the context of a linear regression model with Gaussian noise, we can use the maximum likelihood estimation (MLE) approach. This involves finding the values of the parameters w (the weights) that maximize the likelihood function $p(t|X, w, \beta)$.

2. Log-Likelihood:

To simplify the maximization, we take the **logarithm** of the likelihood function, which gives the **log-likelihood**:

$$\log p(t|X, w, \beta) = \frac{m}{2} \log \beta - \frac{m}{2} \log 2\pi - \frac{\beta}{2} \sum_{i=1}^m (t_i - X_i w)^2$$

Since the first two terms do not depend on w , we focus on maximizing the last term:

$$\log p(t|X, w) = -\frac{\beta}{2} \sum_{i=1}^m (t_i - X_i w)^2$$

This is equivalent to minimizing the **residual sum of squares (RSS)**:

$$\text{RSS}(w) = \sum_{i=1}^m (t_i - X_i w)^2$$

3. Set Gradients with Respect to w to Zero:

To maximize the log-likelihood (or equivalently, minimize the RSS), we take the gradient of the log-likelihood with respect to w , and set it equal to zero:

$$\frac{\partial}{\partial w} \left(-\frac{\beta}{2} \sum_{i=1}^m (t_i - X_i w)^2 \right) = 0$$

Expanding the gradient:

$$\frac{\partial}{\partial w} \left(\sum_{i=1}^m (t_i - X_i w)^2 \right) = -2X^\top(t - Xw)$$

Set this gradient equal to zero:

$$-2X^\top(t - Xw) = 0$$

Simplifying:

$$X^\top Xw = X^\top t$$

4. Solve for w :

The solution to this equation is the **normal equation** for linear regression:

$$w = (X^\top X)^{-1} X^\top t$$

This gives the value of w that maximizes the likelihood, which is the **least squares solution**.

BIAS - VARIANCE TRADEOFF MEAN SQUARED ERROR

1. Expected Loss (MSE)

The Mean Squared Error (MSE) measures the average squared difference between the predicted values \hat{y} and the true values y . Mathematically, the MSE for a regression model is given by:

$$\text{MSE} = \mathbb{E}[(\hat{y} - y)^2]$$

Where:

- \hat{y} is the predicted value from the model,
- y is the true value.

The MSE is an important metric because it quantifies the accuracy of the model by averaging the squared errors over all data points.

2. Bias-Variance Decomposition

The bias-variance decomposition helps explain the sources of error in a model's predictions. The total expected error (MSE) can be decomposed into three terms:

$$\text{MSE} = \underbrace{\mathbb{E}[\hat{y}] - y}_{\text{Bias}^2} + \underbrace{\mathbb{E}[(\hat{y} - \mathbb{E}[\hat{y}])^2]}_{\text{Variance}} + \underbrace{\sigma^2}_{\text{Irreducible Error}}$$

Where:

- **Bias:** Measures how far the average prediction $\mathbb{E}[\hat{y}]$ is from the true value y . It represents the error introduced by the model's assumptions. A high bias indicates that the model is too simplistic and underfits the data.

$$\text{Bias}^2 = (\mathbb{E}[\hat{y}] - y)^2$$

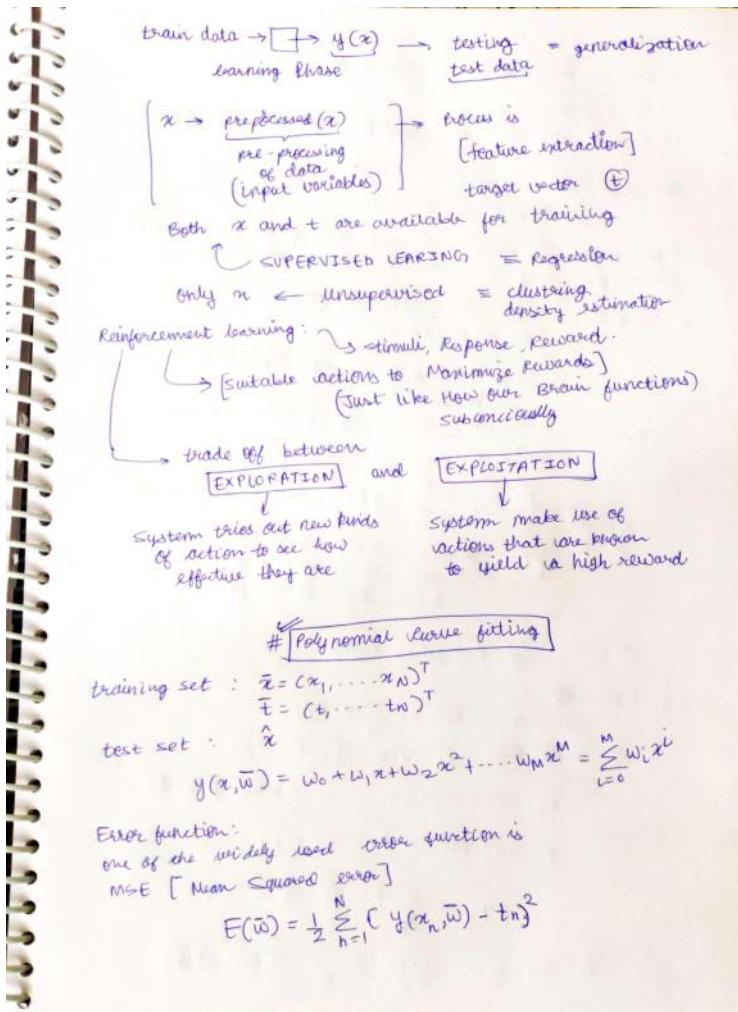
- **Variance:** Measures how much the predictions \hat{y} vary around their expected value $\mathbb{E}[\hat{y}]$. High variance indicates that the model is too sensitive to the specific training data and overfits the data.

$$\text{Variance} = \mathbb{E}[(\hat{y} - \mathbb{E}[\hat{y}])^2]$$

- **Irreducible Error:** Represents the inherent noise σ^2 in the data that no model can capture. This error is due to randomness or measurement errors and cannot be reduced by improving the model.

Bias-Variance Tradeoff

- A model with **high bias** (e.g., a simple model) may underfit the data, meaning it misses important patterns and relationships, resulting in higher errors.
- A model with **high variance** (e.g., a complex model) may overfit the training data, meaning it captures noise or fluctuations specific to the training set, which increases error when making predictions on unseen data.
- The goal in model selection is to balance bias and variance to minimize the total error (MSE).



concept
we choose value of \bar{w} such that it minimizes the $E(\bar{w})$

but there is one problem \rightarrow choosing M .

\hookrightarrow order [Model selection]

\hookrightarrow there is also case of overfitting ; when $E(w^*) \neq 0$

\Rightarrow Root Mean Squared error

$$E_{\text{rms}} = \sqrt{\frac{1}{N} E(w^*)}$$

\hookrightarrow More flexible polynomial with large value of M are becoming increasingly fitted to random noise or target vectors

\hookrightarrow One of the technique to control overfitting

[REGULARIZATION]

\hookrightarrow Involves adding a penalty term to the error function

simplest such penalty is sum of squares of all coeff.

$$\|w\|^2 = w^T w = w_0^2 + w_1^2 + \dots + w_n^2$$

$$\text{new error term } \hat{E}(\bar{w}) = \frac{1}{N} \sum_{i=0}^N (y_i - \bar{w}_0 - \bar{w}_1 x_i)^2 + \frac{\lambda}{2} \|w\|^2$$

$[\lambda \Rightarrow$ governs relative importance]

\hookrightarrow Controls the complexity of model

determine degree of overfitting

Probability theory

Rules of Probability

$$\text{Sum Rule } P(X) = \sum Y P(X|Y)$$

$$\text{product Rule } P(X,Y) = P(Y|X) \cdot P(X)$$

$$P(X,Y) = P(Y,X) \quad \text{symmetry property}$$

$$P(Y|X) \cdot P(X) = P(X|Y) \cdot P(Y)$$

$$P(Y) = \sum X P(Y|X) \cdot P(X)$$

$$\Rightarrow P(X,Y) = P(X) \cdot P(Y) \quad \xrightarrow[X \perp Y]{} \quad P(X|Y) = P(X)$$

Probability density

$$P(x \in (a, b)) = \int_a^b p(x) dx$$

$p(x) \geq 0$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

Change of variable: $x = g(y)$
 $f(x) = f(g(y)) = f'(g(y))$

$$p_y(y) = p_x(g(y)) |g'(y)|$$

$$\rightarrow P(z) = \int_{-\infty}^z p(x) dx$$

[CDF]

$$\rightarrow p'(x) = p(x)$$

Joint PDF $\rightarrow \bar{x} = (x_1, x_2, \dots, x_n)$
 $p(\bar{x}) = p(x_1, x_2, \dots, x_n)$
 $p(\bar{x}) \geq 0$
 $\int_{-\infty}^{\infty} p(\bar{x}) d\bar{x} = 1$

for continuous PDFs: $p(x) = \int_y p(x, y) dy$
 $p(x, y) = p(y|x) \cdot p(x)$

EXPECTATION

expectation of some $f(x)$ is average value of $f(x)$ denote by

$$E(f) = \sum_x p(x) f(x)$$

$$E(s) = \int_x p(x) f(x) dx$$

$$E_x(f(x, y)) = \int_x p(x) \cdot f(x, y) dx = g(y)$$

$$E_x(s|y) = \int_x p(x|y) \cdot f(x) dx$$

VARIANCE

$$\text{Var}[s] = E[(E(f(x)) - s)^2]$$

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

COVARIANCE

$$\begin{aligned}\text{cov}[x,y] &= \mathbb{E}_{x,y}[(x - \mathbb{E}(x))(y - \mathbb{E}(y))] \\ &= \mathbb{E}_{x,y}[xy] - \mathbb{E}(x)\mathbb{E}(y) \\ \text{cov}[\bar{x},\bar{y}] &= \mathbb{E}_{\bar{x},\bar{y}}[\bar{x}\bar{y}^T] - \mathbb{E}[\bar{x}]\mathbb{E}[\bar{y}^T]\end{aligned}$$

$$\rightarrow \text{cov}[x,x] \triangleq \text{var}[x]$$

Bayesian Probability

$$P(\bar{\omega} | D) = \frac{P(D|\bar{\omega}) \cdot P(\bar{\omega})}{P(D)}$$

likelihood function

$D \rightarrow$ some observed data
 $D = \{t_1, t_2, \dots, t_n\}$

[posterior \propto likelihood \times prior]
 Bayes' theorem

$$P(D) = \int P(D|\bar{\omega}) \cdot P(\bar{\omega}) d\bar{\omega}$$

17/9/2024

Gaussian distribution

$$\begin{aligned}N(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\ \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) dx &= 1 \rightarrow \text{normalized}\end{aligned}$$

$$\begin{aligned}E[x] &= \mu \\ E[x^2] &= \mu^2 + \sigma^2 \\ \text{var}(x) &= \sigma^2\end{aligned}$$

D-dimensional:

$$N(\vec{x}|\vec{\mu}, \vec{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\vec{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} (\vec{x} - \vec{\mu})^T \vec{\Sigma}^{-1} (\vec{x} - \vec{\mu})\right)$$

 $\vec{\mu} \rightarrow$ mean vector $\vec{\Sigma} \rightarrow$ covariance matrix

i.i.d (independent & identically distributed)

$$\underset{\text{likelihood}}{\rightarrow} P(\vec{x} | \vec{t}, \sigma^2) = \prod_{i=1}^N N(x_i | t_i, \sigma^2) \quad \text{given } \vec{t}, \sigma^2$$

$$\log \text{likelihood} \ln P(\vec{x} | \vec{t}, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - t_i)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

→ Maximise this w.r.t to \vec{w} , σ^2 we get

Minimum likelihood solution

curve fitting

$$\vec{x}, \vec{t} \rightarrow \text{assume gaussian distribution with mean } f(\vec{x}; \vec{w})$$

$$\text{input} \quad \text{target} \quad P(\vec{t} | \vec{x}, \vec{w}, \beta) = N(t | f(\vec{x}; \vec{w}), \beta^{-1})$$

$\Rightarrow \beta = \text{precision parameter}$
(inverse of variance)
of distribution

Step ① Now we use training data (\vec{x}, \vec{t}) to determine values of \vec{w} and β by Maximum LIKELIHOOD.

$$P(\vec{t} | \vec{x}, \vec{w}, \beta) = \prod_{i=1}^N N(t_i | f(\vec{x}_i; \vec{w}), \beta^{-1})$$

Step ② Maximising log likelihood is equivalent to Minimising $\sum_{i=1}^N \text{quadratic}$

BAYESIAN APPROACH

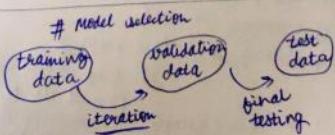
→ introduce polynomial distribution over \vec{w}

$$P(\vec{w} | \vec{x}) = N(\vec{w} | \vec{\theta}, \vec{\Sigma}) = \left(\frac{\alpha}{2\pi}\right)^{\frac{m}{2}} \exp\left(-\frac{\alpha}{2} \vec{w}^\top \vec{\Sigma} \vec{w}\right)$$

[$\alpha \rightarrow \text{Hyperparameters}$] field of Hyperparameter tuning

given: training data \vec{x} and \vec{t}
new test point x] → ignore

goal: predict the value of t
 $P(t | x, \vec{x}, \vec{w}) = \int P(t | \vec{x}, \vec{w}) \cdot P(\vec{w} | \vec{x}, \vec{t}) d\vec{w}$



⇒ Maximizing misclassification rate [DECISION THEORY]

→ divide input space \hat{x} onto K regions R_k are assigned to class c_k . [decision boundaries]

$$P(\text{mistake}) = P(x \in R_1, c_2) + P(x \in R_2, c_1)$$

$$P(\text{correct}) = \sum_{k=1}^K P(x \in R_k | c_k) d\omega$$

⇒ Minimizing expected loss

↳ Mistake 1 → Mistake 2 ↳ healthy person
weight cancer ↓

↳ loss function → minimize this
cost function → minimize this
utility function → minimize this

⇒ rejection option

↳ introducing θ_k and rejecting those input x for which largest of posterior prob. $P(c_k|x)$ is less than θ_k .

↳ extend the rejection criteria to minimize loss.

⇒ Inference and decision classification

Inference stage → Decision stage

use training data to learn a model $P(c_k|x)$ → use posterior probabilities to make optimal class assignment

→ alternative approach

↳ simply learn a function that map x to decision directly [DISCRIMINANT FUNCTIONS]

(a) Inference, decision → posterior probability

(b) Discriminant Model

Better → ① Minimizing Risk → Loss Matrix

② Reject options →

③ Compensating for class prior

(c) Model Combinations

x-ray image (x_1) blood sample (x_B) Need of balanced data set (equal no. of data sets corresponding to each class)

$P(x_1, x_B | c_k) = P(x_1 | c_k) P(x_B | c_k)$ (conditional independence)

Loss functions for Regression
consider Basic Squared loss:

$$L(t, y(x)) = \frac{1}{2} (y(x) - t)^2$$

Average loss $E[L(t, y(x))] = \int \int \frac{1}{2} (y(x) - t)^2 p(x, t) dx dt$
 $\frac{\partial E[L]}{\partial y(x)} = 2 \int (y(x) - t) p(x, t) dt = 0$

$M(x) = \int t \cdot p(x, t) dt = \int t \cdot p(t/x) = E[t/x]$
Minkowski loss \hookrightarrow generalized version of squared loss. \hookrightarrow Regression function

* $L(t, y(x)) = \frac{1}{2} (y(x) - t)^2$

[INFORMATION THEORY]

$p(x) \rightarrow$ probability distribution
 $h(x) \rightarrow$ quantity \in measure of $p(x)$
 \Rightarrow we have 2 events x and y , and they are unrelated
then information gained from both of them is just sum
& thus $h(x,y) = h(x) + h(y)$
similarly $p(x,y) = p(x)p(y)$ (Independent)

$$H(x) = -\log_2 p(x)$$

\rightarrow average amount of information transmitted is $(w \cdot H(x))$
 $H(x) = -\sum_x p(x) \log_2 p(x)$ unit \rightarrow bits
 $\frac{\text{expectation of } h(x) \text{ wrt } p(x)}{p(x)}$ called ENTROPY of random var

\hookrightarrow Huffman Encoding

\rightarrow Entropy is a lower bound on no. of bits to transmit state of Random Variables.

\rightarrow in case of natural log: units \rightarrow nat
 \hookrightarrow choosing objects from N objects in order; $N!$ ways
In i th bin; ordering objects N_i !
 $W = \frac{N!}{\prod N_i!} \rightarrow$ Multiplicity

entropy is defined as log of multiplicity scaled to some const.

$$H = \frac{1}{N} \log W = \frac{1}{N} \ln N! - \frac{1}{N} \sum_i \ln n_i!$$

Hill Sterling approx. $\log N! \approx N \log N - N$ as $N \rightarrow \infty$

$$H = -\lim_{N \rightarrow \infty} \sum_i \left(\frac{n_i}{N} \right) \log \left(\frac{n_i}{N} \right) = -\sum_i p_i \ln p_i$$

specific arrangement of object in bins \rightarrow MICROSTATE
overall distribution of N objects \rightarrow MACROSTATE
 $\rightarrow W$ is also weight of Macrostate.

$H[P] = -\sum_i p(x_i) \log p(x_i) \quad P(X=x_i) = p(x_i)$

$H[X] = \int p(x) \ln p(x) dx$ continuous

\rightarrow In discrete distro \rightarrow Max. Entropy [Uniform distro]
 \rightarrow continuous \rightarrow [Gaussian distro]

$\int_{-\infty}^{\infty} p(x) dx = 1$
 $\int_{-\infty}^{\infty} x p(x) dx = \mu$
Mean

$\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2$
Variance

$H[y|x] = -\iint p(y|x) \ln(p(y|x)) dy dx$

conditional entropy $H[x,y] = H[y|x] + H[x]$

$\rightarrow f(x) \rightarrow$ CONVEX function
property that every chord lies on the function.

Jensen's Inequality $f(E(x)) \leq E[f(x)]$

$f(\int x p(x) dx) \leq \int f(x) p(x) dx$

Ch-2: Probability Distro

learn: model probability distribution $p(x)$ of random var. X
given finite set of $x \in \{x_1, x_2, \dots, x_N\}$ of observation

[DENSITY ESTIMATION]

① Bernoulli distro.

$$\text{Bern}(x|p) = p^x (1-p)^{1-x}$$

$$E(x) = p$$

$$\text{Var}(x) = p(1-p)$$

construction: $P(D|x) = \prod_{n=1}^N p(x_n|x) = \prod_{n=1}^N p^{x_n} (1-p)^{1-x_n}$

log likelihood function: $\ln p(D|x) = \sum_{n=1}^N x_n \ln p + (1-x_n) \ln (1-p)$

$\downarrow d/dp$

get Maximum likelihood estimator $\hat{x}_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$ same as sample mean

② Binomial distro.

$$\text{Bin}(m|N,p) = \binom{N}{m} p^m (1-p)^{N-m}$$

$$E[m] = Np, \text{Var}[m] = Np(1-p)$$

③ Beta distro.

$$\text{Beta}(x|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$$

$$\int \text{Beta}(x|a,b) dx = 1$$

gamma function: $\Gamma(\alpha) = \int_0^\infty x^{\alpha-1} e^{-x} dx$

for natural number: $\Gamma(n) = (n-1)!$

$$\Gamma(\alpha+1) = \alpha \Gamma(\alpha) \quad \Gamma(\frac{1}{2}) = \sqrt{\pi}$$

$$\rightarrow E[x] = \frac{a}{a+b}, \text{Var}[x] = \frac{ab}{(a+b)^2(a+b+1)}$$

$a, b \rightarrow$ hyperparameters

$$N(\bar{x} | \bar{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu}) \right\}$$

$\downarrow \rightarrow$ 0-dimensional Mean vector
 $\det \Sigma \rightarrow$ D-dimensional Covariance Matrix
 $\Delta^2 = (\bar{x} - \bar{\mu})^T \Sigma^{-1} (\bar{x} - \bar{\mu})$
 $\downarrow \rightarrow$ Mahalanobis distance
 $\Sigma \rightarrow$ symmetric (full rank)
 $\Sigma \mathbf{u}_i = \lambda_i \mathbf{u}_i$ (eigen vectors) $i=1 \rightarrow D$
 As Σ is symmetric and real
 $\mathbf{u}_i^T \mathbf{u}_j = I_{ij}$; $J_{ij} = \begin{cases} 1 & i=j \\ 0 & \text{otherwise} \end{cases}$
 Orthonormal set
 $\Sigma = \sum_{i=1}^D \lambda_i \mathbf{u}_i \mathbf{u}_i^T \rightarrow \Sigma^{-1} = \sum_{i=1}^D \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$
 $(\Delta^2 = \sum_{i=1}^D \frac{y_i^2}{\lambda_i}) \quad y_i = \mathbf{u}_i^T (\bar{x} - \bar{\mu})$
 New coordinate system defined by orthonormal vectors \mathbf{u}_i
 $y = (y_1, y_2, \dots, y_D)^T \rightarrow y = U(\bar{x} - \bar{\mu})$
 $U \rightarrow$ orthogonal matrix
 $U = [u_1, u_2, \dots, u_N]^T \quad UU^T = I$
 \rightarrow A matrix whose eigen values are strictly +ve, then matrix \rightarrow positive definite $[\lambda_i > 0] \forall i=0 \rightarrow D$
 $\lambda_i > 0 \rightarrow$ semi-positive definite
 $J_{ii} = \frac{\partial g_i}{\partial y_i} = u_{i,i} = U_{ii}^T$
 $| \Sigma | = \prod_{i=1}^D \lambda_i \rightarrow$ in y_i coordinate system $- \frac{y_i^2}{2\lambda_i}$
 $p(y) = p(x|I) = \prod_{i=1}^D \frac{1}{\sqrt{2\pi\lambda_i}} e^{-\frac{y_i^2}{2\lambda_i}}$
 independent, gaussian RV.
 univariate
 $E[x] = \bar{\mu}$
 $\text{Var}[x] \neq \Sigma$
 $E[xx^T] = \bar{\mu}\bar{\mu}^T + \Sigma$
 $\Rightarrow \text{Cov}[x] = E[(x - E[x])(x - E[x])^T] = \Sigma$

conditional gaussian dist.

