

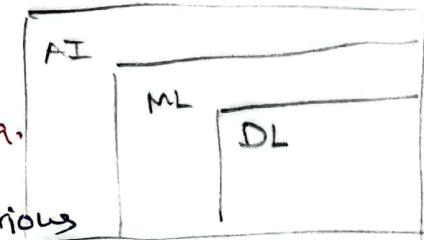
Li: what & why?

- * machine learning is a subfield of computer science, that evolved from the study of pattern recognition & computational learning theory.
 - use Algos to iteratively learn from data.
 - Allow computers to discover patterns, without explicitly programming where to look.

Q: What is machine learning?

Ability of computers to learn from data.
from various sources.
(with noise redundancy)

- supervised learning
- unsupervised learning
- Reinforcement learning



i) supervised learning:-

we got training data; with labelled datums.

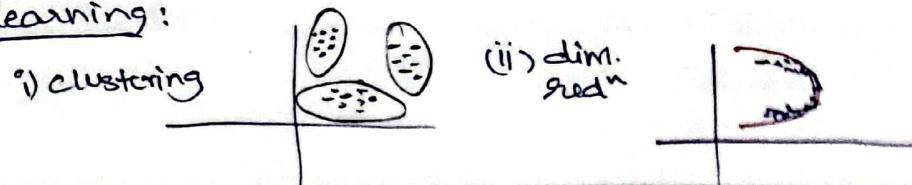
Goal is to predict labels on new datums.

- ### ii) Unsupervised :- data studied through clustering, mining associations, dimen. redn!
- we know nothing about fruits - only their parameters.
 - color
 - size
 - shape.
 - we got to group fruits of same kind.
 - no meaning of labels.

- * both approaches have tradeoffs.
 - ① → human annotated data
 - ② → more computation.

→ Three canonical learning problems:-

- 1) Regression - Supervised :- when predicting a continuous value.
not only linear, even quadratic.
- 2) Classification - Supervised :- predict color of fruit; translation of english.
- 3) Unsupervised learning :

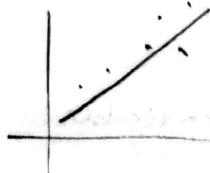


L2.1

- * In supervised learning, the desired outputs are provided on training data. where-in unsup. learning; data is analysed & studied through clustering, mining associations, dimensionality reduction etc. into different classes.

- * 3 canonical learning problems!

1) Supervised regression

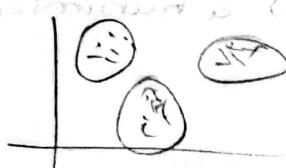


2) Supervised classification:



3) Unsupervised learning:-

• clustering



• dimensionality reduction



→ Supervised learning:-

Training data $\{ (x_i, y_i) \in X^*Y \}$

Function $F: X \rightarrow Y$

Learning:- finding $\hat{f} \in F$

s.t. $\hat{f}(x_i) \approx y_i$ here we have concept of error.

Prediction:- $y = \hat{f}(x)$ \rightarrow new data 'x'.

1) Noise in data:-

- data might have noise, random human errors.
- - Most methods deal with expectations as they minimize effect of error.
- Better to find outliers & clean data in 1st step.
"data cleaning."

2) How to predict?

- Curve fitting! Linear or non-linear.
- Criteria to compare two curves for a dataset.
Errorfunction!

∴ we describe an Errorfunction $E(f, D)$ which returns real no.

Examples for some E:-

$$\sum_D f(x_i) - y_i \times$$

$$L_1 \xrightarrow{\text{norm}} \sum_D |f(x_i) - y_i| \text{ okay; not differentiable}$$

$$L_2^2 \xleftarrow{\text{norm}} \sum_D (f(x_i) - y_i)^2 \text{ okay; differentiable!}$$

cool.

Hence we choose this squared sum as $E(f, D)$ & minimize this.

"method of least squares!"

L2.2: Regression

Linear regression as a canonical example:

- ↳ we got
• regularization → Ridge, Lasso, Bayesian
• non-parametric estimation
• non-linearity through kernels.

* Regression is about learning to predict a set of output vars (dependent) as a function of set of input variables (independent).

Eg: $x_i \rightarrow$ money spent on advertising

$y_i \rightarrow$ sales increased.

then we model $y = \beta_0 + \beta_1 x$.

↳ we find β_0, β_1 by least squares.

then $y^* = \beta_0 + \beta_1 x^*$ ✓ no least squares to

Basic Notation:

* Dataset $D = \langle x_1, y_1 \rangle \langle x_2, y_2 \rangle \dots \langle x_m, y_m \rangle$

$m =$ no. of training samples

x 's = input / training vars

y 's = output / target vars

sales as a non-linear function of investing!

* ϕ_i 's are the attribute/basis functions. & let

$$\Phi = \begin{bmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_p(x_1) \\ \phi_1(x_2) & \dots & \dots & \phi_p(x_2) \\ \vdots & \ddots & \ddots & \vdots \\ \phi_1(x_m) & \dots & \dots & \phi_p(x_m) \end{bmatrix}_{m \times p}$$

↳ basis functions can be different types:
• radial basis
• polynomial basis
• Fourier basis
• wavelet basis.

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}_{m \times 1} \cdot \Phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_p(x) \end{bmatrix}_{p \times 1}$$

↳ not linear in x .

* imp

not linear. hmmm...

General regression problem:-

- * Determine f^* such that $f^*(x)$ is the best predictor of y wrt D .

$$f^* = \underset{f \in F}{\operatorname{argmin}} E(f, D)$$

$f \in F$

set of all possible

functions over which error minimized.
(Linear/nonlinear)

Parameterized Regression problems not linear.

- * Determine parameter w for function $f(\phi(x), w)$ which minimizes $E(f(\phi(x), w), D)$.

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{bmatrix}$$

$$w^* = \underset{w}{\operatorname{argmin}} E(f(\phi(x), w), D)$$

$$\phi(x) = \begin{bmatrix} \phi_1(x) \\ \phi_2(x) \\ \vdots \\ \phi_p(x) \end{bmatrix}$$

Types of regression:-

- * classified based on function class and error function.

heat!

- * Linear regression $\Rightarrow F$ is space of linear functions.

$$f(\phi(x), w) = w^T \phi(x) + b$$

problem is

$$w^* = \underset{w}{\operatorname{argmin}} E(\phi(x), w, D)$$

we introduce

$$\tilde{w} = \begin{bmatrix} w \\ b \end{bmatrix}$$

- * Ridge regression: $\rightsquigarrow E(\cdot) = \sum (y_i - \tilde{w}^T \phi(x_i))^2 + \lambda \|\tilde{w}\|_2^2$

$$\tilde{\phi} = \begin{bmatrix} \phi \\ 1 \end{bmatrix}$$

A shrinkage parameter is added to error function to reduce discrepancies due to variance.

Logistic regression:-

models conditional probability of output variable, given independent var. & used extensively in classification tasks.

$$f(\phi(x), w) = \log \frac{\Pr(y|x)}{1 - \Pr(y|x)} = b + w^T \phi(x)$$

odds

y is $\{0, 1\}$

but

$\log \frac{\Pr(y|x)}{1 - \Pr(y|x)}$ is real valued.

more: lasso regression $E = \text{sq. sum.} + \lambda \|w\|_1$
Stepwise regression
etc.

Least Squares Solution:

$$E(f, D) = \sum_{j=1}^m (f(x_j) - y_j)^2$$

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{j=1}^m \left(\left[\sum_{i=1}^p w_i \phi_i(x_j) \right] - y_j \right)^2$$

→ case-i)

* if minimum value of squared loss is 0; $\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_p \end{pmatrix}$

then $\underline{\underline{\phi(x_u) \cdot w^* = y_u}}$ or equivalently $\underline{\underline{\Phi \cdot w^* = y}}$

* $\Phi w^* = y$ has a solution, if

y is in the column space of Φ $\begin{pmatrix} \phi_1(x_1) & \phi_2(x_1) & \cdots & \phi_p(x_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_1(x_m) & \phi_2(x_m) & \cdots & \phi_p(x_m) \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$

→ case-ii)

* if 0 loss were not attainable?

Geometric interpretation:

• let y^* be a solution in column space of Φ .

• then we need y^* closest to y .

• Foot of perpendicular! you!

∴ $(y - y^*)^T \cdot \Phi = 0$, since $y - y^*$ is perpendicular to every column vector.

$$(y^*)^T \cdot \Phi = y^T \cdot \Phi$$

$$(\Phi \cdot w^*)^T \cdot \Phi = y^T \cdot \Phi$$

$$(w^*)^T \cdot \Phi^T \cdot \Phi = y^T \cdot \Phi$$

$$\Phi^T \cdot \Phi \cdot w^* = \Phi^T \cdot y$$

$$w^* = (\Phi^T \cdot \Phi)^{-1} \cdot \Phi^T \cdot y$$

→ some notion of pseudo inverse.

** $\Phi^T \cdot \Phi$ is invertible iff Φ has full column rank.

PTO for proof:

Proof:-

- 1. say Φ doesn't have full column rank.

$$\Rightarrow \exists x \neq 0 \text{ s.t. } \Phi x = 0$$

$$\Rightarrow \Phi^T \Phi x = 0$$

1. $\Phi^T \Phi$ is not invertible.

- 2. say $\Phi^T \Phi$ is not invertible

$$\Rightarrow \exists x \neq 0 \text{ s.t. } \Phi^T \Phi x = 0$$

$$\Rightarrow x^T \Phi^T \Phi x = 0$$

$$\Rightarrow (\Phi x)^T (\Phi x) = 0$$

$$\Rightarrow \Phi x = 0 \text{ with } x \neq 0$$

$\Rightarrow \Phi$ doesn't have full column rank.

* Therefore

"least square solution y^* is the orthogonal project of y onto the column space of Φ ".

In linear regression.

→ How about analytical derivation?

- More generally, how to minimize a function

- level curves & surfaces

- Gradient vector

- Directional derivative

- Hyperplane

- Tangential hyperplane

- Gradient descent algorithm.

L3.1.1- Linear regression, probabilistic interpretation & regularization:-

- * Owing to basis functions ϕ ; "linear regression" is linear in w but not in x .

↳ radial basis

fourier basis

wavelet basis

domain specific basis

like $\phi_1(x) = 3x^2$

is perfectly

alright.

$\phi \rightarrow$ polynomial basis.

we found $w^* = \underset{w}{\operatorname{argmin}} E(f(w, \phi), D)$.

$$w^* = (\phi^T \phi)^{-1} \Phi^T y$$

→ probabilistic interpretation:- of ordinary least squares :-

- * Gaussian error; ML estimation; addressing overfitting.

→ linear models :-

- * y is a linear function of $\phi(x)$ subjected to a random noise variable ε which we believe is 'mostly bounded' by threshold σ .
modelled output

$$\underbrace{y}_{\text{expected output}} = \underbrace{w^T \phi(x)}_{\text{modelled output}} + \varepsilon$$

$$\varepsilon \sim N(0, \sigma^2)$$

$$\therefore Y \sim N(w^T \phi(x), \sigma^2).$$

- * Now we need w ; that most likely generated the data. \hat{w}_{ML}

$$\Pr(D|w) = \Pr(y|x, w)$$

$$= \frac{m}{\pi} \prod_{j=1}^m \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y_j - w^T \phi(x_j))^2}{2\sigma^2}}$$

$$\hat{w}_{NL} = \underset{w}{\operatorname{argmax}} \Pr(D|w) = L(w|D)$$

$$\hat{w}_{ML} = \underset{w}{\operatorname{argmax}} LL(w|D)$$

$$= \underset{w}{\operatorname{argmin}} \sum_{j=1}^m (y_j - w^T \phi(x_j))^2$$

motivation:-

- * $N(\mu, \sigma^2)$ has maximum entropy among all real-valued distributions with specific variance σ^2 .

3σ rule:-

1σ away - 68%.

2σ away - 95%.

3σ away - 99.7%.

"we want a bounded; but as 'free' as possible distribution" for error.

Same least squares solution.

→ Redundant ϕ and overfitting:-

lets consider polynomial basis $\phi_i(x) = x^i$. with $\phi_1, \phi_2, \dots, \phi_t$.

- * higher t can represent lower t' .

But as $t \uparrow$, we try to fit every data point.

RMS error for train set decreases but test set increases.

"over fitting" → address via bayesian, or via regularization.

(add $\|w\|$ in error term) !!

- * Empirically we observe; too many bends in curve

= high value of w_i 's.

∴ $\|w\|$ increases with increase in t .

- * we want to address overfitting by restricting $\|w\|$.

add $\|w\|$ in error term. → regularization! suppress $\|w\|$.

By using Bayesian Linear Regression:-

- reasoning with some prior belief ($\|w\|$ is restricted in prior).

use posterior over w as result.

we could find $\hat{w} = \underset{\substack{\text{MAP} \\ w}}{\operatorname{argmax}} \Pr(w|D)$.

\max_w a posterior probability

each w_i bounded by

Intuitive prior:

$$\pm \frac{3}{\sqrt{\lambda}} \text{ (99%)}$$

$\Pr(w) \text{ s.t. } \|w\| \leq \frac{1}{\lambda}$. Gaussian again?

well have $w = [w_1 \ w_2 \ \dots \ w_t]$ where $w_i \sim N(0, \frac{1}{\lambda})$ Normal

λ is called precision of the distribution.

- * lets digress for bayesian estimation:-

parameter θ

data (D) = D

$$P(\theta|D) = \frac{P(D|\theta) \cdot P(\theta)}{\int P(D|\theta) \cdot P(\theta) d\theta}$$

posterior prior

λ is precision for w , in prior
 σ^2 is error variance in data.

- * conjugate prior means, if $P(\pi|\theta)$ has d_1 dist. & $P(\theta)$ has d_2 & $P(\theta|D)$ is also d_2 distribution; then d_2 is conjugate prior for d_1 .

Eg: bernoulli r.v has Beta -distribution as conjugate prior

$$P(x|\theta) = \theta^x (1-\theta)^{1-x}$$

$$\therefore P(D|\theta) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i}$$

$$\text{Beta}(p; \alpha, \beta) \sim \frac{p^{\alpha-1} (1-p)^{\beta-1}}{B(\alpha, \beta)} \quad \text{where } B(\alpha, \beta) = \int_{p=0}^1 p^{\alpha-1} (1-p)^{\beta-1} dp$$

$$= \frac{\Gamma(\alpha) \cdot \Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

$$\therefore P(\theta|D) = \frac{\theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \cdot \theta^{\alpha-1} (1-\theta)^{\beta-1}}{\text{normalizing const.}}$$

this is also Beta!

$$P(\theta|D) = \text{Beta}(\theta; \alpha + \sum x_i, \beta + n - \sum x_i)$$

Note:- for a beta distribution $\text{Beta}(p; \alpha, \beta)$

$$E[p] = \frac{\alpha}{\alpha + \beta} \quad \text{model}(p) = \frac{\alpha-1}{\alpha+\beta-2}$$

* $\text{Beta}(p; 1, 1)$ is uniform distribution.

→ conjugate priors for gaussian dist:-

* say $P(x|\mu) \propto e^{-\frac{(x-\mu)^2}{2\sigma^2}}$; σ is a known value.

then conjugate prior is also a gaussian!

$$P(\mu) \propto e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}$$

* now $P(\mu|D) \propto P(D|\mu) \cdot P(\mu)$

"product of gaussians is a gaussian!"

$$\text{say } P(\mu|D) \propto e^{-\frac{(\mu-\bar{x})^2}{2\sigma_m^2}}$$

then we'll get

$$\frac{1}{\sigma_m^2} = m \cdot \frac{1}{\sigma^2} + \frac{1}{\sigma_0^2} \quad (\text{precisions are added})$$

MAP is \bar{x}

$$\frac{1}{\sigma_m^2} \cdot \bar{x} = m \cdot \frac{\hat{\mu}_{MLE}}{\sigma^2} + \frac{\mu_0}{\sigma_0^2}$$

$$\text{where } \hat{\mu}_{MLE} = \text{mean}(x_1, x_2, \dots, x_n)$$

* we'll see multi variat gaussian next.

* Multivariate Gaussian is

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

old friend

a vector of values.

Σ is $n \times n$ covariance matrix & positive-definite. $\mu \in \mathbb{R}^n$.

• for our case

$$\begin{aligned}\hat{\mu}_{MLE} &= \frac{1}{m} \sum_{i=1}^m \phi(x_i) \\ \hat{\Sigma}_{MLE} &= \frac{1}{m} \sum_{i=1}^m (\phi(x_i) - \hat{\mu}_{MLE})(\phi(x_i) - \hat{\mu}_{MLE})^T\end{aligned}$$

outer product. } whatever!!

* Bayesian multivar Gaussian:

$$\mathcal{N}(x; \mu, \Sigma) \propto e^{-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)}$$

$$\& \mathcal{N}(\mu; \mu_0, \Sigma_0) \propto e^{-\frac{1}{2} (\mu - \mu_0)^T \Sigma_0^{-1} (\mu - \mu_0)}$$

then $P(\mu | D)$ is also a multivar Gaussian with μ_m, Σ_m .

$$\begin{aligned}\hat{\Sigma}_m &= m \cdot \hat{\Sigma} + \hat{\Sigma}_0 \\ \hat{\Sigma}_m \mu_m &= m \cdot \hat{\Sigma} \hat{\mu}_m + \hat{\Sigma}_0 \mu_0\end{aligned}$$

multivariate
covariance resulting from
correspondence of $\hat{\Sigma} = \frac{1}{m} \Sigma$

→ coming back to our regressions:

$$w_i \sim \mathcal{N}(0, \frac{1}{\lambda})$$

$$y_i = w^T \phi(x_i) + \epsilon$$

$$\therefore \text{for } w: \mu_0 = 0; \Sigma_0 = \frac{1}{\lambda} I$$

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\therefore \hat{\mu}_{MLE} = \frac{1}{m} \sum \phi(x_i)$$

precision error term.

we'll get

for me;
not clear how.