$N(\bar{x}, \bar{\Sigma})$ \bar{x} partitioned into M components x_a & remaining $D-M \rightarrow x_b$

$$\bar{x} = \begin{pmatrix} x_a \\ x_b \end{pmatrix} \Rightarrow \bar{\Sigma} = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}$$

$$\Sigma^T = \Sigma \quad \Sigma_{ba}^T = \Sigma_{ab}$$

precision Matrix

$$\Lambda = \Sigma^{-1} \quad \Lambda = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$P(x) = r(x_a, x_b)$, let say x_b is observed value
Joint distribution
Now find prob. distn. Since x_a

coefficient in exponent is

$$(x - \mu)^T \Sigma (x - \mu) = (\mu_a - x_a)^T \Lambda_{aa} (\mu_a - x_a) + (x_a - \mu_a)^T \Lambda_{ab} (x_a - \mu_a) + (x_b - \mu_b)^T \Lambda_{ba} (x_a - \mu_a) + (x_b - \mu_b)^T \Lambda_{bb} (x_b - \mu_b).$$

$\checkmark P(x_a | x_b) \rightarrow$ gaussian

Now find its Mean & Variance $(\mu_{a|b}, \Sigma_{a|b})$

conditional prob.

$P(x_a | x_b)$

$\mathcal{N}(\bar{x}_a | \mu_{a|b}, \Sigma_{a|b}) \rightarrow$ gaussian dist.

$X = \begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMB^{-1} \end{pmatrix}$

$M = (A - BD^{-1}C)^{-1}$

$M^{-1} = A - BD^{-1}C$

Shur complement of matrix X

$$(\Sigma_{aa} \Sigma_{ab})^{-1} = \begin{pmatrix} \Lambda_{aa} & \Lambda_{ab} \\ \Lambda_{ba} & \Lambda_{bb} \end{pmatrix}$$

$$(\Sigma_{ba} \Sigma_{bb})^{-1} = \begin{pmatrix} \Lambda_{ba} & \Lambda_{bb} \\ \Lambda_{ab} & \Lambda_{aa} \end{pmatrix}$$

$$\Lambda_{aa} = (\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1}$$

$$\Lambda_{ab} = -(\Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba})^{-1} \Sigma_{ab} \Sigma_{bb}^{-1}$$

put these

$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (x_b - \mu_b)$

$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}$

\rightarrow mean of conditional dist $P(x_a | x_b)$ is linear function of x_b
and covariance Matrix is independent of x_a .

\checkmark [Linear Gaussian Model]

Marginal Gaussian distro.

$$p(x_a) = \int p(x_a, x_b) dx_b$$

↓

$E[x_a] = \mu_a$ Gaussian Marginal distribution
 $\text{cov}[x_a] = \Sigma_{aa}$ Gaussian Conditional distribution

\rightarrow linear gaussian model

$$\hookrightarrow y = Ax + b \quad p(x) = N(x|z, \Lambda^{-1})$$

$$p(y|x) = N(y|Ax+b, L^{-1})$$

$Z = \begin{pmatrix} x \\ y \end{pmatrix} \rightarrow E[Z] = \begin{pmatrix} \mu_x \\ A\mu_x + b \end{pmatrix} \quad E[y] = A\mu_x + b$

$\text{cov}[y] = L^{-1} + A\Lambda^{-1}A^T$

1813/24

Maximum likelihood for gaussian

$$\ln p(\vec{x}|z, \Sigma) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln|\Sigma| - \frac{1}{2} \sum_{n=1}^N (\vec{x}_n - z)^T \Sigma^{-1} (\vec{x}_n - z)$$

$\vec{x} = (x_1, x_2, \dots, x_N)^T$

each obs x_n assumed to be drawn independently from gaussian distribution.

$\sum_{n=1}^N \vec{x}_n$ depends on 2 terms
 $\sum_{n=1}^N \vec{x}_n \vec{x}_n^T$ sufficient statistics for gaussian distro.

$\frac{\partial \ln p(\vec{x}|z, \Sigma)}{\partial \Sigma} = \sum_{n=1}^N (\vec{x}_n - z)^T = 0$

\rightarrow Maximum likelihood estimate of mean

$$\hat{z}_{MLE} = \frac{1}{N} \sum_{n=1}^N \vec{x}_n$$

$$\Sigma_{MLE} = \frac{1}{N} \sum_{n=1}^N (z_n - \hat{z}_{MLE})(z_n - \hat{z}_{MLE})^T$$

Bayesian Inference

$$p(z|X) \propto p(X|z) \cdot p(z)$$

Student's t distribution

we have univariate $N(x|z, \tau^{-1})$ together with prior gaussian $\text{Gam}(z|a, b)$

$$p(x|z, a, b) = \int_0^\infty N(x|z, \tau^{-1}) \cdot \text{Gam}(z|a, b) dz$$

New parameter $v = za$ $\nu = a/b = St(x|z, \lambda v)$

$\text{st}(x|\mu, \lambda, v)$
 degrees of freedom
 for $v=1$ st \rightarrow Cauchy distro
 for $v \rightarrow \infty$ st \rightarrow Gaussian $N(x|\mu, \lambda^2)$

Periodic Var.
 $P(\theta) > 0$
 $2\pi \int_0^{2\pi} P(\theta) d\theta = 1$
 $P(\theta + 2\pi) = P(\theta)$

$\vec{x} = (\vec{x} \cos \theta, \vec{x} \sin \theta)$
 \Rightarrow Mixtures of Gaussians
 Superposition of K Gaussian densities
 $P(\vec{x}) = \sum_{k=1}^K \pi_k N(\vec{x}|\mu_k, \Sigma_k)$
 $(\sum_{k=1}^K \pi_k = 1)$
 $P(\vec{x}) \geq 0$
 $\pi_k > 0$ # Exponential family
 $p(\vec{x}|\vec{\eta}) = h(\vec{x}) g(\vec{\eta}) \exp(\vec{\eta}^\top \vec{\mu}(\vec{x}))$
 $g(\vec{\eta}) \int h(\vec{x}) \exp(\vec{\eta}^\top \vec{\mu}(\vec{x})) d\vec{x} = 1 \rightarrow \vec{\eta} \rightarrow \text{parameters}$
 \downarrow Normalization factor

① Bernoulli: $p(x|d) = d^x (1-d)^{1-x}$
 $p(x|d) = \exp(x \ln d + (1-x) \ln(1-d))$
 $= (1-d) \exp \left[\ln \left(\frac{d}{1-d} \right) x \right] \rightarrow \eta = \ln \left(\frac{d}{1-d} \right)$
 $y = \sigma(\eta)$ \rightarrow Logistic Sigmoid function
 $\eta = \frac{t}{1 + \exp(-t)}$
 $p(x|\eta) = \sigma(-\eta) \exp(\eta x)$
 $1 - \sigma(\eta) = \sigma(-\eta)$
 $\eta(\eta) = \eta$
 $h(\eta) = 1$
 $g(\eta) = \sigma(-\eta)$ \rightarrow for Bernoulli

② Consider Multinomial distro.
 $p(x|\vec{\eta}) = \prod_{k=1}^M \eta_k^{x_k} \exp \left(\sum_{k=1}^M x_k \ln \eta_k \right) \rightarrow p(x|\vec{\eta}) = \exp(\vec{\eta}^\top \vec{x})$
 $\sum_{k=1}^M \eta_k = 1$
 $\eta_k = \ln t_k$
 $t = e^{\eta_1, \eta_2, \dots, \eta_M}$
 $\ln \left(\frac{\eta_k}{1 - \sum_j \eta_j} \right) = \eta_k$
 $\eta_k = \sigma(\eta_k) \rightarrow \eta_k = \frac{e^{\eta_k}}{1 + \sum_j e^{\eta_j}}$ # Softmax function

Maximum likelihood & sufficient statistics

$$P(\hat{\theta}^* | \mathbf{x}) = h(\mathbf{x}) g(\mathbf{x}) \exp(\mathbf{x}^T \hat{\theta}^*(\mathbf{x}))$$

$$\nabla g(\mathbf{x}) \int h(\mathbf{x}) e^{\mathbf{x}^T \hat{\theta}^*(\mathbf{x})} d\mathbf{x}$$

$$g(\mathbf{x}) \int h(\mathbf{x}) e^{\mathbf{x}^T \hat{\theta}^*(\mathbf{x})} \cdot u(\mathbf{x}) d\mathbf{x} = 0$$

$$\boxed{-\nabla \ln g(\mathbf{x}) = \mathbb{E}[\hat{\theta}^*(\mathbf{x})]}$$

$\sum_n \hat{\theta}^*(x_n)$
sufficient stats

$$-\nabla \ln g(\hat{\theta}_{MLE}) = \frac{1}{N} \sum_{n=1}^N \hat{\theta}^*(x_n)$$

Non parametric methods

Parametric \rightarrow might give quite poor model result in poor prediction

Histogram density models

\hookrightarrow partition \mathbf{x} into distinct bins with width Δ & then count

number n_i of obs. \mathbf{x} falling into bin i .

(number too small)
 \hookrightarrow $p_i = \frac{n_i}{N\Delta}$ for normalization
 \hookrightarrow often bins chosen are of same width $\Delta_i = \Delta$.

$\Delta \rightarrow$ very small [spiky too]
 $\Delta \rightarrow$ very large [too smooth] \Rightarrow this model also dependent of choice of edge location for bins.

Limitations of Histogram

\hookrightarrow scaling with dimensionality.
 \hookrightarrow divide each variable with 1-dim space into M bins,
the total no. of bins = $(M)^D$ \rightarrow exponential scaling

Kernel Density Estimators

\rightarrow set we have counted a data set comprising N observations drawn from $P(\mathbf{x})$.

\rightarrow each data point has prob. P of falling inside within R (region).

\rightarrow total number of point k that lie inside R . [distr. to binomial dist.]

$$\text{Bin}(k|N, P) = \frac{N!}{k!(N-k)!} P^k (1-P)^{N-k}$$

for large N ; $k = NP$
 \rightarrow Assume R is sufficiently small that prob. density $P(\mathbf{x})$ is roughly constant

$$P \approx p(\mathbf{x}) \cdot V_R$$
 volume of R .

$$\boxed{P(\mathbf{x}) = \frac{k}{NV}}$$

$p(x) = \frac{K}{NV}$
 fix K and determine value of V
K-nearest neighbour
 In order to count number K of points falling into region;
 $\hookrightarrow f^n: K(u) = \sum_{\substack{i=1 \\ \|x_i - u\| \leq h/2}}^N$
 kernel function Parzen windows
 $K((x-x_n)/h) \rightarrow$ we use x_n if x_n lies inside hypercube of side h
 centered on x .
 total number of data points lying inside this cube
 $K = \sum_{n=1}^N K\left(\frac{x-x_n}{h}\right)$
 $p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{x-x_n}{h}\right); V = h^D$
 If we choose smoother $\Rightarrow p(x) \rightarrow$ smoother
 kernel function \Rightarrow kernel function of side h in dimension
 consider $K(x) \sim$ Gaussian
 $p(x) = \frac{1}{N} \sum_{n=1}^N \frac{1}{(2\pi h^2)^{1/2}} \exp\left(-\frac{\|x-x_n\|^2}{2h^2}\right)$
 $h \rightarrow$ smoothing parameter Bandwidth
small h [sensitivity to noise]
large h [oversmoothing]
 In analogous to choice of bin in histogram
 # Nearest neighbour methods
 Problem in KDE: parameter h is fixed for all kernels
 In region of high density; large value of h may lead to oversmooth
 In region of low density; \hookrightarrow noisy estimate
 general result of location density estimation
 $p(x) = \frac{K}{NV}$
 we consider fixed value of K and find appropriate value of V using data.
 Conditions:
 $k(x_n) \geq 0$
 $\int k(x) dx = 1$

Process: we consider a small sphere centered on point \hat{x} , where
 we wish to estimate $f(\hat{x})$
 ↳ we grow radius of sphere until it contains precisely K data points.
 ↳ estimate of density $f(x)$ is given by setting V to volume of resulting sphere.

K Nearest neighbours

⇒ value of $k \approx$ degree of smoothing.
 ↳ for optimal $k \approx$ neither too large nor too small.

⇒ Solving: Problem of CLASSIFICATION.
 Let us suppose we have data set comprising of N points initial
 and N_k points in each class C_k

Given:

problem: classify new point x
 ↳ draw a sphere centered on x containing precisely K points
 ↳ suppose this sphere has volume V and contains K_k points from each class C_k

$$P(x|C_k) = \frac{K_k}{N_k V}$$

$$P(x) = \frac{K}{NV}$$

$$P(C_k) = \frac{N_k}{N}$$

using Bayes theorem: $P(C_k|x) = \frac{P(x|C_k)P(C_k)}{P(x)}$

→ to minimize probability of misclassification \Rightarrow
 this is done by assigning the test point x to class having minimum probability.

KNN & KDE
 ↳ require entire training set to be stored
 ↳ leading to expensive computation.

Direct models for regressions :

Supervised learning
goal: predict value of one or more continuous target var given
D dimensional x of input vector

linear regression \Rightarrow linear combination of basis functions.
(fixed set of some non-linear filters)

linear Basis function models \Rightarrow Minimize value of Loss functions
reported value of Loss functions

simplest linear Model

$$y(x, w) = w_0 + w_1 x_1 + \dots + w_D x_D$$

$$x = (x_1, x_2, \dots, x_D)^T$$

$$y(x, w) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(x)$$

allows for any fixed offset
in data \hookrightarrow [BIAS] parameter

additional dummy Basis fn $\phi_0(x) = 1$

$$\boxed{y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x)} = w^T \phi(x)$$

$$w = (w_0, w_1, \dots, w_M)^T \quad \phi = (\phi_0, \phi_1, \dots, \phi_{M-1})^T$$

for smooth pre-processing
or feature extraction

Similarly; polynomial regression:
consider single input var $x \Rightarrow \phi_j(x) = x^j$

other Basis functions: Gaussian BF $\phi_j(x) = e^{-\frac{(x-x_j)^2}{2s^2}}$ $s \Rightarrow$ spatial scale
 $x_j \Rightarrow$ location of BF

sig mooidal BF $\phi_j(x) = \sigma\left(\frac{x-x_j}{s}\right) \quad \sigma(a) = \frac{1}{1+e^{-a}}$

Another Basis: Fourier basis: expansion in sinusoidal fn:
 \hookrightarrow each represent specific frequency and has a spatial extent.

Maximum likelihood & least square

$$t = y(x, w) + \epsilon$$

\uparrow
deterministic
function

\hookrightarrow additive noise

$\epsilon \sim \text{Zero mean Gaussian RV with precision } \beta$
 $P(t|x, \omega, \beta) = N(t | y(x, \omega), \beta^{-1})$
 $E[t|x] = \int t P(t|x) dt = y(x, \omega)$
 Gaussian Noise assumption implies that conditional distribution of t given x is unimodal.
 Consider: $X = \{x_1, x_2, \dots, x_N\}$
 $\vec{t} = \{t_1, t_2, \dots, t_N\}$
 $P(\vec{t} | X, \omega, \beta) = \prod_{n=1}^N N(t_n | \omega^\top \phi(x_n), \beta^{-1})$
 MLF: $\ln P(\vec{t} | \omega, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \beta E_0(\omega)$
Sum of Squared Error Functions
 $E_0(\omega) = \frac{1}{2} \sum (t_n - \omega^\top \phi(x_n))^2$
 Gradient of MLF
 $\nabla \ln(P(\vec{t} | \omega, \beta)) = \sum_{n=1}^N 2t_n - \omega^\top \phi(x_n) \phi(x_n)^\top \beta = 0$
 $0 = \sum_{n=1}^N t_n \phi(x_n)^\top - \beta \sum_{n=1}^N \phi(x_n) \phi(x_n)^\top$
 $(\underline{\omega_{MLE}}) = (\phi^\top \phi)^{-1} \phi^\top t$
 $\phi \rightarrow N \times M$ (Design Matrix)
 $\phi = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_M(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_M(x_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_M(x_N) \end{pmatrix}$ Normal eqn for least square problems
 $\phi^\dagger = (\phi^\top \phi)^{-1} \phi^\top$ Pseudo Inverse of Matrix ϕ
 $\underbrace{A A^\top}_{N \times N} = A^\dagger$
 $A^\dagger A = I$
 $A^\dagger A \cdot A^\top = A^\top$
 $\boxed{A^\dagger = A^\top (A^\top A)^{-1}}$
 $A \cdot A^\dagger = I$
 $(A^\top A) A^\dagger = A^\top$
 $\boxed{A^\dagger = (A^\top A)^{-1} A^\top}$

$$A \cdot A^{-1} = I \quad (AB)^{-1} = B^{-1}A^{-1}$$

$$E_D(w) = \frac{1}{2} \sum_{n=1}^N \|t_n - w^T \sum_{j=1}^{M-1} w_j \phi_j(x_n)\|^2$$

$$\frac{\partial}{\partial w_0} \downarrow \quad \text{get } w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j$$

$$\bar{t} = \frac{1}{N} \sum_{n=1}^N t_n \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(x_n)$$

\Rightarrow Bias w_0 compensates for differences in average (over training set) of target values and weighted sum of basis functions.

$$\log P(t|w, \beta) \rightarrow \frac{\partial}{\partial \beta} \rightarrow \text{gives}$$

$$\text{Maximize log likelihood w.r.t. } \beta \quad P_{MLE} = \frac{1}{N} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

\rightarrow Inversion or noise precision is given by residual variance of target value around regression function.

Geometry of least squares

$$y = y(x_n, w)$$

Sum of squares.

Nth dimensional vector
whose n value is

\rightarrow Sum of least squares is then equal to squared Euclidean distance between y and t .

The least square solution w corresponds to that choice of y that lies in subspace S and that is closest to t .

Solution corresponds to orthogonal projections of t onto subspace S .

$y = \phi w_{ML}$ [If ϕ is singular, then adding regularization term is beneficial]

$$t = (t_1, t_2, \dots, t_N)$$

$$y = y(x_i, w)$$

$$y_i = y(x_i, w)$$

Sequential learning

\rightarrow If data set is very large (sequential algo./on-line algo.)

\rightarrow stochastic gradient descent

\hookrightarrow update parameter vector

$$E = \sum_n E_n \quad \hookrightarrow \text{no of iterations}$$

$$w^{(t+1)} = w^{(t)} - \eta \nabla E_t$$

$\eta \rightarrow$ learning rate parameter

$$w \rightarrow w^{(0)} \quad (\text{some starting vector})$$

for case of sum of squared error

$$w(t+1) = w(t) + \eta (t_n - w(t)^T \phi_{t_n}) \phi_{t_n}$$

$$\phi_{t_n} = \phi(x_n)$$

least mean squares or the LMS alg.

value of η \rightarrow algo converges

Regularized least squares

\hookrightarrow Idea of adding a regularization term to an error function in order to control (overfitting)

$$\text{final total error function} \rightarrow E_b(w) + \lambda E_w(w)$$

regularization coefficient

Simple form of regularizer \rightarrow sum of squares of weight vector element

relative importance of data dependent error $E_b(w)$ and regularization term $E_w(w)$

$$E_w(w) = \frac{1}{2} w^T w$$

If we consider sum of squared error function

$$E_b(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

total error: $E_b(w) + \lambda E_w(w)$

$$\frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 + \frac{\lambda}{2} w^T w$$

$\frac{d}{dw} = 0$ \hookrightarrow as it shrinks parameter values toward zero.

so get $w = (\lambda I + \phi^T \phi)^{-1} \phi^T t$

particular choice of regularizer
weight decay in squared learning alg.
parameter shrinkage

\Rightarrow more general regularizers

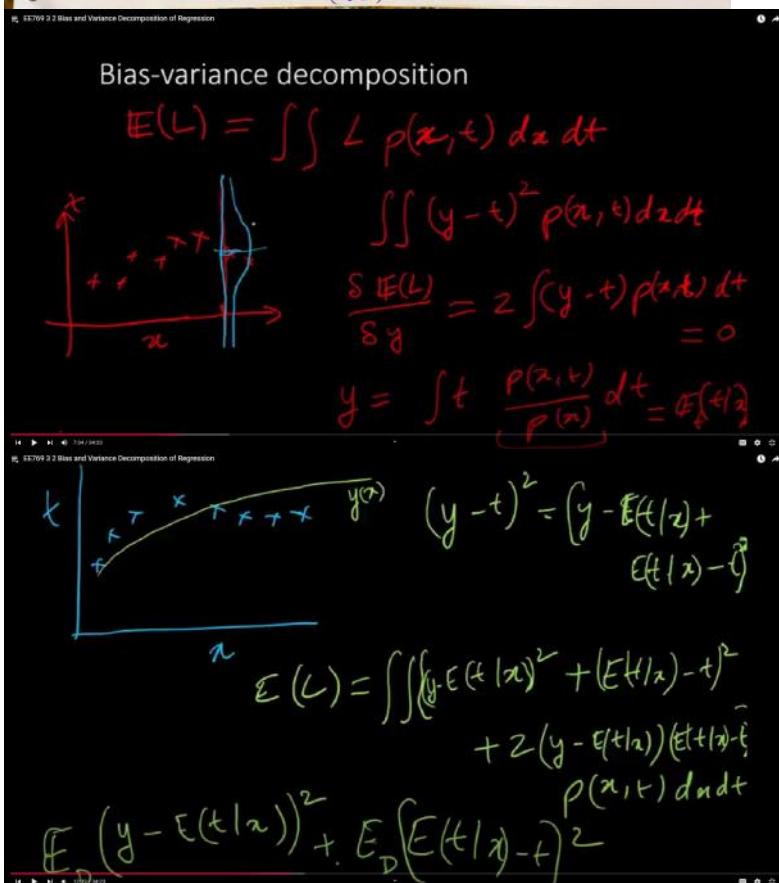
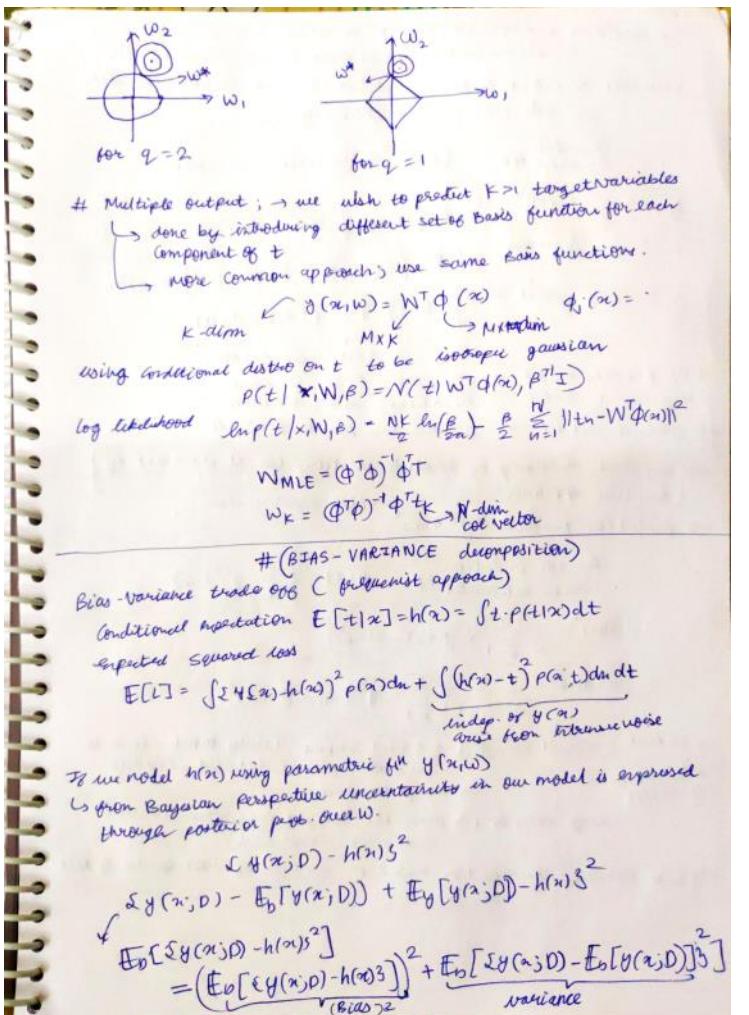
$$E_w(w) = \sum_{j=1}^M |w_j|^q$$