$$\left\{ \begin{array}{l} \hat{\Sigma}_m = \hat{\Sigma}_0 + \frac{1}{\sigma^2} \Phi^T \Phi \\ \hat{\Sigma}_m \mu_m = \hat{\Sigma}_0 \mu_0 + \frac{\Phi^T y}{\sigma^2} \end{array} \right.$$

$$\left\{ \begin{array}{l} \hat{\Sigma}_m = \lambda I + \frac{\Phi^T \Phi}{\sigma^2} \\ \hat{\Sigma}_m \mu_m = 0 + \frac{\Phi^T y}{\sigma^2} \end{array} \right.$$

$$y_i = w^T \phi(x_i) + \eta_i$$

$$P(y_i | w) = e^{-\frac{(y_i - w^T \phi(x_i))^2}{2\sigma^2}}$$

ridge regression!

$$\sum_{i=1}^n \frac{(y_i - w^T \phi(x_i))^2}{2\sigma^2} + \frac{\lambda}{2} \|w\|_2^2$$

$$= (w - \mu_{MLE})^T \Sigma (w - \mu_{MLE})$$

Since

$$(y_i - w^T \phi(x_i))^2 \rightarrow (y_i - w^T \phi(x_i))^2$$

$$(\Sigma)^{-1} (\Sigma^T \Sigma) = \Sigma^{-1}$$

from MAP part

DE: 23.09.2021

* hence for multivar, bayesian linear regression:

- $P(w|D)$ is a multivar gaussian $\mathcal{N}(\mu_m, \Sigma_m)$ such that

$$\begin{aligned}\Sigma_m &= \lambda I + \frac{1}{\sigma^2} \Phi^T \Phi \\ \mu_m &= (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T Y\end{aligned}$$

(MAP)

* Here MAP estimate = mode $[P(w|D)] = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T Y$

Bayesian estimate = mean $[P(w|D)] = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T Y$

Hence

method	estimate	$P(x D)$
MLE	$\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} LL(D \theta)$	$P(x \theta_{MLE})$
Bayes	$\hat{\theta}_B = E_{P(\theta D)}[\theta]$	$P(x \theta_B)$
MAP	$\hat{\theta}_{MAP} = \underset{\theta}{\operatorname{argmax}} p(\theta D)$	$P(x \theta_{MAP})$
Pure Bayesian	θ is a distribution, rather than single value.	$P(x D) = \int_{-\infty}^{\infty} P(x \theta) \cdot P(\theta D) d\theta$ a pdf for x is image. wow such a pain...

→ Pure Bayesian Regression Summarized:-

$$* P(y|x,D) = \int_w \frac{Pr(y|w,x)}{\mathcal{N}(w^T x, \sigma^2)} \cdot \frac{Pr(w|D)}{\mathcal{N}(\mu_m, \Sigma_m)} dw \quad w \text{ is a vector here.}$$

$$= \mathcal{N}(\mu_m^T \phi(x), \sigma^2 + \Phi^T \Sigma_m \Phi(x))$$

pure bayesian result.
Here mode, mean is $\mu_m^T \phi(x)$

Here; $y = w^T \phi(x) + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$

$w \sim \mathcal{N}(0, \frac{1}{\lambda} I)$ • Prior on w .

$w|D \sim \mathcal{N}(\mu_m, \Sigma_m), \mu_m = (\lambda \sigma^2 I + \Phi^T \Phi)^{-1} \Phi^T y$

hmm...
this we do!
anyways

→ Regularised least squares Regression:-

- * Regularization is done to prevent overfitting & keep the model relevant for test samples.
 - In overfitting, w 's. coefficients become large. Hence include them in loss term.
 - * we got two forms
 - penalized
 - constrained
 } one form, can be converted to another.
 - a) Penalized linear reg. :-
 - * generally, we take
- $$w_{reg} = \underset{w}{\operatorname{argmin}} \left[\| \Phi w - y \|_2^2 + \lambda \Omega(w) \right]$$
- if $\Omega(w) = \|w\|_2^2$: Ridge regression } coincides with Bayesian Lin. Reg.
- if $\Omega(w) = \|w\|_1$: Lasso regression } coincides with Gaussian Prior
- if $\Omega(w) = \|w\|_0$: Support based penalty } Laplacian prior.

- * Some $\Omega(w)$ correspond to priors that are inclosed form.

Some give good working solutions.

Some norms are easier to handle.

b) Constrained Linear Regression:

* generally

$$w_{\text{reg}} = \underset{w}{\operatorname{argmin}} \| \Phi w - y \|_2^2 \text{ such that } g(w) \leq 0.$$

* claim: Any penalized formulation with a ' λ ' can be equivalently converted into constrained formulation with a ' λ '.

Proof requires tools of optimization / duality.

* complex models lead to overfitting & has high ||w|| value.

Increasing the penalizing coefficient " λ " makes the model more linear. eigen values of $(\lambda I + \Phi^T \Phi)$ are indicative of curvature of model. more λ means less curved models.

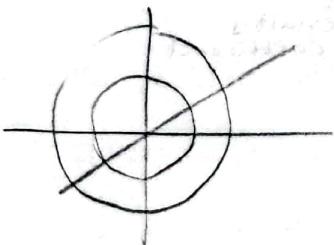
Ridge

$$* w_{\text{reg}} = \underset{w}{\operatorname{argmin}} \| \Phi w - y \|_2^2 + \lambda \| w \|_2^2 = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T y$$

closed form!

obtain by differentiating.

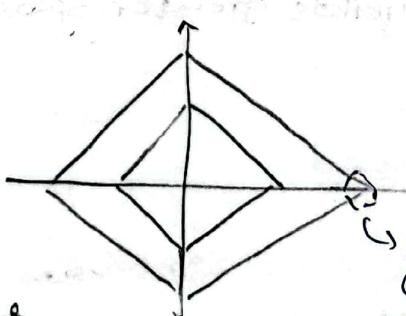
* $\|w\|_2^2$ ridge



level curves

{
we won't get
sparse w's.

$\|w\|_1$ lasso



↳ corners are more prominent.
(coefficients of w are 0)
we can get sparse w's.

Hence; we can find out irrelevant features.!

(whose $w_i = 0$).

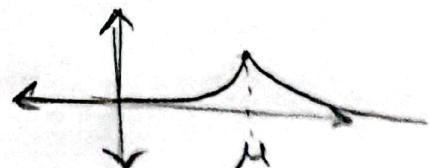
Lasso

$$* w_{\text{reg}} = \underset{w}{\operatorname{argmin}} \| \Phi w - y \|_2^2 + \lambda \| w \|_1 \quad \{ \text{closed form solution not there} \}$$

→ Here Laplacian prior on w_i 's is used.

$$\text{Laplace}(w | \mu, b) = \frac{1}{2b} \cdot \exp\left(-\frac{|w - \mu|}{b}\right)$$

to get Lasso.



→ Lasso: ISTA: iterative Soft thresholding Algorithm :-

* Gaussian → quick closed form, less sparse

Lasso Reg. → slow iterative, more sparse

* lets see the algorithm:-

$$w_{\text{Lasso}} = \underset{w}{\operatorname{argmin}} [E_{\text{LS}}(w) + \lambda \|w\|_1]$$

while drop in $E_{\text{Lasso}}(w^t)$ across $t=k, k+1$ is significant:

1. LS iterate: $w_{\text{LS}}^{K+1} = w_{\text{Lasso}}^K - \eta \nabla E_{\text{LS}}(w_{\text{Lasso}}^K)$

2. Proximal Step:

• for every $i \in [p]$

$$[w_{\text{Lasso}}^{K+1}]_i = \begin{cases} [w_{\text{LS}}^{K+1}]_i - \lambda \eta & \text{if } [w_{\text{LS}}^{K+1}]_i > \lambda \eta \\ [w_{\text{LS}}^{K+1}]_i + \lambda \eta & \text{if } [w_{\text{LS}}^{K+1}]_i < -\lambda \eta \\ 0 & \text{else} \end{cases}$$

soft thresholding

* This soft thresholding yields greater sparsity as $\lambda \uparrow$.

λ = penalty coefficient.

By now, we discussed

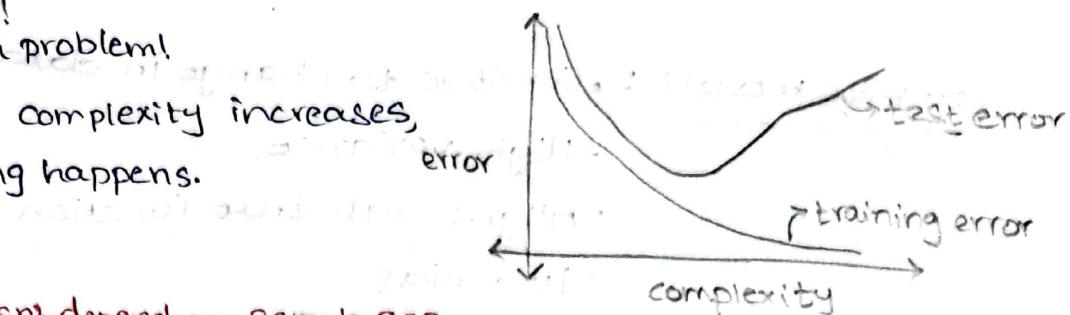
- + OLS, ordinary least squares.
- + MLE interpretation of OLS
- + Bayesian linear regression & its closed form
- + Ridge linear reg. & its closed form
- + Lasso linear regression & its ISTA approach.

L6: Understanding Bias, Variance:-

underfitting too!

* overfitting is a problem!

when model complexity increases,
overfitting happens.

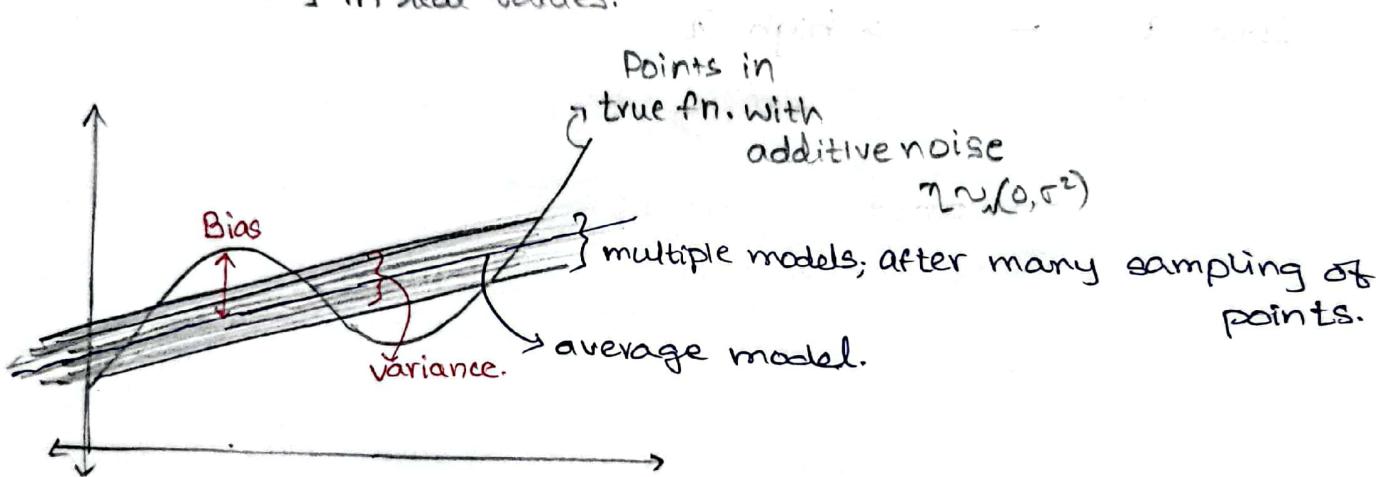


doesn't depend on sample size.

* The three sources of error are

- Bias } b/w avg model & real values
- Variance } b/w avg model & every model
- Noise } in real values.

*



* Say true function is $g(x)$. So, $y = \underbrace{g(x)}_{\text{true fn.}} + \underbrace{\epsilon}_{\text{noise}}$ in sampling.

$$\epsilon \sim N(0, \sigma^2)$$

- let $f_D(x) = w^T x$ be our model.

- given a point \hat{x} , we got $\hat{y} = g(\hat{x}) + \epsilon$ in measurement & $f_D(\hat{x})$ via model.

$$\therefore \text{Err}(\hat{x}) = \underbrace{\mathbb{E}_{D, \epsilon}[(\hat{y} - f_D(\hat{x}))^2]}_{\text{average over D & noise } \epsilon}$$

we manage to show

$$\begin{aligned} * \text{Err}(\hat{x}) &= \underbrace{\mathbb{E}[(f_D(\hat{x}) - \hat{f}_D(\hat{x}))^2]}_{\text{variance}} \\ &\quad + \underbrace{(\hat{g}(\hat{x}) - \hat{f}_D(\hat{x}))^2}_{\text{Bias}^2} \\ &\quad + \underbrace{\sigma^2}_{\text{noise}} \end{aligned}$$

$$\begin{aligned} &= \mathbb{E}[(f_D(\hat{x}) - \hat{f}_D(\hat{x}))^2] + (\mathbb{E}[f_D(\hat{x})])^2 \\ &= \text{Variance} + \mathbb{E}(\epsilon^2) + \mathbb{E}(g(\hat{x}))^2 \\ &\downarrow \\ &= \mathbb{E}[2\hat{g}(\hat{x})f_D(\hat{x})] - \mathbb{E}[2\mathbb{E}f_D(\hat{x})] \\ &\quad \text{proof doable.} \end{aligned}$$

* Here we note:

Complex models: • Sensitive to change in data
- High variance

- Aligns with true function
- Low bias

Simple models: • Low variance
• Can't Align. Has straight line shape
- High bias

* low λ \longrightarrow high λ
complex simple

\longrightarrow
variance decreases
Bias increases.

L7,8 - Linear models for classification. Perceptrons

not for predicting a continuous value!
for classifying!

* we are given data x_1, x_2, \dots, x_m .

& possible class choices C_1, C_2, \dots, C_k .

we need to model a mapping/classifier

$$f: X \rightarrow \{C_1, C_2, C_3, \dots, C_k\}$$

that predicts class labels for each data.

→ we could try

$$1) y = i \text{ for } C_i \in \text{apply usual linear reg.}$$

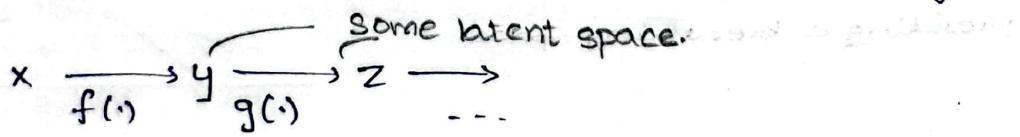
Not we are imposing unnecessary ordering among classes.

$$2) y = [1 \ 0 \ 0]^T \text{ for } C_1 \quad \& \text{ apply usual linear regression.}$$

All are
equidistantly
no ordering

imposed. NO! calls for more machinery.

* instead of single mapping, we'll try for a series of mappings



* this can decrease bias as complexity increases.

→ Perceptron classifier

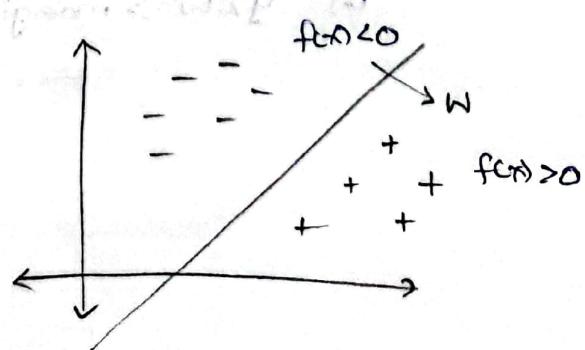
* lets consider a binary classification problem:

$f(x) \in \{+1, -1\}$ this ordering doesn't mean much.

our $f(x) = w^T \Phi(x) + b$

if $f(x) < 0 \rightarrow y = -1$ prediction

if $f(x) > 0 \rightarrow y = +1$ prediction



* How do we find w ?

we iteratively update it

$$w^{(k+1)} = w^{(k)} + \frac{\text{?}}{?}$$

update rule

$= 0$ when point is rightly classified.

$= \eta y \phi(x)$ when otherwise.

!

* for a rightly classified point $y \cdot w^T \phi(x) > 0$.

if $y \cdot w_{(k)}^T \phi(x) < 0$; then unsigned distance.

$$w_{(k+1)} = w_{(k)} + \eta y \phi(x) \quad \text{update rule of perceptron.}$$

bcoz, now

$$y w_{(k+1)}^T \phi(x) = y w_{(k)}^T \phi(x) + \eta y^2 \phi(x) \cdot \phi(x)$$

always ≥ 0

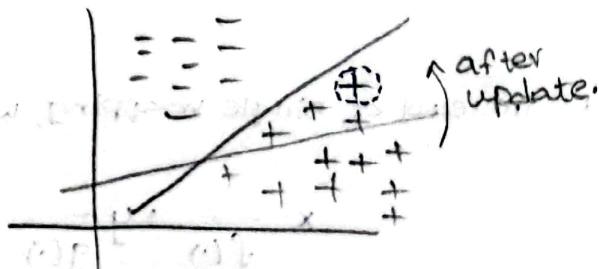
\therefore we trying to make unsigned distance positive.

for every datapoint x_i

* This perceptron does not find the best

separating hyperplane. It finds any.

- SVMs address this problem by providing a breathing space.



* Perceptron can model a linearly separable function.

* if $\exists w^*$ s.t. points are all satisfying $y w^* x > 0$ then wlog we can assume

a) $\|w^*\|_2 = 1$ or $\|x\|_2 < 1$ } use a scaling transformation

b) $y x^T w^* > \text{margin}$.

Sure! our set is finite.

→ convergence of perceptron

- * lets say the given data is linearly separable.
- Hence u be the unit weight vector.
- Also take that $\|x\|_2 < 1$ for all datums.
- let $\gamma = \min_{x \in D} |u^T x|$ is the min. distance of $x \in D$ from 'u' line.