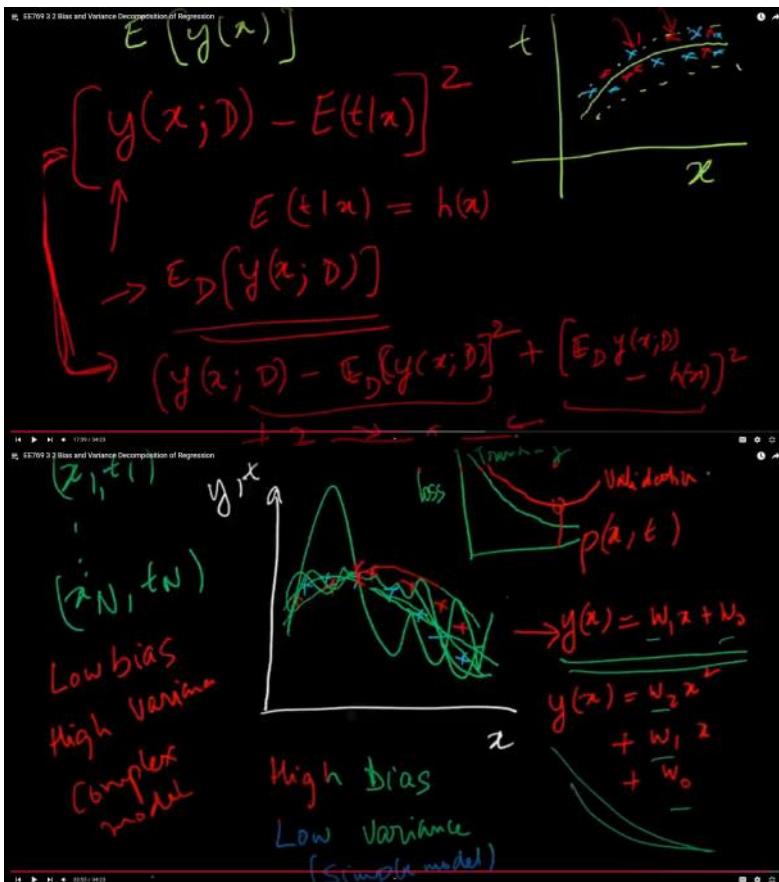
$q=1 \rightarrow$ lasso
 $q=2 \rightarrow$ quadratic

\Rightarrow as λ is increased, some of the coefficients are driven 0.

\Rightarrow regularization allows complex model to trained without severe overfitting.

Optimal Model
problem: finding appropriate number of basis functions.
determining suitable value of λ .





Squared Bias \rightarrow represent extent to which avg prediction over all data set differs from desired regression function

Variance \rightarrow extent to which solutions of each individual data sets vary around their avg.

$$\text{Expected loss} = (\text{bias})^2 + (\text{variance}) + (\text{noise})$$

$$(\text{bias})^2 = \int \mathbb{E}_\theta [y(x; D)] - h(x) \mathbb{P}(x) dx$$

$$(\text{variance}) = \int \mathbb{E}_\theta [\mathbb{E}[y(x; D)] - \mathbb{E}[y(x; D)^2]] \mathbb{P}(x) dx$$

$$\text{noise} = \int (h(x) - t)^2 p(x, t) dx dt$$

very flexible model \rightarrow low bias, high variance

very rigid model \rightarrow high bias, low variance

\Rightarrow optimal model \rightarrow Balance b/w bias & variance

\Rightarrow weighted averaging of multiple solution lie at the heart of

[Bayesian approach

\Rightarrow prediction function $\hat{y}(x)$ after regularization

$$\text{Average prediction can be estimated by } \bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x)$$

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L \hat{y}(x_n) - h(x_n) \mathbb{P}(x_n)$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \sum_{l=1}^L \hat{y}(x_n) - \bar{y}(x_n) \mathbb{P}(x_n)$$

\Rightarrow In book ; small value of $\lambda \rightarrow$ model become finely tuned to noise on each individual data set, leading to large variance

large value of $\lambda \rightarrow$ puts the weight parameters toward 0, leading to large bias.

\Rightarrow Bias-Variance decomposition based on averages of next sentence of data

Linear model for classification

↳ goal: we model: class conditional density $p(x|C_k)$ and class prob. $p(C_k)$

$$P(C_k|x) = \frac{P(x|C_k) \cdot P(C_k)}{\sum_j P(x|C_j) P(C_j)}$$

\leq shaped sigmoid function $= \frac{1}{1+e^{-a}} = \sigma(a)$

Symmetry property $\sigma(-a) = 1 - \sigma(a)$

$$a = \ln\left(\frac{P(x|C_1) \cdot P(C_1)}{P(x|C_2) \cdot P(C_2)}\right)$$

logit fn log odds

→ generalized linear model

$$P(C_k|x) = \frac{P(x|C_k) \cdot P(C_k)}{\sum_j P(x|C_j) P(C_j)}$$

Normalized exponential $= \frac{e^{a_k}}{\sum_j e^{a_j}}$

multiclass generalization of (Logistic Sigmoid)

where $a_k = \ln(P(x|C_k) \cdot P(C_k))$

→ also known as softmax function [smooth version of max function]
if $a_k \gg a_j$; then $P(C_k|x) = 1$ and $P(C_j|x) \approx 0$.

we assume class conditional densities are Gaussian assumption

$$P(x|C_k) = \frac{1}{(2\pi)^{D/2} |\Sigma|^1/2} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right)$$

consider case of 2 classes:

$$P(C_k|x) = \sigma(W^T x + w_0)$$

$$W = \Sigma^{-1} (\mu_1 - \mu_2)$$

$$w_0 = -\frac{1}{2} \mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2} \mu_2^T \Sigma^{-1} \mu_2 + \ln\left(\frac{P(C_1)}{P(C_2)}\right)$$

general case of K classes:

$$a_k(x) = W_k^T x + w_{k0}$$

assumption that all share a common covariance matrix

Linear functions and one optimal solution

$$y_i = W^T x_i + t_i$$

$$t = y + n \quad \text{E}(t \cdot y) \equiv \text{E}(n)$$

$$n \sim N(0, \sigma^2) \quad X = \begin{bmatrix} x_{11} & \dots & x_{1D} \\ \vdots & & \vdots \\ x_{N1} & \dots & x_{ND} \end{bmatrix}$$

$$p(t|x, W, \sigma^2)$$

$$\hat{y} = X^T W$$

Maximize $p(t|x) = \prod_{i=1}^N p(t_i|x_i)$

I. I. D.

$$= \prod_{i=1}^N \mathcal{N}(W^T x_i, \sigma^2)$$

$$= \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (t_i - W^T x_i)^2\right)$$

Maximizing $\ln(p(t|x)) \equiv \text{Maximize} -\sum_{i=1}^N (t_i - W^T x_i)^2$

$\equiv \text{Minimize} \sum_{i=1}^N (t_i - W^T x_i)^2$

$$W_{ML} = \underbrace{(X^T X)^{-1}}_{D \times D} \underbrace{X^T t}_{D \times N \quad N \times 1}$$

$N > D$

$$(X^T X)^{-1} X^T \quad \text{Pseudo inv. of } X$$

$$\underbrace{(X^T X)^{-1} X}_{D \times D} \quad O(D^3)$$

Gradient Descent $O(D^2)$

Iteratively $\underbrace{W \leftarrow W - \eta \nabla_w L}_{D \times 1} \quad \underbrace{W \leftarrow W + \eta \sum_{i=1}^N (t_i - \underbrace{W^T x_i}_{D \times 1}) x_i}_{D \times N}$

If each $P(x|c_k)$ has its own Σ_k , the model is quadratic discriminant.

consider 2 classes c_1, c_2

$P(c_1) = \pi_1, P(c_2) = 1 - \pi_1$

we have data set $(x_n, t_n), n = 1, \dots, N$; $t_n = 1 \rightarrow \text{Class 1}$
 $t_n = 0 \rightarrow \text{Class 2}$

$P(x_n, c_1) = P(c_1) \cdot P(x_n|c_1) \sim \pi_1 N(x_n|\mu_1, \Sigma_1)$
 $P(x_n, c_2) = P(c_2) \cdot P(x_n|c_2) \sim (1 - \pi_1) N(x_n|\mu_2, \Sigma_2)$

log likelihood

or $P(t|t, \pi, \mu, \Sigma) = \prod_{n=1}^N \left(\frac{1}{\sqrt{2\pi}} \sqrt{\pi} N(x_n|t, \Sigma) \right)^{t_n} \left(\frac{1}{\sqrt{2\pi}} \sqrt{1-\pi} N(x_n|\mu, \Sigma) \right)^{1-t_n}$

log likelihood that depends on π are

$\sum_{n=1}^N t_n \ln \pi + (1-t_n) \ln (1-\pi)$

get $\pi = \frac{1}{N} \sum_{n=1}^N t_n = \frac{N_1}{N} = \frac{N_1}{N_1 + N_2}$

generalization.

↳ maximization w.r.t. Σ , in log likelihood are terms depend on Σ , in log likelihood are

$\sum_{n=1}^N t_n \ln (N(x_n|\mu, \Sigma)) \rightarrow \text{find } \Sigma$

$\Sigma = \frac{1}{N_1} \sum_{n=1}^N t_n x_n$

similarly $\Sigma = \frac{1}{N_2} \sum_{n=1}^N (1-t_n) x_n$

Consider MLE for shared covariance Σ

$$\begin{aligned} & -\frac{1}{2} \sum_{n=1}^N t_n \ln(\pi/\Sigma) - \frac{1}{2} \sum_{n=1}^N t_n (x_n - \mu_1)^T \Sigma^{-1} (x_n - \mu_1) \\ & -\frac{1}{2} \sum_{n=1}^N (1-t_n) \ln((1-\pi)/\Sigma) - \frac{1}{2} \sum_{n=1}^N (1-t_n) (x_n - \mu_2)^T \Sigma^{-1} (x_n - \mu_2) \\ & = -\frac{N}{2} \ln(\pi/\Sigma) - \frac{N}{2} \text{tr}(\Sigma^{-1}) \end{aligned}$$

defn: $S = \frac{N_1}{N} S_1 + \frac{N_2}{N} S_2$

where $S_1 = \frac{1}{N_1} \sum_{n=1}^N (x_n - \mu_1)(x_n - \mu_1)^T$
 $S_2 = \frac{1}{N_2} \sum_{n=1}^N (x_n - \mu_2)(x_n - \mu_2)^T$

$\frac{d}{d\Sigma} = 0$

using we get $\Sigma = S$ for MLE

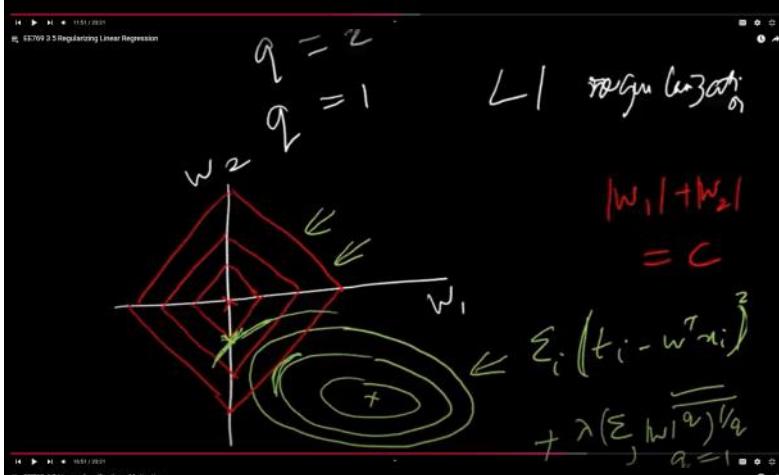
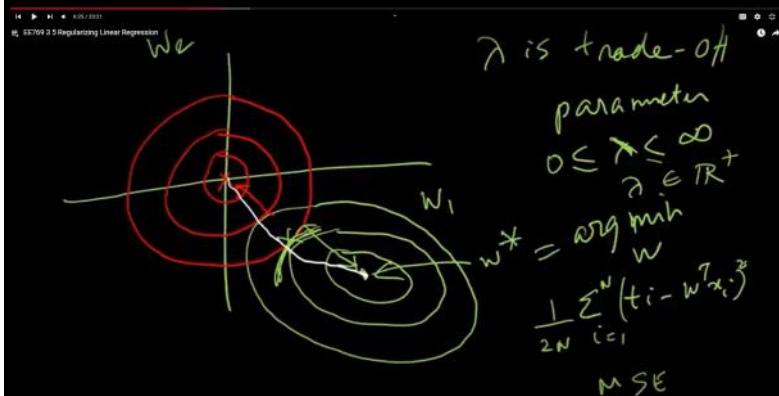
$$y = w_D x^D + \dots + w_1 x^1 + w_0 x^0$$

$$\sum_{j=1}^D w_j^2 \leq \eta$$

w_1

$\|w\|_2$ is penalized

Minimizing $\frac{1}{2N} \sum_{i=1}^N (t_i - w^T x_i)^2 + \frac{\lambda}{2} \sum_{j=1}^D |w_j|^q$ for $q=2$



Bayesian decision rule for classification

Cats Dogs $C = \{C_1, C_2\}$

$P(C) \leftarrow$

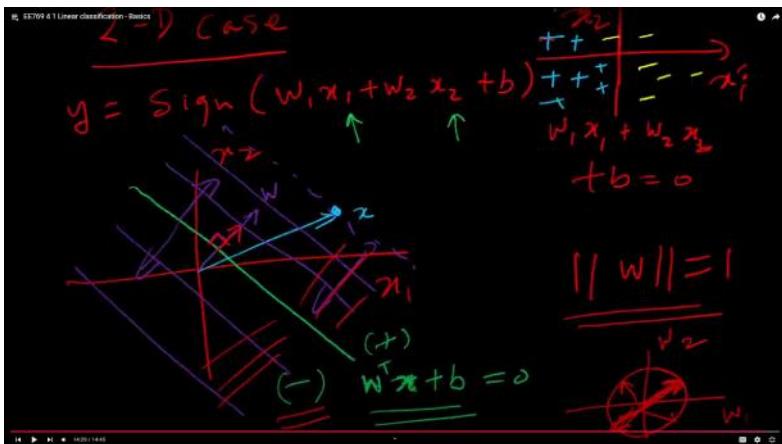
$P(x|C)$

$P(x|C_1) / P(x|C_2) = \frac{P(x|C_1) P(C_1)}{P(x|C_2) P(C_2)}$

$P(C_1|x) = \frac{P(x|C_1) P(C_1)}{P(x)}$

$P(C_2|x) = \frac{P(x|C_2) P(C_2)}{P(x)}$

likelihood \times prior
evidence



LOGISTIC REGRESSION

Circling back to Bayesian decisions

$$a = \log \frac{p(x|c_1)p(c_1)}{p(x|c_2)p(c_2)}$$

$$p(c_i|x) = \sigma(a) = \frac{1}{1 + \exp(-a)} \leftarrow$$

$$= \frac{p(x|c_i)p(c_i)}{\sum_j p(x|c_j)p(c_j)}$$

$$\max_w P(t|w) = \prod_{i=1}^N y_i^{t_i} (1-y_i)^{1-t_i}$$

$$y_i = p(t_i=1|x_i, w)$$

$$\min_w \left[-\sum_{i=1}^N t_i \log(y_i) + (1-t_i) \log(1-y_i) \right]$$

$$y_i = w^\top x_i + b \quad \text{Cross entropy}$$

Logistic regression as a convex problem

$$\nabla_w L = \sum_i (y_i - t_i) x_i$$

$$y_i = \sigma(w^\top x_i + b)$$

$$+ \lambda_1 \|w\|_2^2 \leftarrow \text{G.L}$$

$$+ \lambda_2 \|w\|_1 \leftarrow$$

Decision Rule:

Class x as c_k if $P(c_k|x) > P(c_j|x)$ for all $j \neq k$

$$P(x|c_k) \cdot P(c_k) > P(x|c_j) \cdot P(c_j) \quad \forall j \neq k$$

Intuition:

- Likelihood $P(x|c_k)$ reflects how probable the feature vector x is given class c_k .
- Prior $P(c_k)$ how probable is class c_k given no data.
- Posterior $P(c_k|x)$ update probability of class after observing x .

→ for binary classification

Class is 1 if $P(t=1|x) > P(t=0|x)$
 $P(x|t=1) \cdot P(t=1) > P(x|t=0) \cdot P(t=0)$

Gaussian Assumption with same Cov. Matrix Σ for all classes

$$P(x|c_k) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right)$$

$d \rightarrow$ dim of x

$\mu_k \rightarrow$ mean vector of class c_k

$\Sigma \rightarrow$ covariance matrix

→ Decision Rule

$$P(c_k|x) \propto P(x|c_k) \cdot P(c_k)$$

$P(x|c_k) = N(x|\mu_k, \Sigma)$ (Class conditional gaussian dist.)

Decision boundary is linear when all have same Cov. Matrix $\Sigma_k = \Sigma$

$$\rightarrow P(x|c_k) = N(x|\mu_k, \Sigma_k)$$

$$P(x|c_j) = N(x|\mu_j, \Sigma_j)$$

Assume $\Sigma_k = \Sigma_j = \Sigma$ (common)

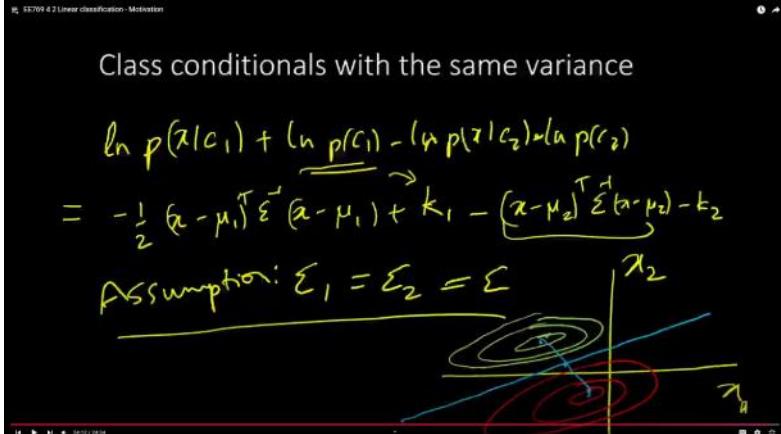
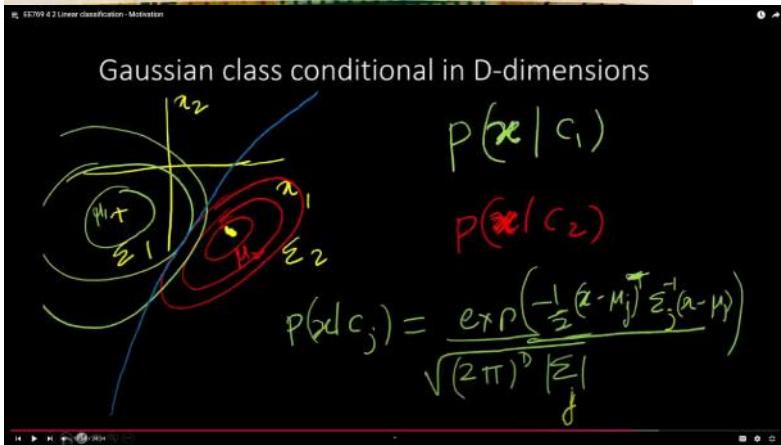
$P(x|c_k) \cdot P(c_k) > P(x|c_j) \cdot P(c_j)$

→ want boundary \nexists put equality and get condition

$$N(\mu_k, \Sigma) P(c_k) = N(\mu_j, \Sigma) P(c_j)$$

After solving we get boundary as $\log\left(\frac{P(c_k)}{P(c_j)}\right)$

$\left(\frac{\log[P(c_k|x)]}{\log[P(c_j|x)]} \text{ is under only if } \Sigma_k = \Sigma_j \right)$



decision boundary pass under gaussian assumption

$$P(t=1|x) > P(t=0|x)$$

$$P(x|t=1) \cdot P(t=1) > P(x|t=0) \cdot P(t=0)$$

$$\frac{P(x|t=1)}{P(x|t=0)} > \frac{P(t=0)}{P(t=1)}$$

$$\log \left(\frac{\exp(-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1))}{\exp(-\frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0))} \right) > \frac{P(t=0)}{P(t=1)}$$

$$\frac{-\frac{1}{2}(x-\mu_1)^T \Sigma^{-1}(x-\mu_1) + \frac{1}{2}(x-\mu_0)^T \Sigma^{-1}(x-\mu_0)}{2(\mu_1 - \mu_0)^T \Sigma^{-1}x + \mu_0^T \Sigma^{-1}\mu_0 - \mu_1^T \Sigma^{-1}\mu_1 + 2\log \left(\frac{P(t=0)}{P(t=1)} \right)} > 0$$

after grouping

$$w^T x + b > 0$$

$$w^T x + b = 0 \rightarrow \text{Boundary condition}$$

Bayesian Classifier Algorithm

① calculate prior probabilities

$$P(t=1) = \frac{n_1}{N}$$

$$P(t=0) = \frac{n_0}{N}$$

$$[n_1 + n_0 = N]$$

\downarrow

sample

where class

$x_0 \rightarrow$ set of x_0
where class is 0

$x_1 \rightarrow$ set of x_1
where class is 1

② estimate conditional prob.

$$P_0 = \frac{1}{|X_0|} \sum_{x \in X_0} x$$

$$P_1 = \frac{1}{|X_1|} \sum_{x \in X_1} x$$

③ calculate shared covariance matrix

$$\Sigma = \frac{1}{n} \sum_{x \in X} (x - \mu_0)(x - \mu_1)^T$$

$$P(x|t=0) = N(\mu_0, \Sigma)$$

$$P(x|t=1) = N(\mu_1, \Sigma)$$

④ apply Bayes' rule

$$P(x|t=0) \cdot P(t=0) > P(x|t=1) \cdot P(t=1)$$

for class 0.

SC109 4.3 Linear classification - Certain Bayesian classifiers are linear

Algorithm to build a Bayesian classifier

(1) Assume parametric forms of
 $P(x|c_i)$

... (2) Set $P(c_i) = \frac{N_i}{N}$

(3) Find the parameters of $P(x|c_i)$
M.L.

(4) Set decision criteria $y = \arg \max_{p(x|c_i)}$

Gradient descent for linear classifier

$$\text{Hypothesis} : h_i = \underbrace{\log \left(\frac{p(t_i=1|x_i)}{p(t_i=0|x_i)} \right)}_{\text{log of odds ratio}} = w^T x_i + b$$

Probability :

$$y_i = p(t_i=1|x_i) = \sigma(h_i) = \frac{1}{1+e^{-h_i}}$$

BCE Loss ;
Binary Cross Entropy

$$L_i = -t_i \log y_i - (1-t_i) \log (1-y_i)$$

\rightarrow find gradient of loss wrt w_k

$$\frac{\partial L_i}{\partial w_k}$$

$$\frac{\partial L_i}{\partial y_i} = -\frac{t_i}{y_i} + \frac{1-t_i}{1-y_i}$$

$$\frac{\partial y_i}{\partial h_i} = y_i(1-y_i)$$

$$\frac{\partial h_i}{\partial w_k} = x_{ik}$$

$$\frac{\partial L}{\partial w_k} = \frac{\partial L}{\partial y_i} \cdot \frac{\partial y_i}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_k} = \left(-\frac{t_i}{y_i} - \frac{1-t_i}{1-y_i} \right) y_i(1-y_i) x_{ik}$$

$$\boxed{\frac{\partial L}{\partial w_k} = (y_i - t_i) \cdot x_{ik}}$$

\rightarrow gradient descent update rule

$$w_{\text{new}} = w_{\text{old}} - \eta \sum_i (y_i - t_i) x_i$$

↳ learning rate

Adding Regularization

$$L_1 (\text{Lasso}) \quad L_{1,i} = L + \lambda \sum_k |w_k|$$

$$\frac{\partial L_{1,i}}{\partial w_k} = \frac{\partial L}{\partial w_k} + \lambda \cdot \text{sign}(w_k)$$

$$w_{\text{new}} = w_{\text{old}} - \eta \left(\sum_i (y_i - t_i) x_{ik} + \lambda \cdot \text{sign}(w_k) \right)$$

$$L_2 (\text{Ridge}) \quad L_{2,i} = L + \frac{\lambda}{2} \sum_k w_k^2$$

$$\frac{\partial L_{2,i}}{\partial w_k} = \frac{\partial L}{\partial w_k} + \lambda w_k$$

$$w_{\text{new}} = w_{\text{old}} - \eta \left(\sum_i (y_i - t_i) x_{ik} + \lambda w_k \right)$$

Elastic Net [$L_1 + L_2$] $L_{EN} = L + \lambda_1 \sum_k |w_k| + \lambda_2 \sum_k w_k^2$

 $w_{\text{new}} = w_{\text{old}} - \eta \left(\sum_i (\hat{y}_i - y_i) x_i + \lambda_1 \frac{\text{sign}(w_i)}{\text{old}} + 2\lambda_2 w_i \right)$

$L_1 \rightarrow$ eliminate var.
 \hookrightarrow doesn't encourage two correlated var to have same weight

$L_2 \rightarrow$ grouping effect on correlated var
 \hookrightarrow encourage two correlated var to have same weight

ROC are not symmetric like FP vs FN
 \hookrightarrow goal: minimize expected risk

- ① Accuracy $= \frac{TP + TN}{TP + TN + FP + FN}$
- ② Precision $= \frac{TP}{TP + FP}$
- ③ Recall $\text{(Sensitivity)} = \frac{TP}{TP + FN}$
- ④ Specificity $= \frac{TN}{TN + FP}$
- ⑤ F1 score $= 2 \cdot \frac{PR}{P+R}$
- ⑥ ROC curve: \rightarrow plot of True Positive Rate (Recall) against False Positive Rate (Specificity)
 \hookrightarrow AUC \rightarrow Area under ROC \square higher AUC \rightarrow better performing model
 $\hookrightarrow \in [0, 1]$
 \hookrightarrow by varying decision threshold

Lagrange multiplier

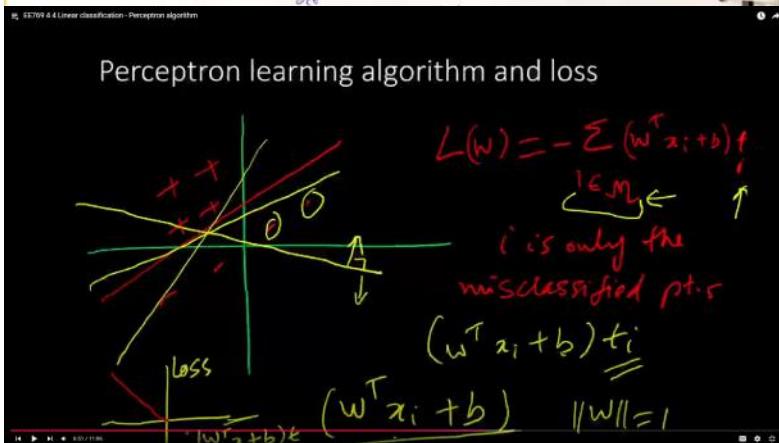
Problem: $f(x_1, x_2)$ minimization subject to $g(x_1, x_2) \leq 0$

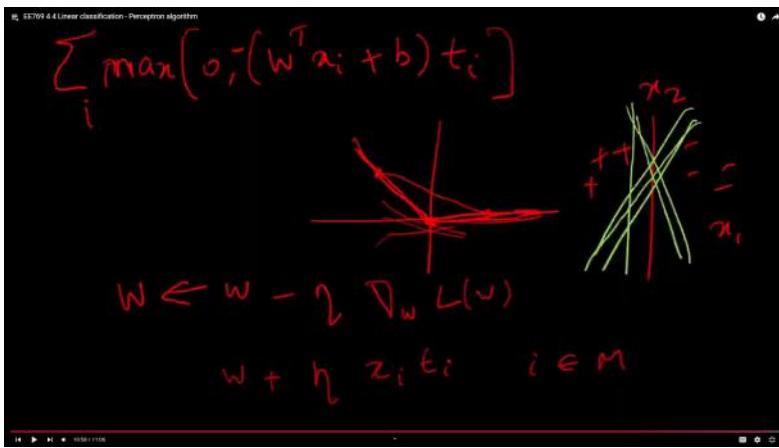
 $\nabla g(x) = \frac{\partial g}{\partial x}$
 $g(x+\epsilon) \approx g(x) + \nabla g(x) \cdot \epsilon^\top$
 $\nabla f + \lambda \nabla g = 0$

λ Lagrange Multiplier

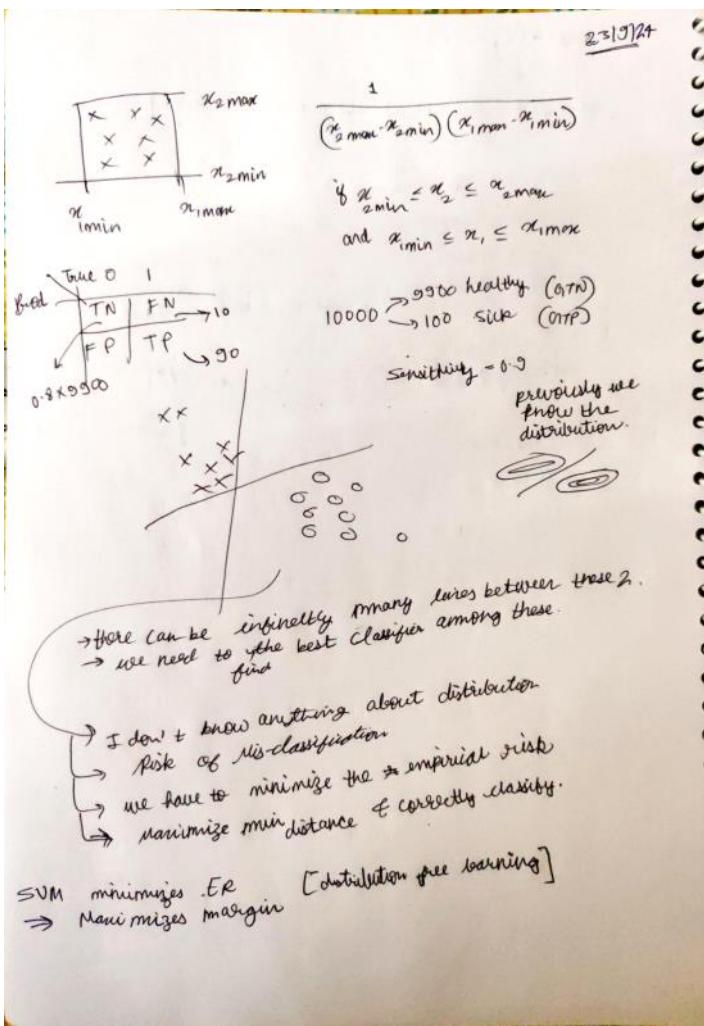
 $L(x, \lambda) = f(x) + \lambda g(x)$

constraint satiability condition $\nabla L = 0$
 $\frac{\partial L}{\partial x} = 0 \rightarrow g(x) = 0$



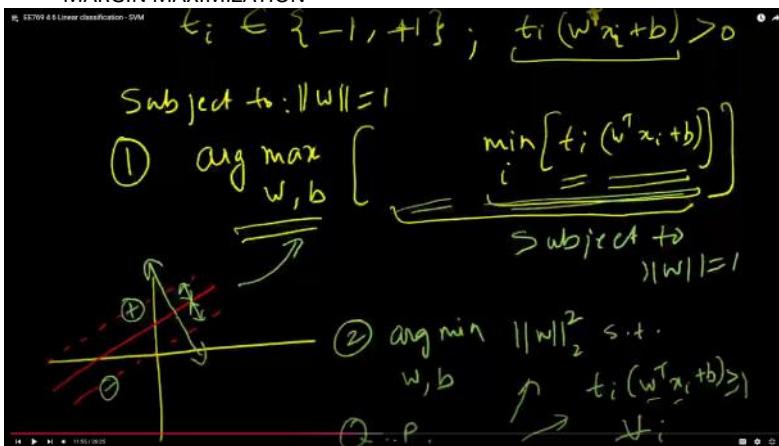


AFTER MIDSEM



SVM- EMPIRICAL RISK MINIMIZATION

MARGIN MAXIMIZATION



EE-557W9 4.5 Linear classification-SVM

$$\min_{w,b} \left(\lambda \|w\|^2 + \frac{1}{N} \sum_{i=1}^N [1 - t_i(w^T x_i + b)]_+ \right)$$

$$[a]_+ = \max(0, a) = \text{ReLU}(a)$$

In Bayesian \Rightarrow we minimize expected risk
[By assuming distro.]

SVM maximizes the separating margin $w^T x + b = \pm \epsilon$

Support vector \rightarrow point closest to hyperplane

No. of support vector in 2-D

accidental setting \rightarrow Non-accidental

Maximize over w and b

① $\text{Min}_{w,b} [\min_i \|w^T x_i + b\|]$ constraint optimization

subject to:

- $t_i \cdot t_i (w^T x_i + b) \geq 0$
- $1 \oplus 1 [Not 0 or 1]$
- $\|w\| = 1$

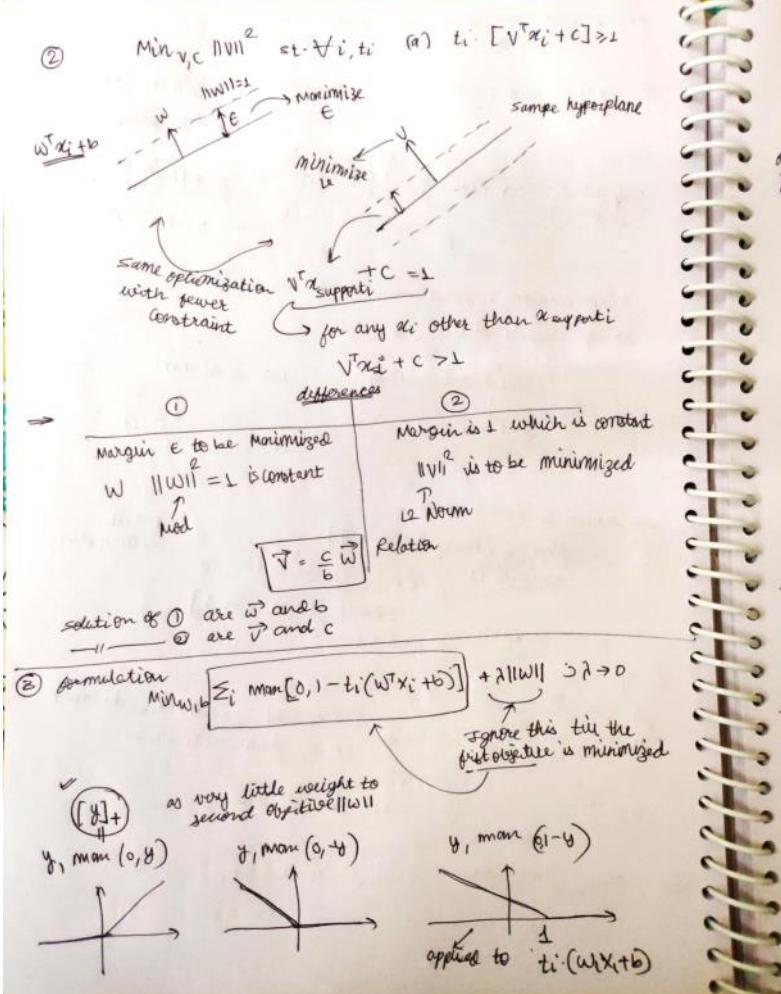
$t_i = +1$

$t_i = -1$

$(w^T x_i + b)$ is euclidean distance of x_i from $w^T x_i + b = 0$

similar formulation

② $\text{Min}_{w,b} \|w\|^2$, st. $t_i \cdot t_i (w^T x_i + b) \geq 1$



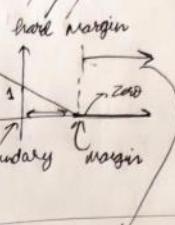
SVM Minimize separating margin

2.6 | 8/29

$$\text{① } \underset{\mathbf{w}, b}{\operatorname{Min}} \left[\underset{i}{\min} \|\mathbf{w}^T \mathbf{x}_i + b\| \right] \text{ s.t.}$$

- (a) $\|\mathbf{w}\| = 1$
- (b) $\forall i, y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 0$

assume linearly separable



- # $\max [0, y] = [y]$
 $\max [0, 1-y] = 1 - \min [y]$
 support \rightarrow defining the decision boundary
 (soft) decision boundary, only (entirely)
 defined by only support vectors
 (but like in other classifiers)
 adding new point can shift
 the decision boundary.

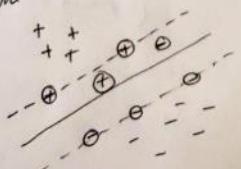
#

soft Margin,

No scope for regularization in Hard SVM.
 Hard SVM \rightarrow want to maximize Hard Margin condition
 $y_i = t_i(\mathbf{w}^T \mathbf{x}_i + b)$
 (linearly separable)

After Optimization
 all points will
 be here
 (Hard SUM)

\Rightarrow soft sum
 we want to minimize $\|\mathbf{w}\|^2 + C \sum_i \xi_i$ slack variable
 $\text{s.t. } \forall i, t_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$
 \Rightarrow for most points we want slack to be zero. $\nexists \xi_i > 0$
 \Rightarrow but for some point slack is greater than zero.
 \hookrightarrow some can include in the margin or outside the margin



$\nexists C \rightarrow \infty$
 all $\xi_i \rightarrow 0$
 \Rightarrow Hard SUM

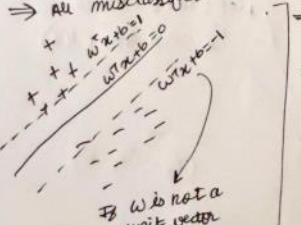
$\nexists C = 0$
 fully regularization
 many misclassification

Minimize $\|w\|^2 + C \sum_i \xi_i$ slack variable
 st. $\forall i t_i(w^T x_i + b) \geq 1 - \xi_i \quad \forall i \xi_i \geq 0$

Minimize $w, b \sum_i \text{hinge}(t_i[w^T x_i + b]) + \frac{\lambda}{2} \|w\|^2$
 $\lambda = 0 \rightarrow \text{Hard SVM}$ $\lambda \rightarrow 0 \rightarrow \text{Fully Regularized}$ $\lambda > 0$ $\lambda \propto \frac{1}{C}$
 In logistic regression: $\min_{w,b} CE + \frac{\lambda}{2} \|w\|^2$ (L2)

(Soft SVM) vs (L2 regularized logistic regression)
 ↓
 Loss function changed
 hinge loss cross entropy loss

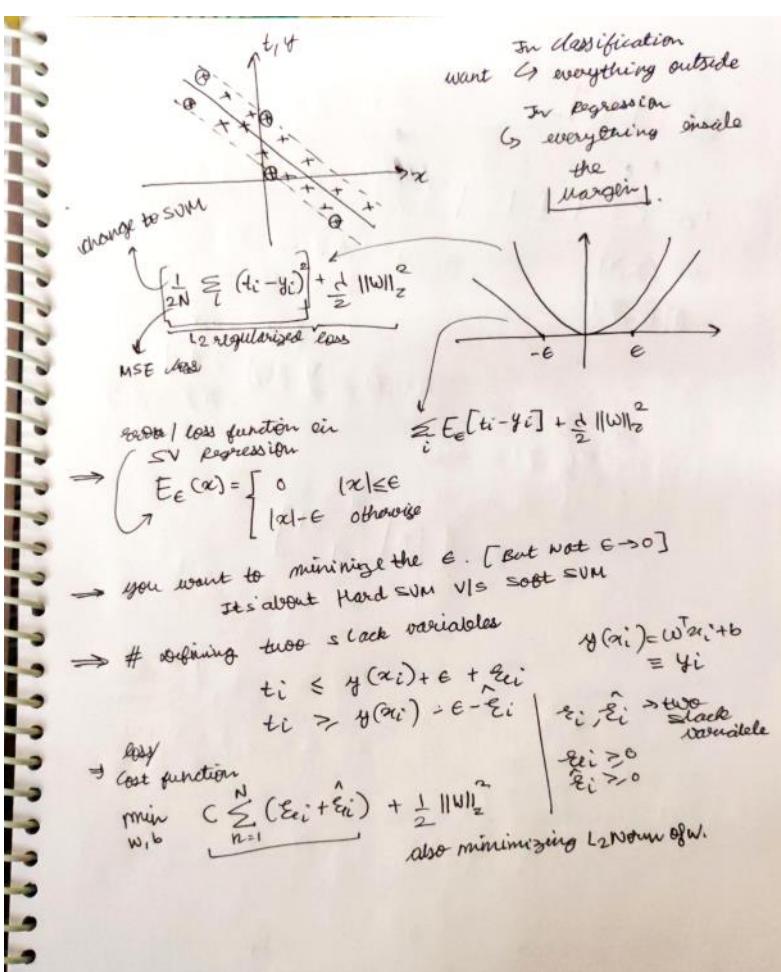
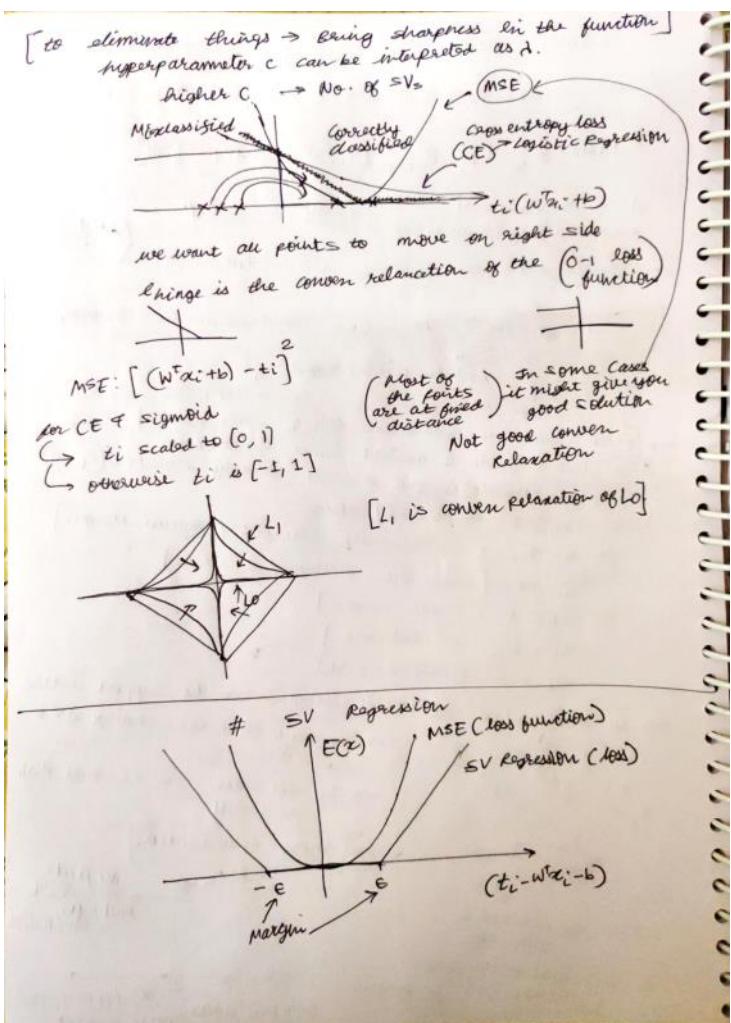
⇒ slack variable is individual loss of each points
 ⇒ slack variable is constant times hinge loss.
 ⇒ slack variable is zero for all correctly classified ($\xi_i = 0$)
 ⇒ $\xi_i > 1 \rightarrow$ misclassification
 ⇒ $0 < \xi_i < 1 \rightarrow$ correctly classified [beyond Margin]
 ⇒ $\xi_i = 0$ [Correct side of Margin / on margin]
 ⇒ $0 < \xi_i < 1$ [Inside Margin]
 ⇒ $\xi_i = 1$ [On Boundary]
 ⇒ $\xi_i > 1$ [Misclassified]

⇒ All misclassified and inside Margin are also support vectors.