* lets watch $\|w_{(k)}^T u\|$ & $\|w_{(k)}\|^2$:

$$a) \underline{w_{(k+1)} = w_{(k)} + (\phi(x))u}$$

take $\gamma = 1$ & $y = 1$

$$\|w_{(k+1)}\|^2 = \|w_k\|^2 + \|x\|^2 + 2w_k^T x.$$

~~so if $w_k^T x < 0$, then $\|w_{(k+1)}\|^2 - \|w_k\|^2 > 1$~~ hence we are updating

$$\therefore \underline{\|w_{(k+1)}\|^2 - \|w_k\|^2 < 1}$$

$$b) \underline{u^T w_{(k+1)} = u^T w_k + u^T x} \geq 0$$

$$2. \underline{u^T w_{(k+1)} - u^T w_k > \gamma}$$

* Say after K updates, $\|w_k\| = 0$

$$\|w_k\|^2 \leq K$$

$$u^T w_k > K\gamma$$

$$\& \text{we know that } \|w_k\|^2 \geq |u^T w_k|^2$$

$$\therefore K > K\gamma^2$$

$$K < \frac{1}{\gamma^2} \quad \text{Hence, we've shown}$$

PROOF & RATE OF CONVERGENCE.

if we had $\|x\|_2 < t$

then

$$K < \frac{t^2}{\gamma^2} \quad (\text{midsem question})$$

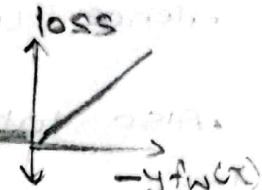
γ definition still same.

→ What are we minimising? Hinge loss:-

* $L_p(f_w(x), y) = \max(0, -yf_w(x))$ for each x_i

lets apply Stochastic gradient descent (SGD):
update for each datapoint

$$w \leftarrow w - \nabla L_p(f_w(x), y)$$



if $yf_w(x) > 0$; $\nabla L_p(\cdot, \cdot) = 0$

if $yf_w(x) < 0$; $\nabla L_p(f_w(x), y) = -y$

$$\therefore w \leftarrow w + yx \quad (\text{same as our update rule})$$

* Always the final perceptron isn't a good idea. overfitting

o Voted perception:

votes = no. of samples correctly classified
wt-ed avg!

o Average perception:

Average different values of w at the end of algorithm.

* Batch stochastic gradient descent:

$$w \leftarrow w - \sum_{(x,y) \in B} \nabla L_p(f_w(x), y)$$

// if we do batch wise,
each batch ^{entry} should use same
old 'w' to calculate $L_p(f_w(x), y)$

if $B = \text{complete dataset}$, we have regular gradient descent.

* SGD → each point (x_i, y_i) should have different updated 'w'.

gradient descent → All points gradient is computed using single, old 'w' value

& 'w' is updated.

Next iteration uses this updated 'w'.

Lecture 9 - Logistic Regression:-

* till now, $f_w(x) = \text{Sign}(w^T \phi(x))$ signum activation.

Now, $f_w(x) = \sigma(w^T \phi(x))$ for activation function σ

σ is sigmoid function \rightarrow Logistic regression.

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

* Interpret $\sigma_w(x)$ as $P(y=1|x)$

which our model predicts.

* $L(D, w) = \prod_{i=1}^m P(y=y_i | x_i, w)$

Aim: find a 'w' such that
our model's predictions are
most accurate.
OR MLE of w.

$$w = \underset{w}{\operatorname{argmax}} \prod_{i=1}^m (\sigma(w^T \phi(x_i)))^{y_i} (1 - \sigma(w^T \phi(x_i)))^{1-y_i}$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} -\log(L(D, w))$$

$$\hat{w} = \underset{w}{\operatorname{argmin}} \left\{ \underbrace{-y_i \log(\sigma(w^T \phi(x_i)))}_{\substack{\text{data} \\ \text{to} \\ \text{code}}} + \underbrace{(1-y_i) \log(1-\sigma(w^T \phi(x_i)))}_{\substack{\text{coding} \\ \text{length}}} \right\}$$

We are minimizing coding required for the data.

* This is the Cross Entropy loss function :-

$$E(w) = \sum \frac{1}{M} \left[y_i \log(\sigma_w(x_i)) + (1-y_i) \log(1-\sigma_w(x_i)) \right]$$

↓ After simplification :-

$$E(w) = - \left[\frac{1}{m} \sum_{i=1}^m \left(y_i w^T \phi(x_i) - \log(1+e^{w^T \phi(x_i)}) \right) \right]$$

unsigned
loss
of normal
perceptron

- * NO closed form for w . do gradient descent:

we get this, rather than hinge loss,
since sigmoid
is smooth

$$\nabla(-\eta E(w)) = \eta \left[\frac{1}{m} \sum_{i=1}^m (y_i - \sigma_w(x_i)) \phi(x_i) \right]$$

old term
- * $w \leftarrow w + \eta \left[\frac{1}{m} \sum_{i=1}^m (y_i \phi(x_i) - \sigma_w(x_i) \phi(x_i)) \right]$

if $y_i = 1$
 $\sigma_w(x_i) \rightarrow 1$
then no update!
good!
meaningful!

Stochastic version:-

$$w \leftarrow w + \eta y_i \phi(x_i) - \eta \sigma_w(x_i) \phi(x_i)$$

Here, sigmoid is smooth. So good exp.
Signum is step function.

→ Regularized LR:-

$$* E(w) = -\frac{1}{m} \sum_{i=1}^m [y_i \log(\sigma_w(x_i)) + (1-y_i) \log(1-\sigma_w(x_i))] + \frac{\lambda}{2m} \|w\|_2^2$$

regularized w .

Probabilistic reason:- A bayesian posterior explanation.
Gaussian prior $w_j \sim N(0, \frac{1}{\lambda})$

→ Multiclass LR:-

- * K classes; each class has different weight vector

$w_1, w_2, w_3, \dots, w_{K-1}$ (one class don't need)

$$P(Y=c | \phi(x)) = \frac{e^{-w_c^T \phi(x)}}{1 + e^{-w_1^T \phi(x)} + e^{-w_2^T \phi(x)} + \dots + e^{-w_{K-1}^T \phi(x)}}$$

$c = 1, 2, \dots, K-1$

$$P(Y=k | \phi(x)) = \frac{1}{1 + \sum_{c \neq k} e^{-w_c^T \phi(x)}}$$

Lecture-10 Kernels:-

- * We observe that, we get a lot of $\phi(x_i)^T \phi(x_j)$ everywhere. So, we introduce:

$$\underline{K(x_i, x_j) = \phi(x_i)^T \phi(x_j)}$$

And, by talking in terms of $K(\cdot, \cdot)$ we can use implicitly large $\phi(x)$ spaces.
even infinite.

→ Kernel perception:-

- * $f_w(x) = \text{Sign} \left(\sum \alpha_i y_i K(x, x_i) + b \right)$

- Initialize : $\alpha_i = 0 \forall i \in [m]$

- Update :

if $f(x) \neq y_i$

then

$\alpha_i \leftarrow \alpha_i + 1$

this is the normal $w \leftarrow w + y_i \phi(x)$.

But kernelised.

Never explicitly compute
 $\phi(x)$ now. It's very costly.

- * Linear Kernel: $K(x, y) = x^T y$

Exponential: $K(x, y) = e^{x^T y}$

Polynomial: $K(x, y) = (1 + x^T y)^d$

(Kinda polynomial)

RBF Kernel: $K(x, y) = e^{-\frac{1}{2} \frac{\|x - y\|^2}{\sigma^2}}$ } Implicit $\phi(x)$ is ∞ .

$K_y(x)$ is Smoothing Kernel.

Eg: What is the implicit $\phi(x)$?

Say $K(x, y) = (1 + x^T y)^2 \quad \& \quad x, y \in \mathbb{R}^2$

$$= (1 + x_1 y_1 + x_2 y_2)^2$$

$$= 1 + x_1^2 y_1^2 + x_2^2 y_2^2 + 2x_1 y_1 + 2x_2 y_2 + 2x_1 x_2 y_1 y_2.$$

Now; $\phi(x)^T \phi(y) =$ ↑

$$\therefore \phi(x) = \begin{bmatrix} 1 \\ x_1^2 \\ x_2^2 \\ \Sigma x_1 \\ \Sigma x_2 \\ \Sigma x_1 x_2 \end{bmatrix}$$

So simple $K(x, y)$ 2D Space

But so huge $\phi(x)$!! 6D Space

So $K(x, y)$ has obvious advantage.

We can get non-linear classifiers !!!

even in CS215??

→ The Gram matrix & Kernel matrix

$$K = \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_m) \\ \vdots & \ddots & \vdots \\ K(x_m, x_1) & \cdots & K(x_m, x_m) \end{bmatrix}$$

if $K(x_i, x_j)$ is valid and
 $= \phi(x_i)^T \cdot \phi(x_j)$

then

K is Symmetric matrix

pos. semidefinite matrix.