 If w is not a unit vector
 Then $w^T x + b = \pm \epsilon$ for Margin

⇒ Any point that can change $w^T b$ are SV.
 ⇒ In soft SVM SVs are more than Hard SVM
 ⇒ Convex optimization

hinge
 Only these matter
 sparse subset

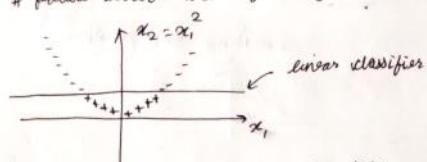
All points here don't matter for classification.



$\dots + + + + - - -$

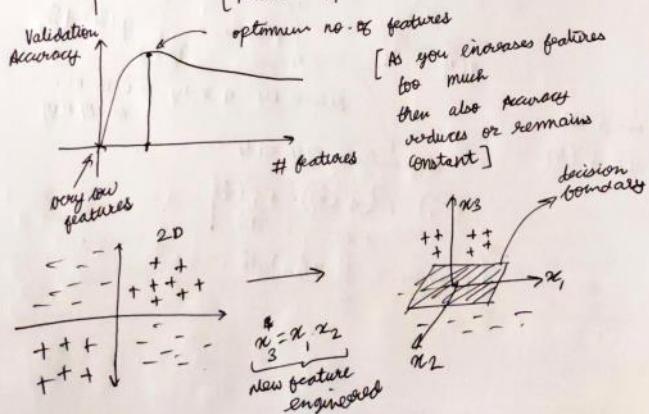
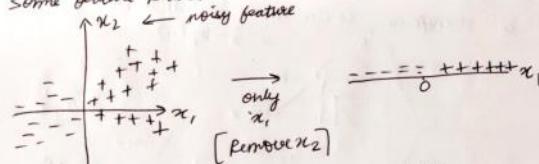
linear classifier in 2D its single threshold

feature selection and engineering

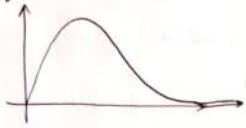


No L2 Regularized's can lead to instability
[in grouping effect]

Some feature reduction is also useful

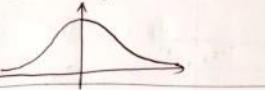


when your feature or distribution is not gaussian
 → you can use different transformation



$$\begin{aligned} \phi(x) &= \log(x+\epsilon) \\ \phi(x) &= x^p \\ \phi(x) &= e^{ax+b} \end{aligned}$$

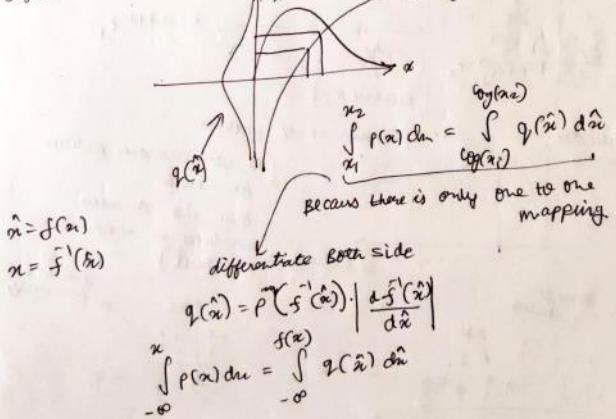
↓
 something like gaussian

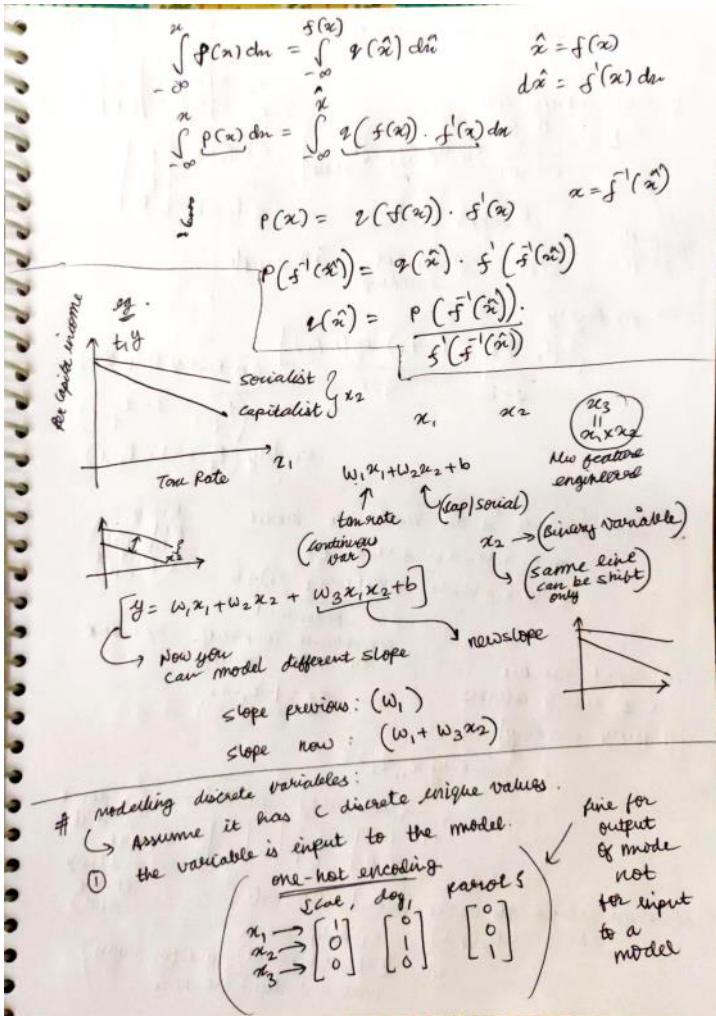


$x \sim p_x(x)$
 $\hat{x} = f(x)$ s.t. $q_{\hat{x}}(\hat{x})$ is a desired distribution

↪ find the transform $f(x)$?

objective is to transform to an easier distribution
 $\hat{x} = \log(x)$





DUMMY ENCODING

→ the variable is an output of model

$$t_{ij} = \begin{cases} 1 & \text{jth class} \\ 0 & \text{else} \end{cases}$$

cross entropy loss

$$L = -\sum_{i=1}^N \sum_{j=1}^C t_{ij} \log p_{ij} \quad \text{(jth class)}$$

only appropriate for o/p

ith sample jth one-hot encoding ith sample

for Binary

$$-\sum_{i=1}^N \left[t_i \log p_i + (1-t_i) \log (1-p_i) \right] \quad \text{(for binary only two j)}$$

$$(t_{j=0} = 1 - t_{j=1})$$

$$\text{similarly } (p_{j=0} = 1 - p_{j=1})$$

when variable is an input to model

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

$$\equiv w_1 x_1 + w_2 x_2 + w_3 (1-x_1-x_2) + b$$

for input we use C-1 binary columns

DUMMY encoding

only x_1, x_2

$x_1 \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad x_2 \begin{bmatrix} 0 \\ 1 \end{bmatrix} \quad x_3 \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

don't use integers to encode discrete variable like $\{1, 0, 2\}$

Simply mean (cat and parrot) are far away (cat and dog) are near.

Redundant
one column completely dependent on other column.
 $x_3 = 1 - x_1 - x_2$.

dummy
just
eliminate
dependent
var

we don't need it (Add adding to the bias)

DOMAIN SPECIFIC FEATURE SELECTION

MFCC- Mel-frequency Cepstral Coefficients

DOMRIN

2. pet type

	x_1	x_1	x_2	x_1	x_2	x_2
dog	1	1	0	1	0	0
dog	1	1	0	1	0	0
cat	2	0	1	0	1	0
parrot	3	0	0	0	0	1

↓
not ordinal
for input

↓
for output

MFCC (coeff.)

- ① windowing
- ② DFT \rightarrow power spectrum density
- ③ filter bank
- ④ DCT \rightarrow discrete cosine transform

features

Text (NLP)

I got lucky in test today

word dictionary

word:
I
got
wordc

Normalized frequency histogram

f_i frequent count of each word

$(\frac{f_i}{\sum f_i})$

fundamental model

→ pre-trained deep Neural network trained on very large data

U80 (480x720x3) → CNN → [] → generate some good feature → simple linear model

TF-IDF remove less redundant words like ["the"] → good feature reduction

FEATURE REDUCTION

$$P_{\hat{x}}(x) = \int_{-\infty}^x P_x(x) dx$$

$$Q_{\hat{x}}(\hat{x}) = \int_{-\infty}^{\hat{x}} g_{\hat{x}}(\hat{x}') d\hat{x}'$$

$$\hat{x} = Q_x^{-1} \left(\int_{-\infty}^{\hat{x}} g_{\hat{x}}(\hat{x}') d\hat{x}' \right)$$

$$\hat{x} = f(x)$$

$$d\hat{x}' = f'(x) dx'$$

result

$$\hat{x} = f(x) = Q_x^{-1}(P_x(x))$$

CDFs

feature reduction

 $x \in \mathbb{R}^{D \times 1} \rightarrow \hat{x} \in \mathbb{R}^{d \times 1}, d < D$

while doing so, we want to improve the accuracy, efficiency of model.

→ best subset selection of size d out of D .

total subsets $\binom{D}{d}$

Method ① Filters Based Method

↳ examine each feature individually $O(D)$

e.g. Classification: loop over j
evaluate utility of just using x_{ij}

Looking at its own validation performance
select top d subset.

train on top d .

Problem

x_j → Both $x_i \notin x_j$ individually
can't classify correctly
 x_i → so uses this combination
it [Student's t test]

↳ In regression: → Correlation with f and predicted f

Method ②: wrapper based

e.g. Forward selection
Backward elimination

* Forward selection \leftarrow set of remaining feature

$S \leftarrow \emptyset ; R \leftarrow \{x_1, x_2, \dots, x_n\}$ [initialization]

empty set

loop over $j \in D$

loop over elements $x_j \notin R$

[evaluate utility of including x_j in S]

pick the best feature, include x_j in S

remove from R

problem: when combⁿ of 2 give best result But individually don't

Alternative stopping criteria (use threshold on Accuracy)

↳ It cannot pick the best feature

Backward elimination

$S = \{x_1, x_2, \dots, x_n\} ; R \leftarrow \emptyset$

loop over j

loop over elements $x_j \in S$

[Evaluate utility of removing x_j from S]

remove the worst feature from S

Add it to R .

(Also called recursive feature elimination.)

↳ greatest ↑ in accuracy → give best feature

lowest ↓ in accuracy → give worst feature

③ # embedded

e.g. LASSO (L1 regularization)
elastic net (L1 + L2²)

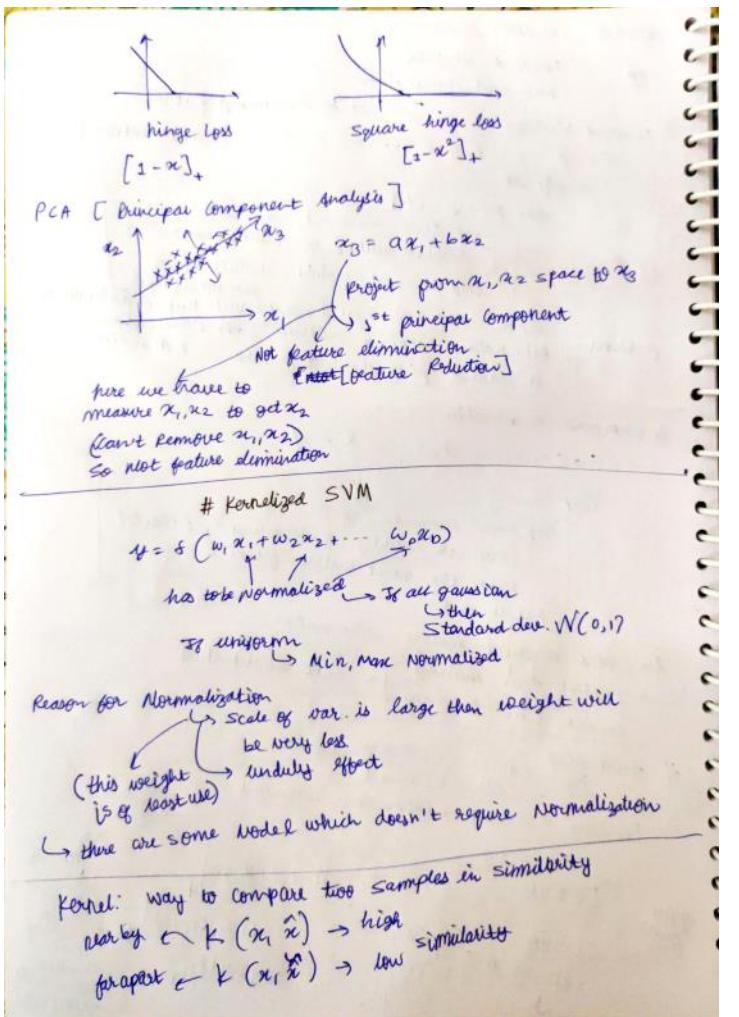
Learning & subset selection are integrated

L2 SVM: $\frac{1}{2} \|w\|_2^2 + \sum_i [1 - t_i(w^T x_i + b)]$ → hinge loss.

L1 SVM: $\frac{1}{2} \|w\|_2^2 + \sum_i [1 - t_i(w^T x_i + b)] + \sum_i |w_i|$ → sparsity in w

Support Vector → (because of sharp corner, only left & right half be support vector)

SUPPORT VECTOR REGRESSION (KERNELIZED SVM)



EE769-6-1 SVM - Introduction to Kernels

$$k(\mathbf{x}_i, \mathbf{x}_j) = k(\mathbf{x}_j, \mathbf{x}_i)$$

Gram matrix K

$$K = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_n, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

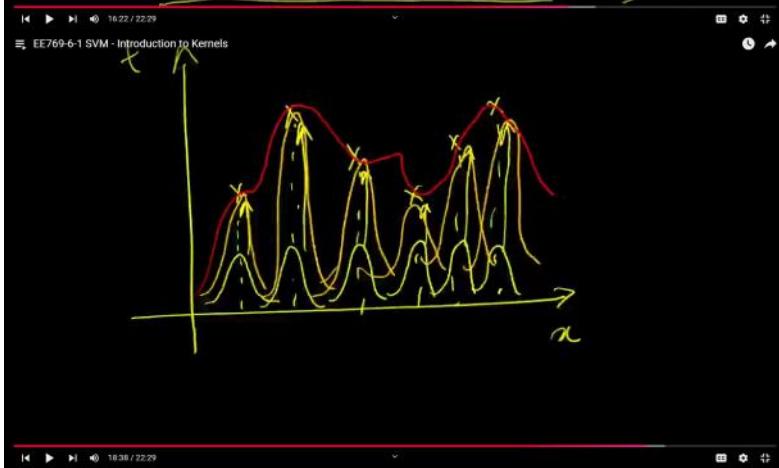
P.S.I

exp(- $\gamma \| \mathbf{x}_i - \mathbf{x}_j \|^2$)

Non-negative eigenvalues

Mercer kernel

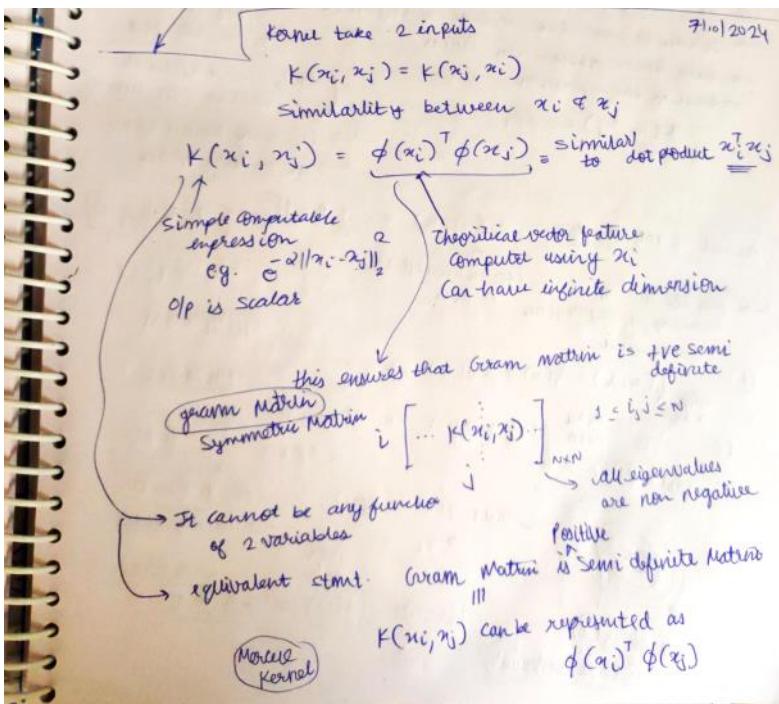
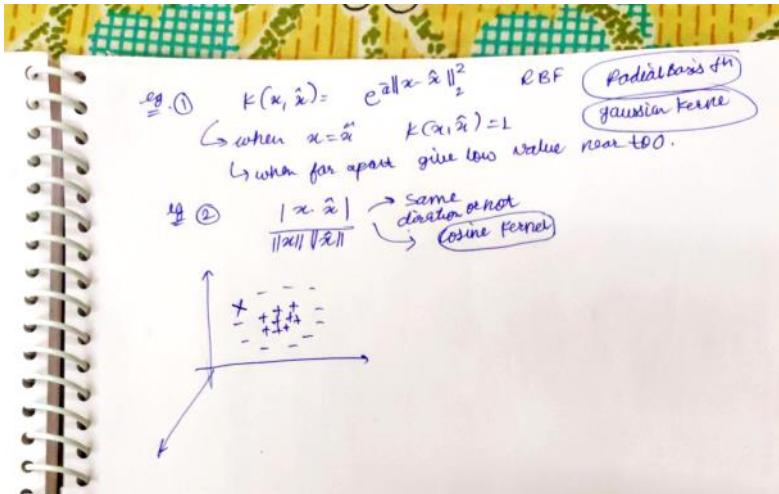
$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i) \phi(\mathbf{x}_j)$$



Kernel trick

- The kernel trick is make linear problems nonlinear
- Try to find a **dual of an expression** such that it uses $\langle x_i, x_j \rangle$ or $x_i^T x_j$ instead of x_i
- Replace $\langle x_i, x_j \rangle$ with kernel function
 $k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \phi(x_i)^T \phi(x_j)$
- The function is some nonlinear function, which need not be known to us so long as
- $k(x_i, x_j)$ is a Mercer kernel
- Implements an inner product in Hilbert space
- Has positive semi-definite Gram matrix

22.08 / 22.29



LAGRANGIAN

- we are interested in defining K not ϕ
 - computing kernel is much easier than $\phi(x_i)^T \phi(x_j)$ than dot product
 - you cannot just take any function of two vector, then say it as Mercer Kernel, you have to prove that gram matrix is PSD.
- # Kernel trick
- ① Take an objective function of x_i and try to change so that it becomes a function of $x_i^T x_j$ → DUAL form
 - ② Replace $x_i^T x_j$ with $\phi(x_i)^T \phi(x_j)$ or $K(x_i, x_j)$
- It can be applied to several optimization
 → Make linear problem non-linear
 → Simplify the expression
- $$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$$
- (Now, in \mathcal{X} space it will be non-linear)
 ↗ linear in Kernel space
 ↗ in feature space

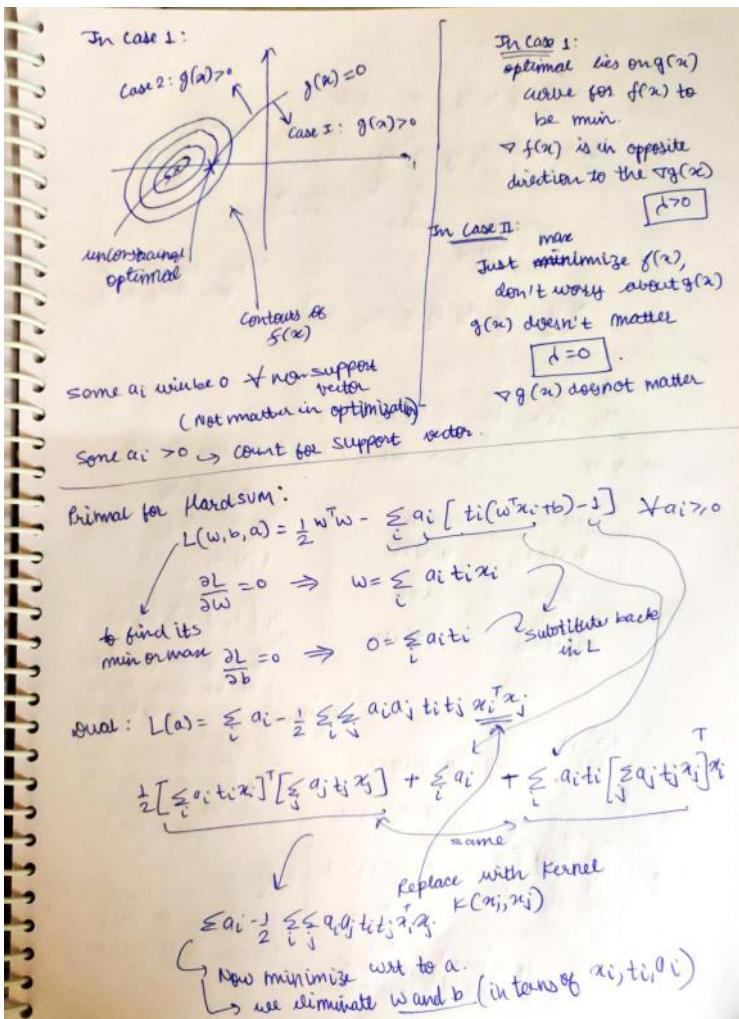
Dual SUM objective: $L(w, b) = \frac{1}{2} \|w\|_2^2 - \sum_i a_i [y_i (w^T q_i + b)]$

↳ this objective is basically a lagrangian loss in weights + bias
 ↳ Hard SUM

① Background
 $L(x, \lambda) = f(x) + \lambda g(x); \lambda > 0$ ↳ equivalent

② $\min_w \max_{\lambda} f(x), \text{ subject to } g(x) \geq 0$
 ↳ x-space ↳ R+0
 ↳ case 1: $g(x) > 0$ (here)
 ↳ case 2: $g(x) = 0$ (for $g(x) = 0$)
 ↳ see already know.
 ↳ n constraints
 $L(w, b) = \min_w \max_{\lambda} f(w) + \lambda g(w, \lambda)$

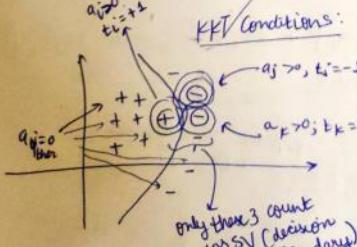
contours of $f(x)$
 (single constraint)



$y_i = w^T x_i + b$
 $y_i = \left[\sum_j a_j t_j x_j \right]^T x_i + b$
 $\underline{y}_i = \sum_j a_j t_j \underline{x}_j^T x_i + b$ replace $x_j^T x_i$
weights of training points test points training points
 $\underline{y} = \sum_j a_j t_j x_j^T x_i + b$ only SV will help you compute this.