* Since K is positive semi definite,
we can write $(SVD?)^H$

$$K = U\Sigma U^T = (U\Sigma^{1/2})(U\Sigma^{1/2})^T$$

$$= R^T R$$

- here rows of U are linearly independent. U is an orthogonal matrix.
- Σ is a positive diagonal matrix.

This rows of R can be seen as feature vectors.

$$\begin{array}{c} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_m) \end{array}$$

→ Mercer's theorem:

defn: $K(x_1, x_2)$ is a mercer kernel, if

$$\underbrace{\int \int K(x_1, x_2) g(x_1) g(x_2) dx_1 dx_2}_{\text{for all square integrable functions } g(x)} \geq 0$$

functions $g(x)$.

$$\text{i.e. } \int_x (g(x))^2 dx < \infty.$$

thm: "For any mercer kernel $K(x_1, x_2)$,

$$\exists \phi(x) : \mathbb{R}^n \mapsto H, \text{ s.t. } K(x_1, x_2) = \phi(x_1)^T \phi(x_2)$$

hilbert space

- H is hilbert space, infinite dimension version of euclidean space.
- Formally, H is a innerproduct Space with associated norm, where every cauchy sequence is convergent.

* prove $\pi_1^T \pi_2$ is a Mercer Kernel :-

$$\iint_{X_1 X_2} x_1^T x_2 g(x_1) \cdot g(x_2) d\mu(x_1) d\mu(x_2)$$

$$= (\int_{X_1} x_1 g(x_1))^T (\int_{X_2} x_2 g(x_2))$$

relabel x_2 as π_1 :

$$= (\int_{X_1} x_1 g(x_1))^T \cdot (\int_{X_1} \pi_1 g(\pi_1) d\mu)$$

$$\geq 0$$

$$\iint_{X_1 X_2} (x_{1a} x_{2a} + x_{1b} x_{2b}) g(x_1) g(x_2)$$

$$= \sum_{x_1, x_2} \iint_{X_1 X_2} x_1 g(x_1) \cdot x_2 g(x_2)$$

$$= \sum_{x_1} \left(\int_{X_1} x_1 g(x_1) d\mu \right)^2$$

$$(\int_{X_1} x_1 g(x_1) d\mu)^2$$

$$\sum_{x_1} \left(\int_{X_1} \pi_1 g(\pi_1) d\mu \right)^2$$

$$\geq 0$$

→ Kernels closure properties:

1) Addition:-

* sum of two kernels is also kernel

* ϕ_{new} is concat (ϕ_1, ϕ_2). ← this is proof.

2) Product:-

* product of two kernels also kernel.

* ϕ_{new} is cross-multiplication of ϕ_1, ϕ_2, \dots

Good!

Lecture-11: Various Kernelisations

→ Kernalized LR:-

* cross entropy loss function was :-

$$E(w) = - \left[\frac{1}{m} \sum_{i=1}^m \underbrace{(y_i w^\top \phi(x_i) - \log(1 + e^{w^\top \phi(x_i)}))}_{\text{cross entropy loss}} \right] + \frac{\lambda}{2m} \|w\|_2^2$$

now, kernalized version is :-

$$E(\alpha) = - \left[\sum_{i=1}^m \left(\underbrace{\sum_{j=1}^m y_j K(x_i, x_j) \alpha_j}_{\text{inner product}} - \frac{\lambda}{2} \alpha_i \alpha_j K(x_i, x_j) \right) \right. \\ \left. - \log(1 + \exp \left[- \sum_{j=1}^m \alpha_j K(x_i, x_j) \right]) \right]$$

easy....

& decision function:

$$\hat{y}_w(x) = \frac{1}{1 + \exp \left[- \sum_{j=1}^m \alpha_j K(x_i, x_j) \right]} // \text{fine}$$

→ Kernalized ridge regression:-

→ now!

* Recall ridge regression:

The bayes and map estimate for Linear regression $y = w^\top \phi(x) + \epsilon$
with gaussian prior $w \sim N(0, \frac{1}{\lambda} I)$ coincide with ridge regression $\underset{w}{\text{min}} \frac{1}{2} \|w\|^2 + \frac{\lambda}{2} \epsilon^2$

$$w_{\text{ridge}} = \underset{w}{\text{arg min}} \left(\|w - y\|_2^2 + \lambda \sigma^2 \|w\|_2^2 \right)$$

with closed form

$$w_{\text{ridge}} = \underbrace{(\Phi^\top \Phi + \lambda \sigma^2 I)^{-1} \Phi^\top y}_{\text{closed form}}$$

So, now the function $f(x)$ will be...

$$f(x) = w^\top \phi(x) = \phi(x)^\top w = \underbrace{\phi(x)^\top (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top y}_{\text{closed form}}$$

NOW.... LOOK for kernalization...

$\Phi^\top \Phi$ is NOT K .

$\Phi^\top \Phi$ is kernel matrix K .

Need to use matrix identity ...

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1} \text{ hahaha...oops}$$

Hence, now $R = I, P = \frac{I}{\lambda}, B = \phi$

$$(\lambda I + \Phi^T \Phi)^{-1} \Phi^T = \frac{1}{\lambda} \Phi^T (\frac{\Phi \Phi^T + I}{\lambda})^{-1}$$

$\therefore w = \Phi^T (\Phi \Phi^T + \lambda I)^{-1} y = \sum_{i=1}^m \alpha_i \phi(x_i), \quad x_i = ((\Phi \Phi^T + \lambda I)^{-1} y)_i$

redge
reg.

so, $f_w(x) = \sum_{i=1}^m \alpha_i \underline{\phi^T(x) \cdot \phi(x_i)}$
Kernel!

$\therefore f_w(x) = \sum_{i=1}^m \alpha_i K(x, x_i)$ where $\alpha_i = ((K + \lambda I)^{-1} y)_i$

Hence, given (x_i, y_i) we can directly use Kernels and calculate a new $f(x)$

there's no w here now.
only α_i 's.

infinitely large
implicit $\phi(x)$

linear + non-linear + infinite dimensional

$$\Phi^T (\Phi \Phi^T + \lambda I)^{-1} \Phi = \text{scaled identity}$$

and more with some scaling factors

linear + non-linear + infinite dimensional

and more with some scaling factors

Lecture-12: Unsupervised learning

PCA and dimensionality reduction

Principle Component Analysis

- * Dimensionality reduction is a feature extraction technique where we want to represent input data with fewer dimensions.

Occam's Razor ? Make Everything Simple.

- * Consider $x \in \mathbb{R}^d$. here, the new axes are eigen vectors of old space.
Project it to a k -dim vector z

$$z = U^T x \quad \text{where } U \in \mathbb{R}^{d \times k} \text{ is projection matrix.}$$

$$\therefore x = U z$$

we want

$U^T U = I$. } like, columns are orthonormal
 $U^T U$ can be anything. eigen vectors.

- What are PCA's objectives?

Perspective-1 :- Minimized sum of squared distances of \hat{x} from x .

$$\hat{x} = U z \quad \text{constructed.}$$

x is given

$$z = U^T x.$$

\therefore we need $U^* = \underset{U}{\operatorname{argmin}} \sum_x \|x - UU^T x\|^2$ such that $U^T U = I$

Perspective-2 :- Maximize variance of projected data.

Assume $E(x) = 0$. (since variance won't change)

$$z = U^T x$$

$$E(\|z\|^2) = \frac{1}{n} \sum_{i=1}^n \|U^T x_i\|^2$$

\therefore we need $U^* = \underset{U, U^T U = I}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \|U^T x_i\|^2$

Both are equivalent

oh lol!

$$\text{since } x = \underbrace{U U^T x}_{A} + \underbrace{(I - U U^T)}_{B} x \quad AB = 0$$

$$\therefore \|x\|^2 = \underbrace{\|U U^T x\|^2}_{E(\|U^T x\|)} + \underbrace{\|x - U U^T x\|^2}_{\text{reconstruction error.}}$$

to minimize

to maximize.

* we are going to work with "maximizing variance" lets do one component

$$\max_{\|u\|=1} E[(u^T x)^2]$$

$$= \max_{\|u\|=1} \frac{1}{n} \|u^T x\|^2 \quad x = \begin{bmatrix} \vdots & \vdots & \vdots \\ x_1 & x_2 & \dots \end{bmatrix}$$

$$= \max_{\|u\|=1} \frac{1}{n} u^T (x x^T) u$$

$$= \max_{\|u\|=1} u^T S u$$

~~this means u is an eigenvector.~~

$$S = \frac{1}{n} (x x^T) \text{ covariance matrix. } \because E(x) = 0.$$

u = eigen vector with largest eigen value.

∴ the $d \times k$ projection matrix $U = \begin{bmatrix} \vdots & \vdots & \vdots \\ e_1 & e_2 & e_3 & \dots \\ \hline \vdots & \vdots & \vdots & \vdots \end{bmatrix}_{d \times k}$

PCA Algorithm

eigenvalues,
in decreasing order.

→ Kernel PCA Algorithm:-

* consider SVD of X . $X = A D B^T$

$$\text{so } nS = X X^T = A D^2 A^T$$

$$X^T X = B D^2 B^T$$

$$X = \begin{bmatrix} \vdots & \vdots & \vdots \\ x_1 & x_2 & \dots & x_m \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

$U = AD$ is the matrix of all principle components.

* Now, after taking the kernel matrix; $K = X^T X = B D^2 B^T$.

the projections XU of data X onto components U can be computed from eigen decomposition of K .

Pending
L12.

Lecture-13: clustering

Unsupervised learning.

* Two paradigms — Bottom-up agglomerative clustering

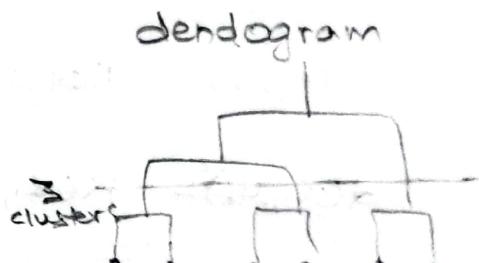
Top-down partitioning "move to nearest"

→ K-means clustering:

Important Step!

1) Initialize $\mu_i^{(0)}$ to random values.

$$2) P_r^{(t+1)} \in \operatorname{argmin}_P \sum_{j=1}^m \sum_{l=1}^K P_{l,j} \| \phi(x_j) - \mu_l^{(t)} \|^2$$



solution:

$$\text{let } i \text{ be } \operatorname{argmin}_l \| \phi(x_j) - \mu_l^{(t)} \|^2$$

$$\text{then } P_{i,j}^{(t+1)} = 1 \text{ and}$$

$$P_{l,j}^{(t+1)} = 0 \quad \forall l \neq i$$

$$3) \mu^{(t+1)} = \operatorname{argmin}_{\mu} \sum_{j=1}^m \sum_{l=1}^K P_{l,j}^{(t+1)} \| \phi(x_j) - \mu_l \|^2$$

solution:

$$\mu_l^{(t+1)} = \frac{\sum_{j=1}^m P_{l,j}^{(t+1)} \cdot \phi(x_j)}{\sum_{j=1}^m P_{l,j}^{(t+1)}}$$

4) if any $P_{l,j}$ is changed, go to step 2.

* K-means will converge, since we have finite labelling for points

Σ we always move to a less error value.

→ Kernel Kmeans

- * Basic idea; replace μ_L by $\phi(\mu_L)$, which is calculated by Pls.

- * we expand $d(x_i, \mu_L)^2 = \|\phi(x_i) - \mu_L\|^2$ as

$$d(x_i, \mu_L)^2 = \|\phi(x_i)\|^2 + \|\mu_L\|^2 - 2 \|\phi(x_i) \cdot \mu_L\|$$

$$d(x_i, \mu_L)^2 = K(x_i, x_i) + K(\mu_L, \mu_L) - 2 K(x_i, \mu_L)$$

and now step-3 won't be there.

in step-2; we calculate

- 1) $K(x_i, \mu_L)$ using

$$\frac{\sum P_{Lj} K(x_i, x_j)}{\sum P_{Lj}}$$

thanks
mukesh!

- 2) $K(\mu_L, \mu_L)$ using

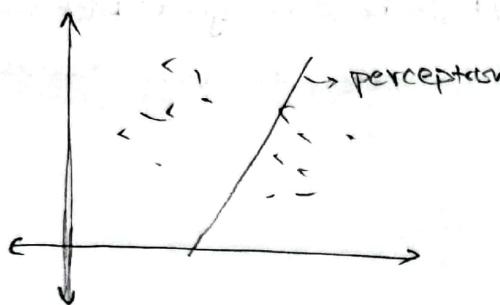
$$\frac{\sum P_{Lj_1} P_{Lj_2} K(x_{j_1}^i, x_{j_2}^i)}{\sum P_{Lj_1} \cdot P_{Lj_2}}$$

- * Note that ϕ doesn't have to be computed/stored for each data instance x_i or for the cluster means μ_L .
- * remember that $\phi(\mu_L)$ is mean of data in feature space, which might not have a pre-image in the dataspace.

Lecture-14: SVMs

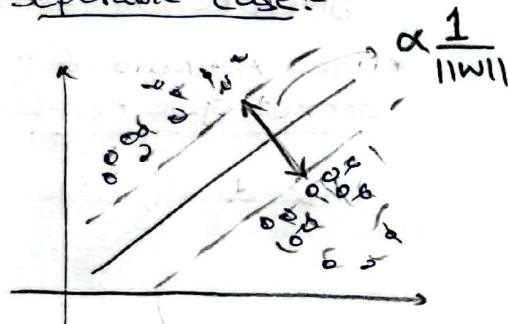
support vector machines.

- * The perception discussed previously doesn't find the Best hyperplane. To address this, we use SVMs.



→ support vector classification:

1) separable case:

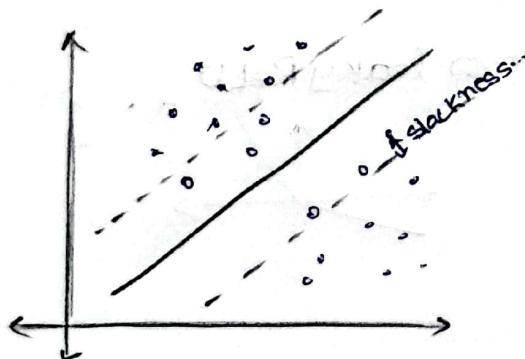


$$y(\mathbf{w}^T \phi(\mathbf{x}) + b) \geq +1$$

margin!

$$\alpha \frac{1}{\|\mathbf{w}\|}$$

2) non-separable case:



when examples are not separable, we need to consider slackness

$$\epsilon_i \geq 0$$

$$y(\mathbf{w}^T \phi(\mathbf{x}) + b) \geq 1 - \epsilon_i$$

Objective:

* i) Maximize the margin

$$= \frac{2}{\|\mathbf{w}\|_2}$$

2) Minimize the slackness

$$= \sum \epsilon_i$$

this min. will maximize margin

3. final objective:

$$(\mathbf{w}^*, b^*) = \underset{\mathbf{w}, b, \epsilon_i}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \epsilon_i$$

such that

$$y_i (\mathbf{w}^T \phi(\mathbf{x}_i) + b) \geq 1 - \epsilon_i \quad \text{and} \quad \epsilon_i \geq 0$$

✓ good!

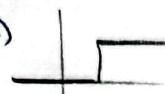
Lecture 15 - Intro to neural networks:-

* Neural networks:- Universal approximation properties

- with a single hidden layer of sufficient size and reasonable choice of non-linearity: Sigmoid/hyperbolic tangent/RBF unit; one can represent any smooth function to any desired accuracy.

• Activation functions

1) Step : $g(s) = 1 \text{ if } s \in [0, \infty)$
 $= 0 \text{ else}$



* neural networks is cascade of perceptions giving non-linearity.

2) Sigmoid: $g(s) = \frac{1}{1+e^{-s}}$



rectified linear unit

3) ReLU: (most popular) $g(s) = \max(0, s)$



* linear Activation fun, cannot give non-linear

* diff. for ReLU \rightarrow Softmax: $g(s) = \ln(1+e^s)$
 Softplus



4) tanh: $g(s) = \frac{2}{1+e^{-2s}} - 1$



cascade of linear is linear.

* softmax

- multi class logistic reg.
- with K weight vectors

$P(Y=i|X) = \frac{w_i^T \phi(x)}{e^{\sum_{i=1}^K w_i^T \phi(x)}}$

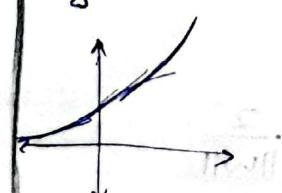
$$\sum_{i=1}^K e^{w_i^T \phi(x)}$$

- one vector is redundant

softplus

Activation function...

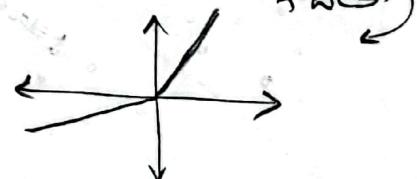
$g(s) = \ln(1+e^s)$



Smooth Version to

ReLU

6) leaky ReLU $f_w(s)$



→ VC-dimension:- Vapnik-Chervonenkis

* measure for separability of a classifier.

* = cardinality of the largest set of points that can be shattered.

a function f_w is said to

shatter a set of points (x_1, x_2, \dots, x_n) if,

for all 2^n assignments of labels to those points,
there exists a 'w' such that f_w makes no errors.

* Hence; if there exist atleast one set of size 'n',

$$\text{then } V.C. \geq n.$$

Examples:

1) For f_w as threshold classifier, $VC = 1$
 in \mathbb{R}

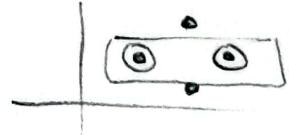
$VC = 4$ for
 parallel rect.
 classifier

2) for f_w as interval classifier, $VC = 2$
 in \mathbb{R}^2

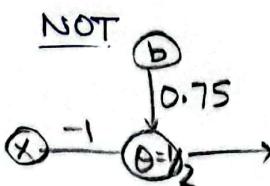
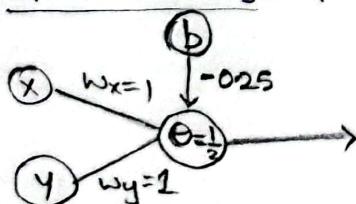
3) for f_w as perceptron in \mathbb{R}^2 , $VC = 3$
 (linear)



linear in \mathbb{R}^n , $VC = n+1$
 perception

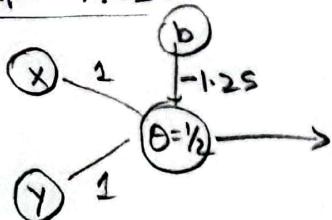


Simple OR using step



Now; whole logic gates are ready.