 $b = \frac{1}{N_S} \sum_{i \in S} [t_i - \sum_j a_j t_j x_i^T x_j]$ similarly
no. of SV S \rightarrow set of support vectors ($a_i > 0$)

Case I: $g(x) = 0$ $\lambda > 0$ $\partial g(\alpha) = 0$	Case II: $g(x) > 0$ $\lambda = 0$ $\partial g(\alpha) = 0$	$a_i > 0 \quad \forall i$ $t_i y(x_i) - 1 > 0$ $a_i [t_i y(x_i) - 1] = 0$ One of the two has to be zero to be in Margin. $= 0$ for Non-SV $= 0$ for SV. [these point defined the decision boundary]
--	---	---


KKT conditions:
Only three count as SV (Decision Boundary)
 $0 = [a_i x_i^T x - a_j x_j^T x - a_k x_k^T x + b]$ (Decision Boundary linear in Kernel space)

EE769-6-2-1 SVM - Derivation of the dual form 1

Combining the objective function and the constraints, we get...

$$\rightarrow L(w, b) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i (w^T x_i + b) - 1]$$

s.t. $a_i \geq 0, \forall i$

① $\min L(x, \lambda) = f(x) + \lambda g(x); \lambda \geq 0$

② $\min f(x), \text{ subject to } g(x) \geq 0$



$\min L(w, b, a) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i (w^T x_i + b)]$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_i a_i t_i x_i$$

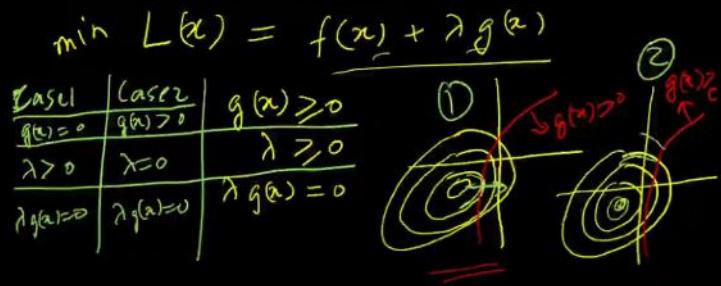
$$\frac{\partial L}{\partial b} = 0 \Rightarrow 0 = \sum_i a_i t_i$$

$$L(a) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j t_i t_j x_i^T x_j$$

$$y = \sum_i a_i t_i x_i^T x + b$$

$\xrightarrow{\text{test } \mathcal{K} \text{ training } k(x, x_i)}$ $\xrightarrow{k(x_i, x)}$

Solving the dual Lagrangian $L(a)$



$$L(a) = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j t_i t_j x_i^T x_j$$

$$L(w, b, a) = \frac{1}{2} \|w\|^2 - \sum_i a_i [t_i (w^T x_i + b) - 1]$$

$$\left(\begin{array}{l} a_i > 0 \\ \exists j \in S \\ \exists x_i \end{array} \right) \left(\begin{array}{l} a_i > 0 \\ t_i y(x_i) - 1 > 0 \\ a_i [t_i y(x_i) - 1] = 0 \end{array} \right) \left(\begin{array}{l} t_i \left[\sum_{j \in S} a_j t_j x_i^T x_j + b \right] \\ b = \frac{1}{N_S} \sum_{i \in S} \left[t_i - \sum_{j \in S} a_j t_j x_i^T x_j \right] \end{array} \right)$$

Analysing the KKT conditions

\leftarrow Not. S. V. $\notin S$

- Interpreting the points where $a_n=0$
 - For these points, $t_n[y(x_n)-1]$ does not need to be zero
 - These points do not define the decision boundary
 - These points can be moved around within most of the decision region without changing the boundary
- Interpreting the points where $a_n>0$
 - For these points, $t_n[y(x_n)-1]=0$
 - These points define the decision boundary
 - These points need to be retained after training for classification of test data, while rest of the points can be discarded
 - These points are called SUPPORT VECTORS

S. V.

decision boundary : $\sum_{i=1}^k \alpha_i K(x_i, x) - \sum_{j=1}^{l-1} \alpha_j K(x_j, x) - \sum_{k=l+1}^m \alpha_k K(x_k, x) + b = 0$
 finally is linear in feature space
 after replacing with kernel
 But our actual boundary is non-linear.

dual of soft SVM: for slack variable

PRIMAL: $L(w, b, \alpha, \xi) = \frac{1}{2} w^T w - \sum_i \alpha_i [t_i(w^T x_i + b) - 1] + \sum_i \xi_i$
 $\forall i; \alpha_i \geq 0, \xi_i \geq 0, t_i$
 $\xi_i \geq 0$
 we want to minimize sum of all slack variables

$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_i \alpha_i t_i x_i$
 $\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_i \alpha_i t_i = 0$
 $\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \alpha_i = C - \xi_i$

New KKT condition
 old KKT condn
 $\alpha_i > 0$
 $t_i [w^T x_i + b] - 1 + \xi_i \geq 0$
 $\alpha_i (t_i [w^T x_i + b] - 1 + \xi_i) = 0$
 Non constraint

two lagrangian with α_i with ξ_i
 total 6 KKT conditions
 New KKT condn

$\forall i; 0 \leq \alpha_i \leq C$
 $\alpha_i = 0 \Rightarrow \text{non-SVM (away from margin)}$
 $0 < \alpha_i < C \Rightarrow \text{SV (but on margin)}$
 $\alpha_i = C \Rightarrow \text{SV (but inside the margin)}$
 or it can be misclassified to
 hyper parameter we choose it.

- Steps:
- ① change objective to a lagrangian (a_i, ν)
 - ② square derivatives wrt. w, b to 0.
 - ③ substitute or w , expression in a_i 's and t_i 's.
 - ④ Manipulate the loss exp. till you get all $x_i^T x_j$.
 - ⑤ Re write the constraints as KKT conditions

Box constraints: 4 types of points

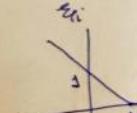
$a_i = 0 \Rightarrow$ Not s.v. \Rightarrow correctly classified and on correct side of margin

$0 < a_i < C \Rightarrow$ s.v. \Rightarrow correctly classified

on the margin

$a_i = C \Rightarrow$ s.v. inside the margin $\Rightarrow \xi_i \leq 1 \Rightarrow$ not misclassified

$\xi_i > 1 \Rightarrow$ misclassified



Hypoparams: ① C

② kernel-hypoparam

↳ scaling factors

↳ power of poly kernel

$\Rightarrow C$ in HardSVM. $C \rightarrow \infty$ \Rightarrow no constraint on a_i

Just $a_i > 0$

forcing $\xi_i = 0$

more regularization \Rightarrow less C
[reduce C more s.v.s]

all $\xi_i = 0$; nothing outside the margin

looking for hard margin

with less C type a_i is limited to only and many of a_i 's will be greater than 0.
So Number of S.V.s increases.

Analysing the KKT conditions

- \leftarrow Not.S.V. $\notin S$
- Interpreting the points where $a_i=0$
 - For these points, $t_i(y(x_i)-1)$ does not need to be zero
 - These points do not define the decision boundary
 - These points can be moved around within most of the decision region without changing the boundary
 - Interpreting the points where $a_i>0$
 - For these points, $t_i(y(x_i)-1)=0$
 - These points define the decision boundary
 - These points need to be retained after training for classification of test data, while rest of the points can be discarded
 - These points are called SUPPORT VECTORS

What if the classes are not separable?

- We need to allow some misclassification
- An error function that gives infinite error upon misclassification needs to be modified
- We introduce a slack variable $\xi_n \geq 0$ for each point n , such that ξ_n is:
 - 0, if the point is on the correct side of margin boundary
 - $|t_n - y(x_n)|$ for points outside the margin

$$L(w, b, \alpha, \xi) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [t_i y(x_i) - 1 + \xi_i] + C \sum_i \xi_i - \sum_i \mu_i \xi_i$$

↑
Lagrange

$\alpha_i \geq 0$

$t_i y(x_i) - 1 + \xi_i \geq 0$

$\alpha_i [t_i y(x_i) - 1 + \xi_i] = 0$

$\mu_i \geq 0$

$\xi_i \geq 0$

$\mu_i \xi_i = 0$

$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_i \alpha_i t_i x_i$

$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_i \alpha_i t_i = 0$

$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow \alpha_i = C - \mu_i$

$\hat{L}(a) = \sum_i \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j t_i t_j \underline{x_i^T x_j}$

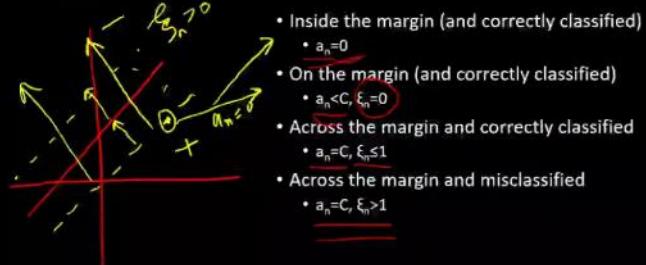
$0 \leq \alpha_i \leq C; \quad \sum_i \alpha_i t_i = 0$

$\alpha_i = 0; \text{Not S.V.}$

$C > \alpha_i > 0; \quad \xi_i \leq 1 \quad \text{correct, but inside margin}$

$\alpha_i = C; \quad \xi_i \geq 1 \quad \Rightarrow \text{incorrectly classified.}$

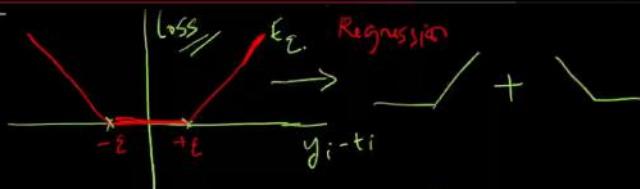
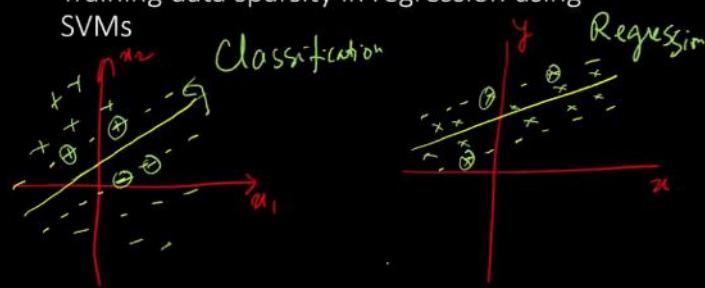
Box constraints: Four types of points



What defines the boundary?

- The following points define the boundary:
 - $\xi_n = 0$, on the margin (same as the separable case)
 - $0 < \xi_n \leq 1$, across the margin, but not misclassified
 - $\xi_n \geq 1$, misclassified points
- For these points: because $a_n > 0$
- Again, $y(x_n)$ calculation remains the same

Training data sparsity in regression using SVMs



$$\text{M.S.E.} : \frac{1}{2} \sum_i (y_i - t_i)^2 + \frac{\lambda}{2} \|w\|^2$$

$$\text{S.V.R.} : C \sum_i E_\epsilon(y_i - t_i) + \frac{1}{2} \|w\|^2$$

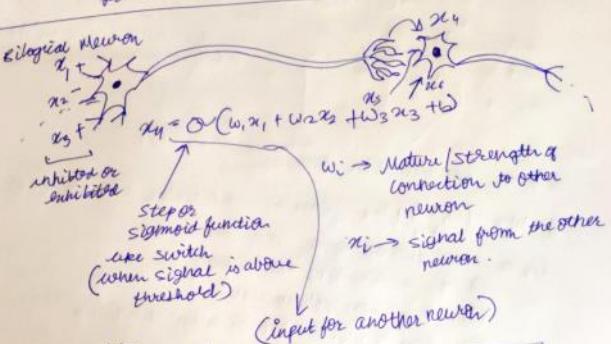
$E_\epsilon(y_i - t_i) = \max(0, |y_i - t_i| - \epsilon)$

Two slack vars. $t_i \leq y_i + \xi_i$; $t_i \geq y_i - \xi_i$

$$C \sum_i (\xi_i + \xi_i^*) + \frac{1}{2} \|w\|^2$$

↑ Non-linearity in models ## Neural Networks ##

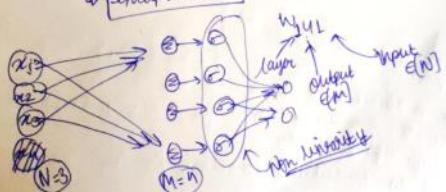
- $x \rightarrow$ Linear model \rightarrow loss
- feature engineering: $x \rightarrow \phi$ \rightarrow Linear model \rightarrow loss
- kernelized: $x \rightarrow k(x, \hat{x}) \rightarrow$ Linear model \rightarrow loss
- SVM: $\phi(x)^T \phi(\hat{x})$
- Neural network: $x \rightarrow$ trainable features \rightarrow Linear Model \rightarrow loss

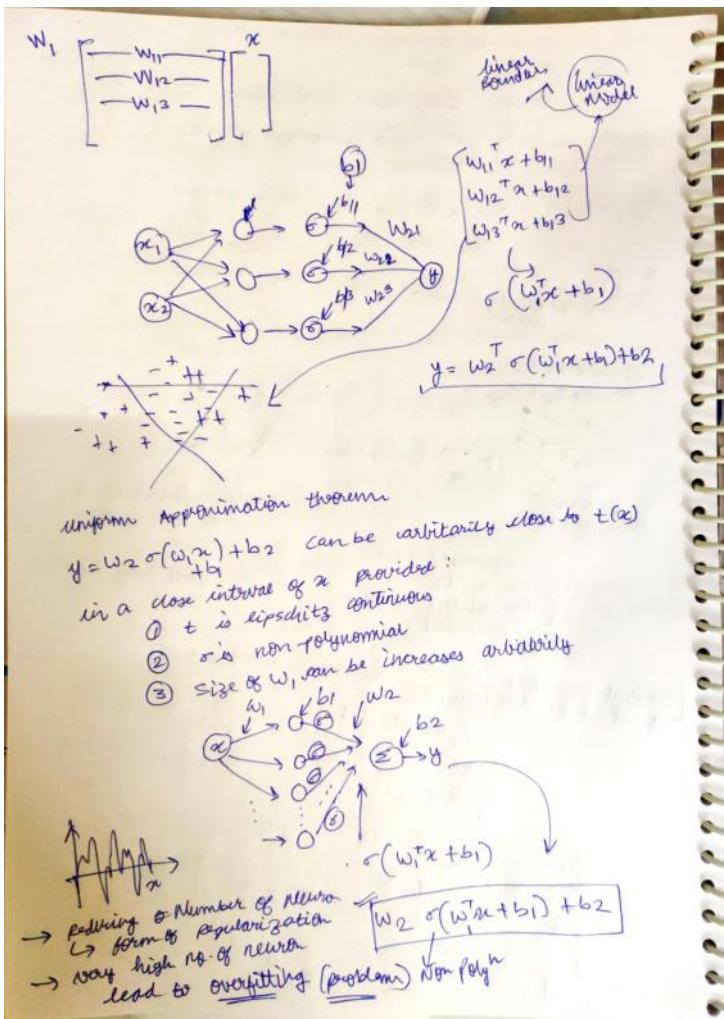


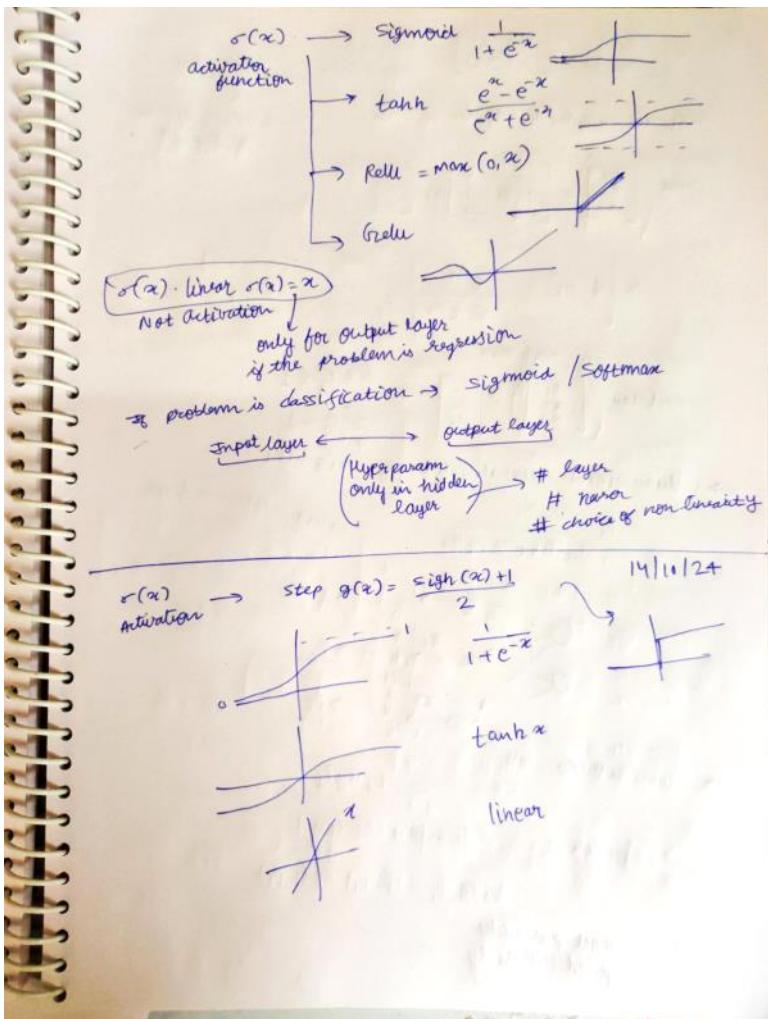
[Layered structure to use space effectively]

$x \rightarrow w \rightarrow z \rightarrow \sigma(w^T z)$ (encapsulation)

$x \rightarrow w \rightarrow z \rightarrow \sigma(w^T z)$ Interpretation







$\sigma\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \sigma(x_1) \\ \sigma(x_2) \end{bmatrix}$ pt wise application sigmoid
 Softmax \rightarrow smooth appear of Max smooth application of step fn

$\max\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

$\Rightarrow g(x_i) = \frac{e^{x_i}}{\sum e^{x_i}}$ Soft Max

$\text{softMax}\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} \quad 0 \leq p_i \leq 1 \quad \sum p_j = 1$

$\Rightarrow c\text{-dimensional generalization of sigmoid}$

$p_i = \frac{e^{x_i}}{e^{x_1} + e^{x_2} + e^{x_3}}$ using softmax layer

\rightarrow cross entropy loss
 $CE = -\sum_{j=1}^c t_j \log p_j$
 when $c=2$ [Binary CE loss]
 $-t_1 \log p_1 - (1-t_1) \log (1-p_1)$

$t \rightarrow$ target probability
 $p \rightarrow$ expected probability

BACK PROPAGATION

Multi-Class SVM

choice 1 class i vs class j
 $\binom{C}{2}$ SVMs

choice 2 class i vs rest SVMs
 $C-SVMs$

arg max $[w_j^T x_i + b_j]$ whichever j it is largest

for choice i $[j \rightarrow \text{class}]$

\rightarrow we usually use choice 2.

decided class for any point x in the space

$j = 1, \dots, C$
total no. of classes.