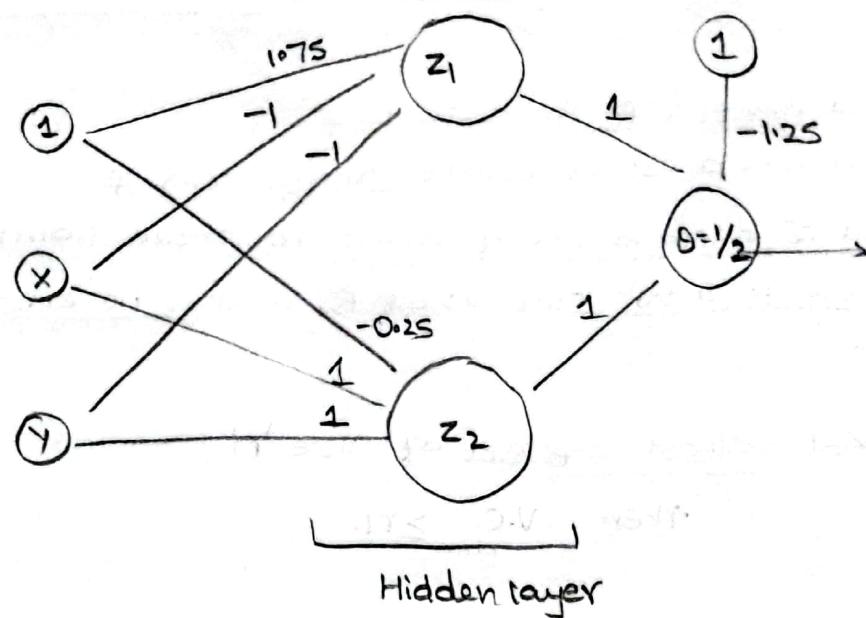
Simple AND



Any logic circuit can be
 built.

Now, we build XOR as

$$(x \oplus y) = (\bar{x} \wedge y) \vee (\bar{y} \wedge x) \quad (x \vee y) \wedge (\bar{x} \vee \bar{y})$$



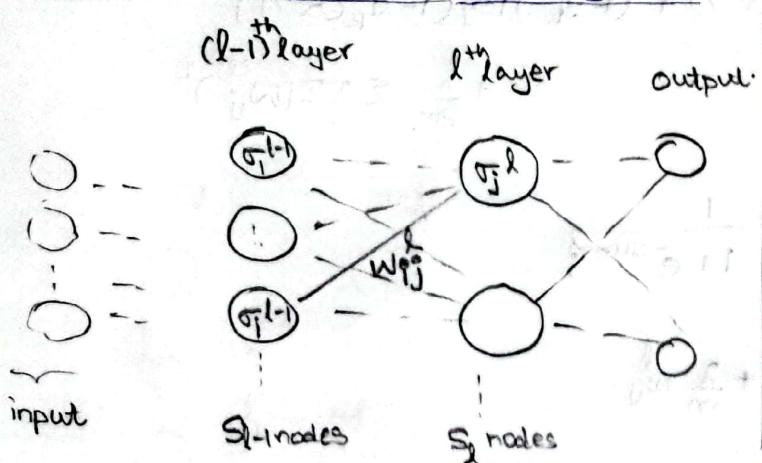
Lecture 16 - Training NNs, Backpropagation Algo

* Picking a network architecture:

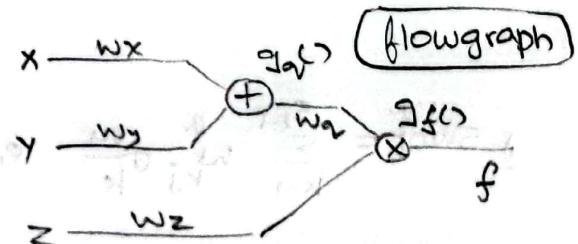
- No. of input units: Dimension of features $\phi(x_i)$
- No. of output units: No. of classes
- Reasonable default: 1 hidden layer or
 > 1 hidden layer having same no. of nodes
- No. of hidden units in each hidden layer is constant factor ($\times 3$ or $\times 4$) of input units.

* We want to apply gradient descent to get optimal weights.
(on each edge).

→ Backpropagation for training NNs:-



$$f(xyz) = g_f(w_x g_g(w_x x + w_y y) + w_z z)$$



$$\text{now, } \frac{\partial f}{\partial w_x} = \underline{\frac{\partial f}{\partial g_f}} \cdot \underline{\frac{\partial g_f}{\partial g_g}} \cdot \underline{\frac{\partial g_g}{\partial w_x}}$$

back prop...

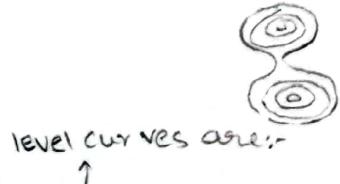
1) Randomly initialize weights w_{ij} for $l=1,2,\dots,L$; $i=1,2,\dots,S_{l-1}$; $j=1,2,\dots,S_l$

2) Implement forward propagation to get $f_w(x)$ for any $x \in D$

3) Execute backward propagation:

a) by computing partial derivatives

$$\frac{\partial E(w)}{\partial w_{ij}^l} \text{ using chain rule}$$



b) and use gradient descent to try to minimize (non-convex) $E(w)$

$$w_{ij}^l = w_{ij}^l - \eta \frac{\partial}{\partial w_{ij}^l} E(w)$$

4) verify $E(w)$ has indeed reduced, else resort to random perturbation of w's to address non-convexity of $E(w)$.

memoize the previous partial derivatives!

* The objective to be minimised is

$$E(w) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\sigma_k^L(x^{(i)})) + (1-y_k^{(i)}) \log(1-\sigma_k^L(x^{(i)}))$$

wrote for
every K
classes.

{ cross entropy

wrt output layer, true!

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{S_{l-1}} \sum_{j=1}^{S_l} (w_{ij}^{(l)})^2$$

* chain rule & memoisation
is bcoz of this part

all weights in

hidden layers too
regularised.

* lets consider sigmoidal activation $g(s) = \frac{1}{1+e^{-s}}$.

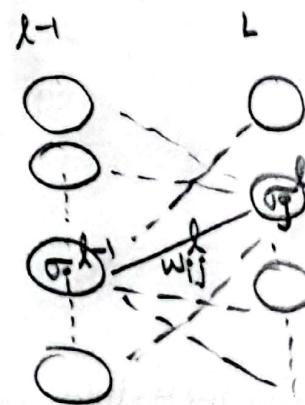
$$\frac{dg(s)}{ds} = g(s)(1-g(s))$$

now, for single point...

$$E(w) = -\sum_{k=1}^K y_k^{(i)} \log(\sigma_k^L(x^{(i)})) + (1-y_k^{(i)}) \log(1-\sigma_k^L(x^{(i)}))$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^L \sum_{i=1}^{S_{l-1}} \sum_{j=1}^{S_l} (w_{ij}^{(l)})^2$$

$$1) \text{sum}_j^l = \sum_{k=1}^{S_{l-1}} w_{kj}^l \cdot \sigma_k^{l-1} \quad \sigma_j^l = \frac{1}{1+e^{-\text{sum}_j^l}}$$



$$2) \frac{\partial E}{\partial w_{ij}^{(l)}} = \underbrace{\frac{\partial E}{\partial \sigma_j^l}}_{\text{our gradient descent value}} \cdot \underbrace{\frac{\partial \sigma_j^l}{\partial \text{sum}_j^l}}_{\sigma_j^l(1-\sigma_j^l)} \cdot \underbrace{\frac{\partial \text{sum}_j^l}{\partial w_{ij}^{(l)}}}_{\sigma_{j-1}^{l-1}} + \frac{\lambda}{m} w_{ii}^{(l)}$$

$$3) * \frac{\partial E}{\partial \sigma_j^l} = \sum_{p=1}^{S_{l+1}} \underbrace{\frac{\partial E}{\partial \sigma_p^{l+1}}}_{\text{use memoization here!}} \cdot \underbrace{\frac{\partial \sigma_p^{l+1}}{\partial \text{sum}_p^{l+1}}}_{\sigma_p^{l+1}(1-\sigma_p^{l+1})} \cdot \underbrace{\frac{\partial \text{sum}_p^{l+1}}{\partial \sigma_j^l}}_{w_{jp}^{l+1}}$$

at outer layer

$$\frac{\partial E}{\partial \sigma_j^L} = -\frac{y_j}{\sigma_j^L} - \frac{1-y_j}{1-\sigma_j^L}$$