Back propagation uses chain rule of derivatives

$0 \rightarrow t \rightarrow \text{loss}$

t : vector for all params (weight, biases)

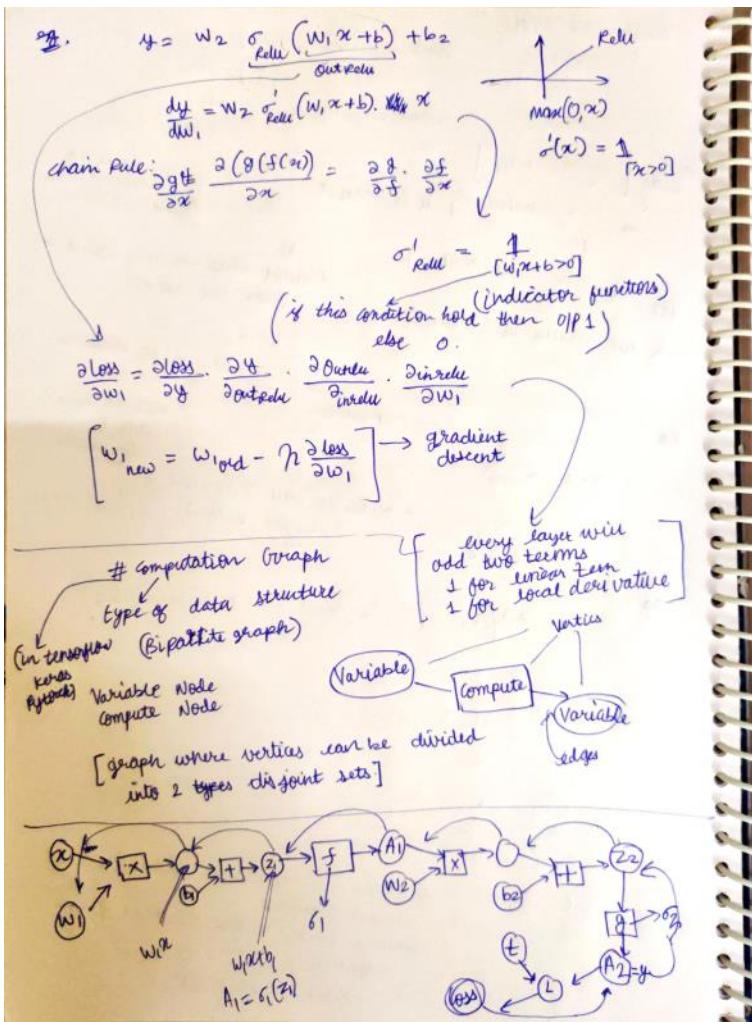
$\nabla_{\theta} t$

e.g. $l = \text{loss}(t, o_3(\underbrace{o_3(o_2(w_2 o_1(w_1 x + b_1) + b_2) + b_3)}_4))$

$o_1 \rightarrow 0.5 \rightarrow o_2 \rightarrow 0.5 \rightarrow o_3 \rightarrow 0$

$w_1 \rightarrow o_1$
 $w_2 \rightarrow o_2$
 $w_3 \rightarrow o_3$

$\frac{\partial l}{\partial y}$ derivative is computed to train the neural network coming back to the signal from output (y)
 \rightarrow BACK PROPAGATION



JACOBIAN
VANILLA GRADIENT DESCENT

[edges corresponds to the local derivatives]

vector valued f^n & JACOBIANS.

$\nabla_{\vec{x}} f \rightarrow \vec{x}^T f$

Vector ip
Scalar op

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

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

\vec{x}^S
 $n+m$

num inputs

functions / outputs

$Z = \vec{W} \vec{x} = f(\vec{x})$

$J_{\vec{x}}(Z) = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \vec{W}$

[Batch of samples]
 $\vec{x} \rightarrow$ Vector x batch size
 $x_i \rightarrow$ Inputs

Vanilla gradient descent

Iteration loop n
→ until some stopping criteria
loop over training sample
loop over i

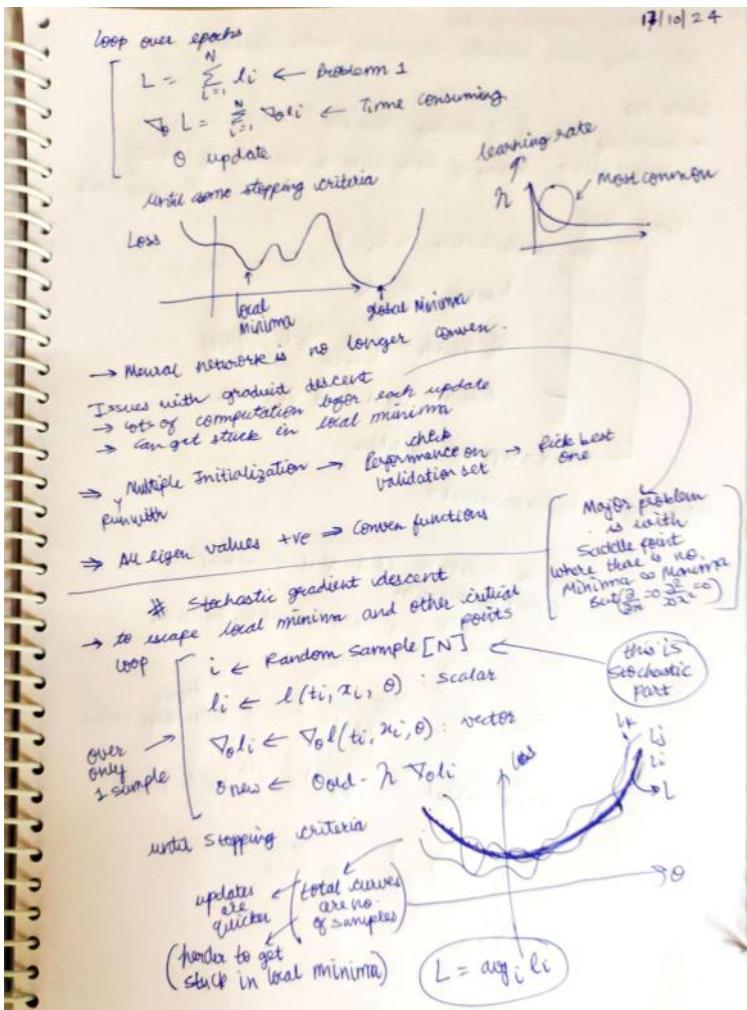
initialize loss $\leftarrow 0$

$l_i \leftarrow \text{loss}(\theta, x_i, t_i)$
 $\nabla_{\theta} l_i \leftarrow \text{grad}(\theta, x_i, t_i)$
 $\nabla_{\theta} l + = \nabla_{\theta} l_i$

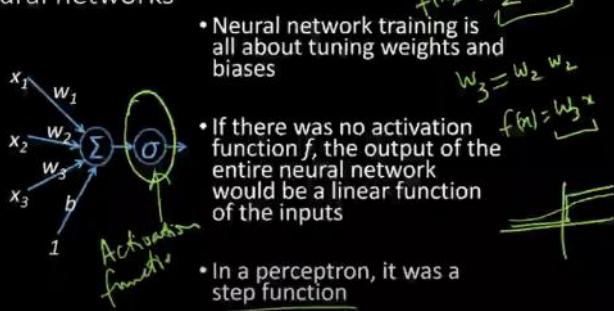
$\theta^D \leftarrow \theta^D - \eta \nabla_{\theta} l$
regular gradient descent

($l_i \leftarrow$ scalar)

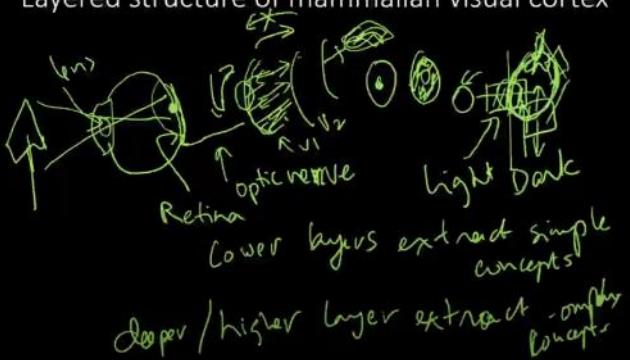
STOCHASTIC GRADIENT DESCENT



Activation function is the secret sauce of neural networks

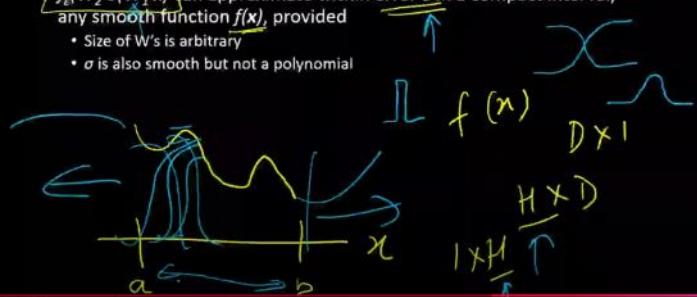


Layered structure of mammalian visual cortex

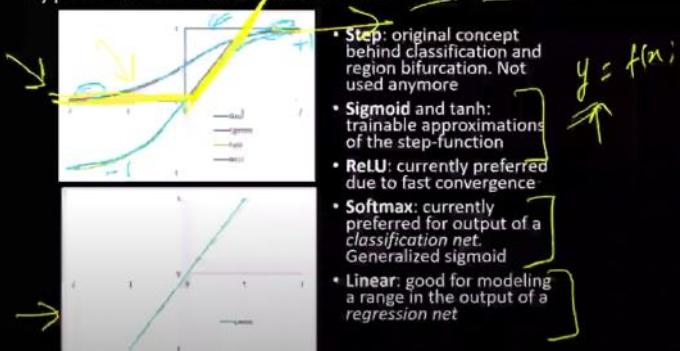


Universal approximation theorem

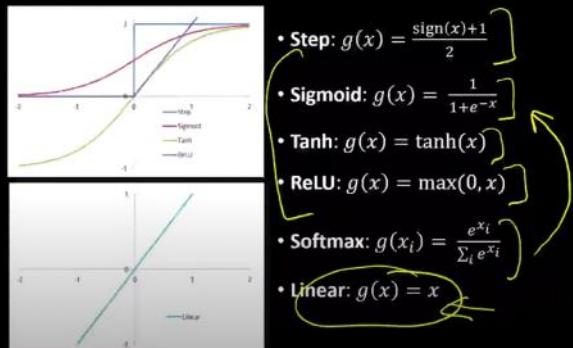
- $f_\epsilon(W_2 \sigma(W_1 x))$ can approximate within error ϵ in a compact interval, any smooth function $f(x)$, provided
 - Size of W 's is arbitrary
 - σ is also smooth but not a polynomial



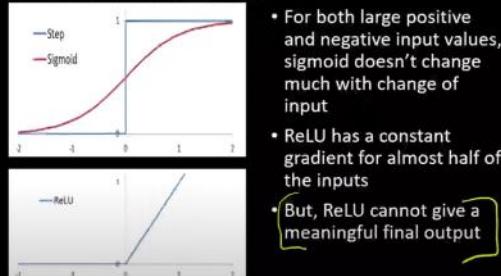
Types of activation functions

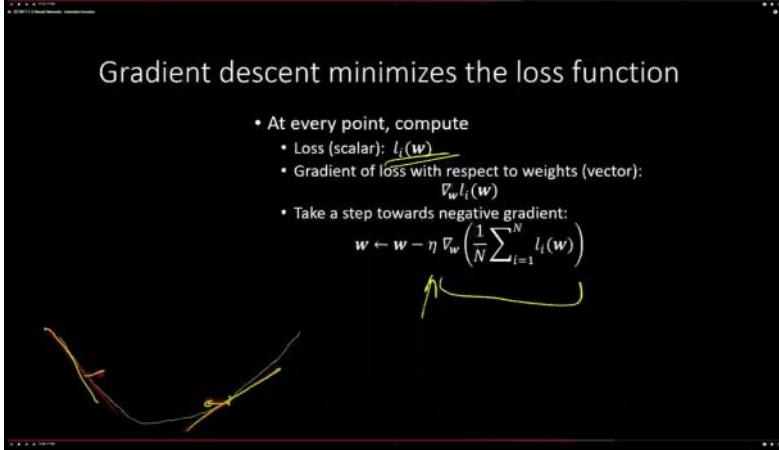
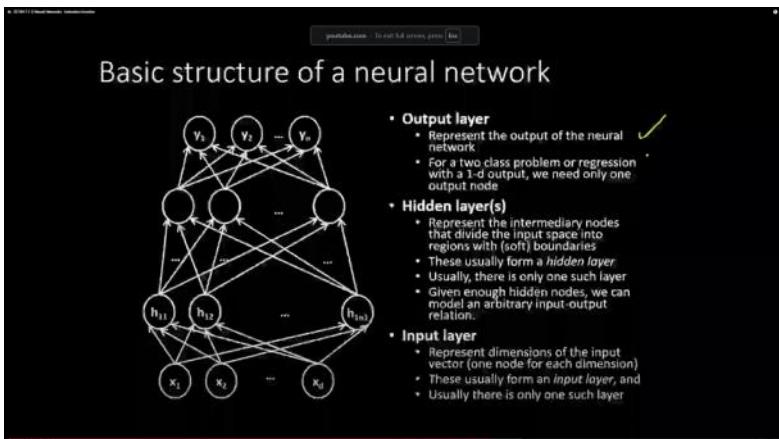


Formulas for activation functions



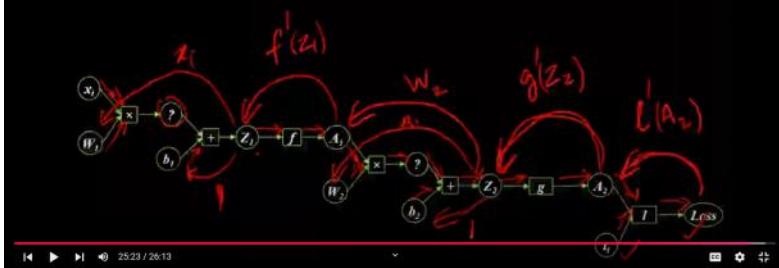
The problem with sigmoid is (near) zero gradient on both extremes





EE769-7-2-1 Neural Networks - Chain rule and backprop

1. Make a forward pass and store partial derivatives $F^{(n)}$
2. During backward pass multiply partial derivatives $F^{(n)}'$

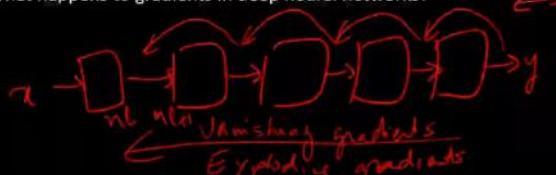


Some questions

- What if we scale all the weights and biases by a factor?

$$(Wx + b) \rightarrow (\alpha Wx + \alpha b) \rightarrow \text{activation}$$

- What happens to gradients in deep neural networks?



The perfect step size is impossible to guess

- Goldilocks finds the perfect balance only in a fairy tale

The step size is decided by learning rate η and the gradient

Learning rate scheduling.

BATCH GRADIENT DESCENT

DOUBLE DERIVATIVES(HESIAN MATRIX)

The perfect step size is impossible to guess

- Goldilocks finds the perfect balance only in a fairy tale

The step size is decided by learning rate η and the gradient

Learning rate scheduling.

Saddle points, Hessian and long local furrows

- Some variables may have reached a local minima while others have not
- Some weights may have almost zero gradient
- At least some eigenvalues may not be negative

Adding momentum in lieu of second derivative

ADDING MOMENTUM IN GRADIENT DESCENT

Advanced neural networks

Learning objectives

- Write NN loss functions for some intermediate ML problems
- Write some properties of natural signals
- List advantages of convolution and pooling layers
- Explain the reason behind exploding and vanishing gradients
- Write how LSTM units solve the exploding gradient problem

0:42 / 12:40 - Learning Objectives

① Noisy weight updates
 ② may leave compute hardware under utilized

Batch GD

→ randomly permute training sample
 → divide into training data into B batches of size $\frac{N}{B}$ samples each

Epoch loop

for each batch $b \in [B]$

$$L_{\text{batch}} = \sum_{i \in b} \ell_i$$

$$\nabla_{\theta} L_{\text{batch}} = \left(\sum_{i \in b} \nabla_{\theta} \ell_i \right) / (Nb)$$

$$\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla_{\theta} L_{\text{batch}}$$

until stopping condition

Batch size → hyper parameter

Intuition hypothesis class of a threshold

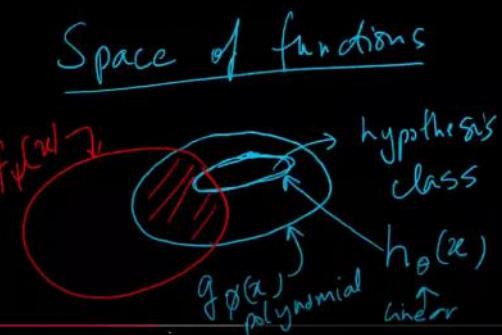
Space of function $f(Y|X)$ → data generation function

Non-linear

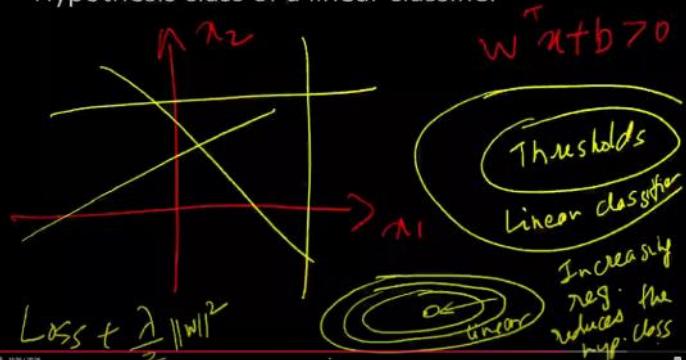
Questions for today's lecture

- Why combine models?
 - Multiple models are available, but none of them are perfect
 - Multiple types of features can be extracted for a given pattern
 - Certain complimentary properties exist among different models and different features.
- Issues
 - How many classifiers are needed?
 - What kind of models should be used?
 - What features to be used in each classifier?
 - How to combine results from different classifiers?

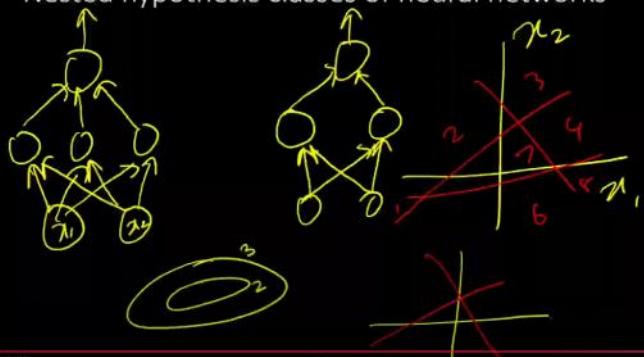
Introduction to hypothesis class



Hypothesis class of a linear classifier



Nested hypothesis classes of neural networks



Ensembles for classification

Assumption of independence of models

Classifier 1



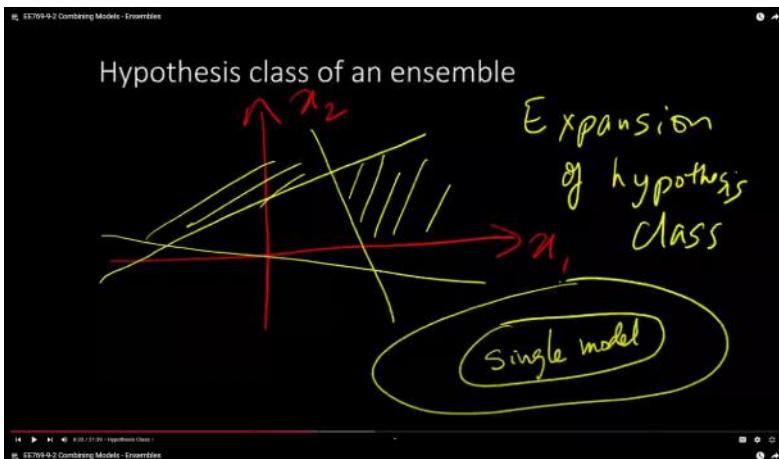
Classifier 2

$$y_i = \text{mode}_{m=1}^3 f_m(x_i)$$

Classifier 3

$$\begin{aligned} p^2(1-p) + p^3 &= \\ \frac{p}{2} & \quad \text{if } p > \frac{1}{2} \end{aligned}$$

Weak Learner

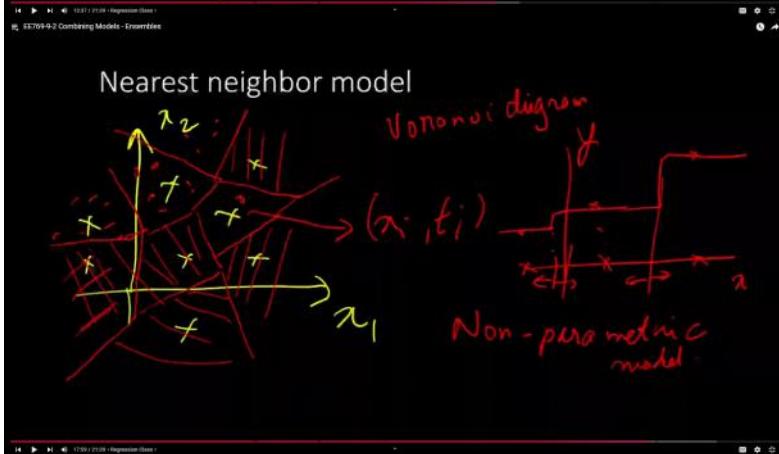


Ensembles for regression

$$y_i = \underbrace{f_1(x_i)}_{E(t_i - y_i)^2} + \underbrace{f_2(x_i)}_{E(f_1(x_i) - t_i)^2} + \underbrace{f_3(x_i)}_{E(f_2(x_i) - t_i)^2} / 3$$

$$= \frac{\sigma_1^2 + \sigma_2^2 + \sigma_3^2}{3^2}$$

$$E[(f_1(x_i) - t_i)(f_2(x_i) - t_i)] = 0$$



very simple model
Randomly select input dim
Find threshold on x_1 for model

different ways of combining model Hypothesis class

Ensembles for classification probabilities

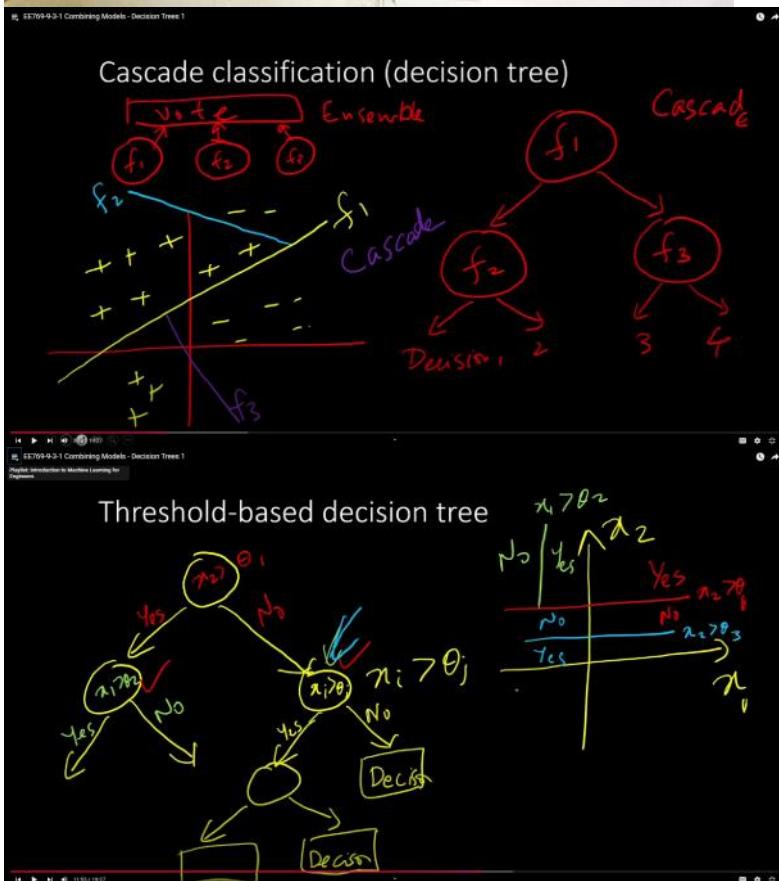
h_1	h_2	h_3	y	probabilities
0	0	0	0	p^3
0	0	1	0	$p^2(1-p)$
0	1	0	0	$p^2(1-p)$
0	1	1	1	$p(1-p)^2$
1	0	0	0	$p^2(1-p)$
1	0	1	1	$p(1-p)^2$
1	1	0	1	$p(1-p)^2$
1	1	1	1	$(1-p)^3$

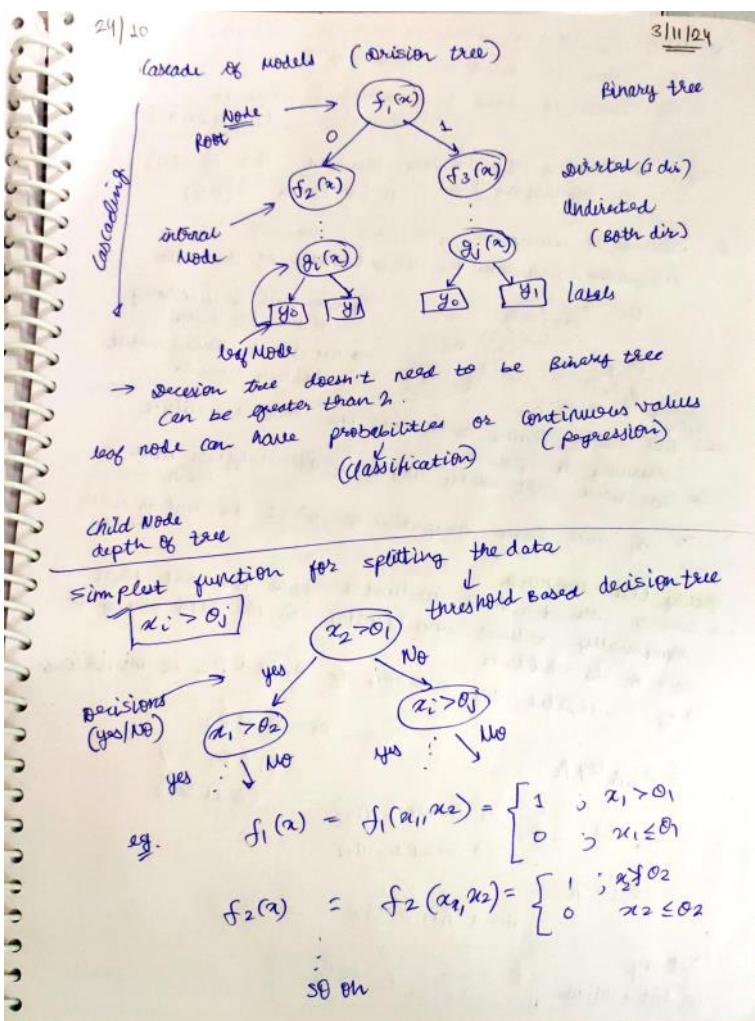
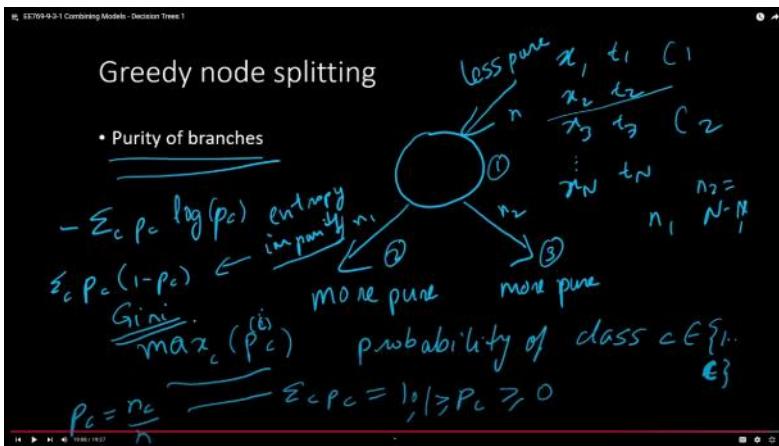
Ensembles for regression

$$h_i = \frac{h_1 + h_2 + h_3}{3}$$

$$\text{Variance} = \frac{\sum (h_i - \bar{h})^2}{n}$$

$$y_i = (w_1^T x_i + b_1 + w_2^T x_i + b_2 + w_3^T x_i + b_3)/3$$

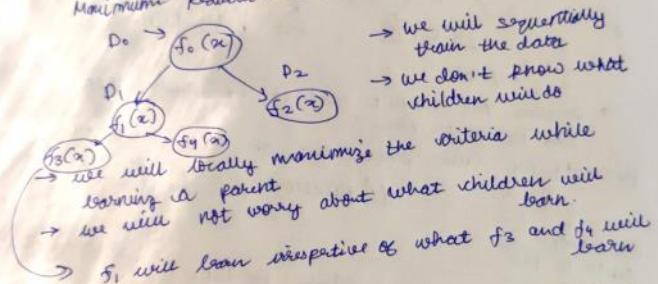
$$E(t_i - y_i)^2 = \frac{\sigma_1^2 + \sigma_2^2 + \sigma_3^2}{3^2}$$




- questions:
- ① How to choose a threshold for decision node?
 - ② How to select a threshold?
 - ③ when to have leaf nodes (when to TERMINATE?)

say N training pts \rightarrow Max. threshold $N-1 \Rightarrow O(N)$
 D dimensions $\rightarrow D \cdot (N-1) \Rightarrow O(ND)$

criteria to compare threshold classifier
 Maximum Reduction in uncertainty of the data



greedy tree learning
 among all threshold evaluated pick the one that minimally reduces uncertainty in the data going to each of its children.
 \Rightarrow e.g. weighted avg. entropy of P_1 and P_2 is minimized.

$$P_0 = P_1 \cup P_2$$

$$P_1 \cap P_2 = \emptyset$$

(Null set)

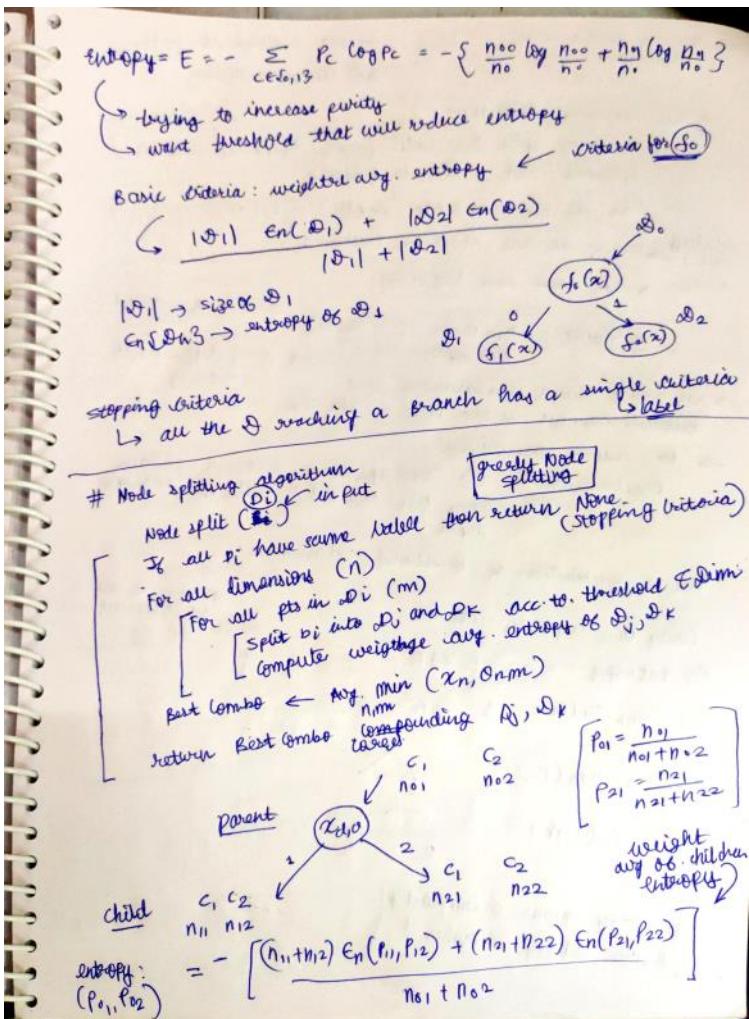
$\oplus = D_0$
 \rightarrow label 0 n_{00}
 \rightarrow label 1 n_{01}

$$|D_0| = n_0$$

$$n_{00} + n_{01} = n_0$$

Entropy for distribution

$$E = -\sum_{c \in \{0,1\}} p_c \log p_c$$



No. of leaf nodes = n \times each sample require its sole decision.
 (minimum) has its own node decision.

* Tree learning Algorithm

recursively split D_o until [none / stopping criteria] is returned and make leaf node

$$\rightarrow \text{Max. possible leaf Node (depth)} = O(N)$$

\rightarrow No. of internal nodes = # leaves - 1
 (including root)

NOTE: deep trees can **OVERFIT**

N leaves: Min Order $O(\log N)$ \rightarrow well structured
 Max depth $O(N)$ \rightarrow very difficult to classify

\rightarrow we don't want imbalanced tree
 But we can get if our data is like this

\Rightarrow to reduce overfitting
 regularization: ① stop splitting below a certain depth
 ② prune tree beyond a certain entropy depth

Tree \rightarrow combination of threshold classifiers

Purity or impurity criteria:

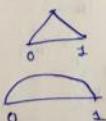
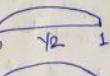
$$① \text{Entropy} = - \sum_c p_c \log p_c$$

$$② \text{Gini Index} = 1 - \sum_c p_c^2$$

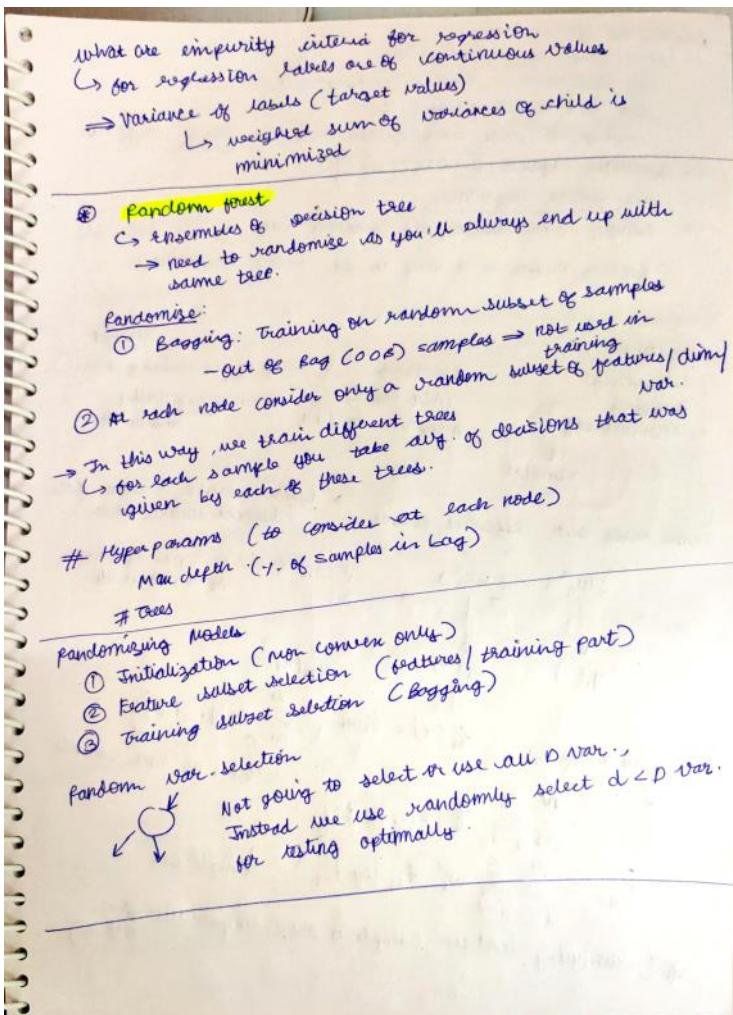
$$③ 1 - \max(p_c)$$

$$④ p_c(1-p_c)$$

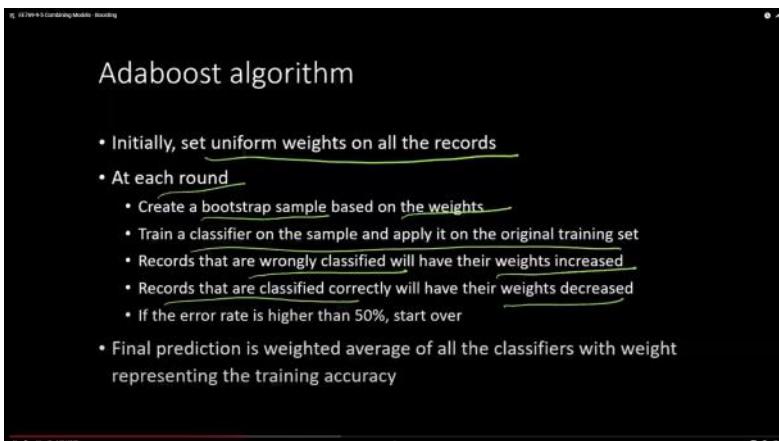
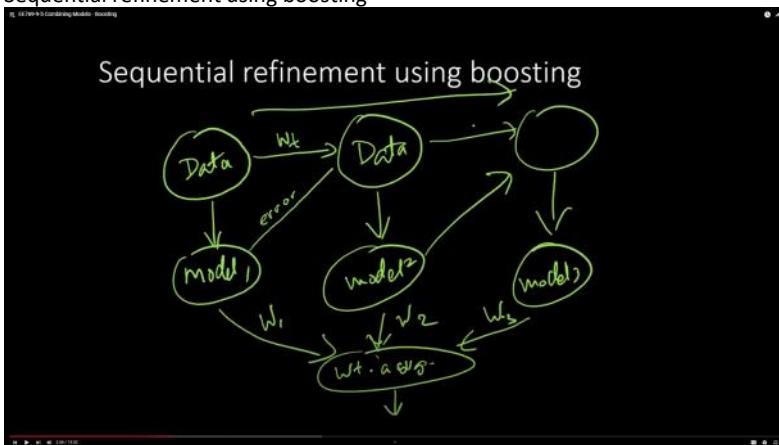
graph of matrix
 for different values
 of p_c .



[\circ entropy means higher entropy \rightarrow higher "impurity"]



Sequential refinement using boosting



Adaboost visualized
 S is training data

Initialize $D^{(1)} = \left[\frac{1}{n}, \frac{1}{n}, \dots, \frac{1}{n} \right]$

For m in $1 \dots M$

Learn $w_L(D^{(m)}, S) \leftarrow$

$\epsilon_m = \sum_i D_i^{(m)} \mathbb{1}[t_i \neq y_i]$

$w_m = \frac{1}{2} \log \left[\frac{1 - \epsilon_m}{\epsilon_m} \right]$

$D_i^{(m+1)} = \frac{D_i^{(m)} \exp(-w_m t_i y_i)}{\sum_j D_j^{(m)} \exp(-w_m t_j y_j)}$

$\text{Sign}\left(\sum_{m=1}^M w_m y_i^{(m)}\right)$

- Properties of RF:
- good idea of generalization
 for each tree \rightarrow OOB \rightarrow testing data
 avg OOB label across tree
 converge to their generalization error \rightarrow law of large no.
 - guarantee against OVERFITTING
 - give feature importance
 - handle both univariate and continuous var
- Random shuffle is a way to get feature importance

28/10/24

Cross-validation
 Bagging
 Ensemble learning

Boothing (Ada Boost)
 Adaptive boosting

3/11/24

(way of combining model)
 sequentially combination

same node but different subset \rightarrow different model, same data (different initialization)

$y_1(x), y_2(x), \dots, y_M(x)$

$y_M(x) = \text{sign} \left(\sum_{m=1}^M w_m y_m(x) \right)$

weights for each sample

Gross entropy loss

$L = - \sum_{i=1}^N \sum_{j=1}^c \mathbb{1}(t_{ij} \log t_{ij})$

$(L_w = - \sum_{i=1}^N \sum_{j=1}^c w_i t_{ij} \log t_{ij})$ Sample Bit

* classifying certain sample is more important than other

At beginning all weights have same importance
 ↳ are equal
 → then decrease the importance for that sample that
 are misclassified and vice-versa.
 (second barrier to focus more on that pattern)

AdaBoost Alg

Set all $w_n^{(1)} = \frac{1}{N}$ for $n=1, \dots, N$ (initialization)

loop over $m=1$ to M

Compute loss $J_m = \sum_{n=1}^N w_n^{(m)} I(y_m(x_n) \neq t_n)$ ↳ indicator function

Normalized error $E_m = \frac{J_m}{\sum_{n=1}^N w_n^{(m)}}$

Compute model weight $\alpha_m = \log \left(\frac{1 - E_m}{E_m} \right)$ ↳ exponential function (can use any other function in that case)

update sample rate $w_n^{(m+1)} = w_n^{(m)} \exp \{ \alpha_m I(y_m(x_n) \neq t_n) \}$

If error made → Indicator function \Rightarrow weight for those where error is large

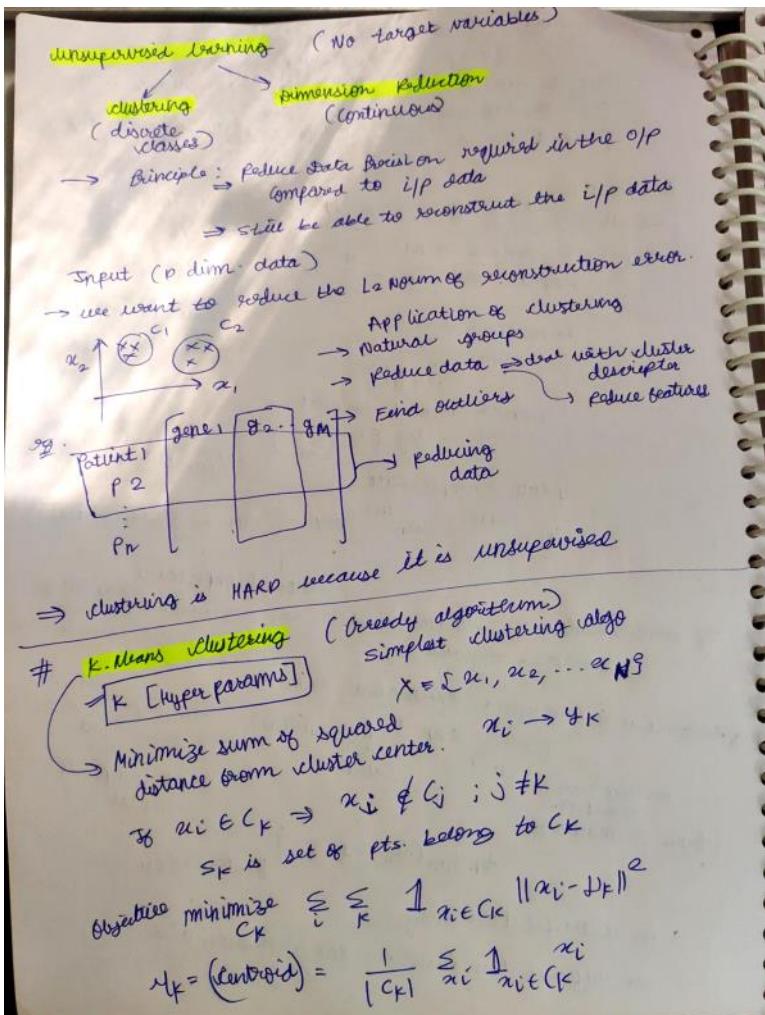
No error → generate = 0

⇒ AdaBoost is minimizing overall loss function $E = \sum_{n=1}^N \exp \{ -t_n \cdot f_m(x_n) \}$ ↳ label

$f_m(x) = \sum_{i=1}^m \alpha_i y_i(x)$ ↳ m^{th} composite model

$t_n y_m(x_n) = 1 - 2 \mathbb{I}_{y_m(x_n) \neq t_n}$

for classification only
 use different loss functions for regression also



k means clustering

