

①  $a_1, \dots, a_n \in \mathbb{R}$ ,  $\sum (a_i - b)$  minimized  
(like loss func.)

(i) Analytical Solution

(Solving it by rigorous mathematical treatment to get exact solution)

$$\text{Total } = \frac{1}{2} \sum_{i=1}^n (a_i - b)^2$$

Total =  $\frac{1}{2} \mathbf{x}^T \mathbf{x}$  (can be written as loss function)

$$\frac{\partial \mathbf{x}^T \mathbf{x}}{\partial b} = 2 \mathbf{x}^T \frac{\partial \mathbf{x}}{\partial b}$$

General result in  
matrix calculus for  
any  $n \times m$  matrix

$$\frac{\partial \mathbf{a}}{\partial n} = \left[ \frac{\partial a_1}{\partial n}, \frac{\partial a_2}{\partial n}, \dots, \frac{\partial a_n}{\partial n} \right]^T, \quad \frac{\partial (\mathbf{a}^T \mathbf{a})}{\partial n} = \left[ \frac{\partial a_1^T}{\partial a_1}, \frac{\partial a_2^T}{\partial a_2}, \dots, \frac{\partial a_n^T}{\partial a_n} \right]$$

vector  
of scalars  
of scalars

$$\frac{\partial \mathbf{a}^T \mathbf{a}}{\partial b} = \left( \frac{\partial a_i}{\partial b} \right) \rightarrow \text{2nd order tensor, } \frac{\partial a_i}{\partial b}$$

vector =  $\sqrt{n}/n$  (divergence)

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{a}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (\mathbf{a}^T \mathbf{x}) = \mathbf{a}$$

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{AB}) = \mathbf{A} \frac{\partial \mathbf{x}}{\partial \mathbf{x}} \cdot \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}} \neq \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}}$$

matrices

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{x}^T \mathbf{Ax}) = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

tensor

$$\therefore \frac{\partial \mathbf{x}^T \mathbf{x}}{\partial b} = 0 \quad \left\{ \begin{array}{l} \text{condition for minimizing} \\ \text{loss} \end{array} \right.$$

$$\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial b} = 0 \iff \frac{\partial \mathbf{x}}{\partial b} = \left[ 1, 1, \dots, 1 \right]^T$$

$$\frac{\partial}{\partial b} \sum (a_i - b)^2 = 0$$

$$\sum a_i = \frac{1}{n} b \Rightarrow b = \frac{\sum a_i}{n}$$

$$(ii) p(n) \frac{2^{(n-b)}}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (n-b)^2\right)$$

(Normal)

Distr

Assume all  $x_i$  are drawn from the same distribution  
minimum form of maximum

for

$$y = w^T n + b + \epsilon, \text{ where } \epsilon \sim N(0, \sigma^2)$$

$$y - w^T n - b \sim \epsilon, y - w^T n \sim \epsilon$$

$$P(\epsilon) = P(y/n) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y - w^T n - b)^2\right)$$

cont. ~~prob (x2)~~

Given  $n$  variable = 'n', what is the prob  
of  $y$  variable =  $y$ :  $\rightarrow P(\epsilon)$

$$\Pr(y/x) = \prod_{i=1}^n P(y_i/x_i)$$

According to the principle of maximum likelihood,  $(w, b)$  are best. choice when it maximizes the likelihood over entire dataset.

$$-\log P(y/x) =$$

$$\sum_{i=1}^n \left[ \log \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) + \frac{(y_i - w^T n_i - b)^2}{2\sigma^2} \right]$$

$$= \left( \sum_{i=1}^n \frac{1}{2\sigma^2} (y_i - w^T n_i - b)^2 \right)$$

minimizing  $(w, b)$  solution

prob of  $y_i, n \rightarrow$

linear relation  $\rightarrow y_i \sim w^T n_i + b$

Likelihood is  $P(y/x)$

Ques: the assumption that ' $\sigma$ ' is fixed.

When minimizing  $-\log P(y/x)$ , the best  $(w, b)$ ,

(1) the first term has no effect,

(2) The second term is equivalent to minimizing the squared sum error, except for the multiplicative constant.

Minimizing  $\sum [x_i - b]^2$ ,  $x_i = y_i - w^T n_i$

$$(iii) \text{ Loss} = \sum_i (n_i - b)^2 = \sum_i j(n_i - b)^2$$

The same approach like (1i) can't be applied here, because the derivative of mod is not defined.

But I don't know, if it can be rigorously solved to get an optimal solution.

→ By intuition, it seems  $b = \frac{\sum n_i}{n}$  will give the minimum value, as

as the task is still to minimize  $\sum_j (n_i - b)^2$ , → for an unbiased term like  $(n_i - b)^2$ ,  $b$  is center (~~mean~~-value) should give the minimized loss

(iv)

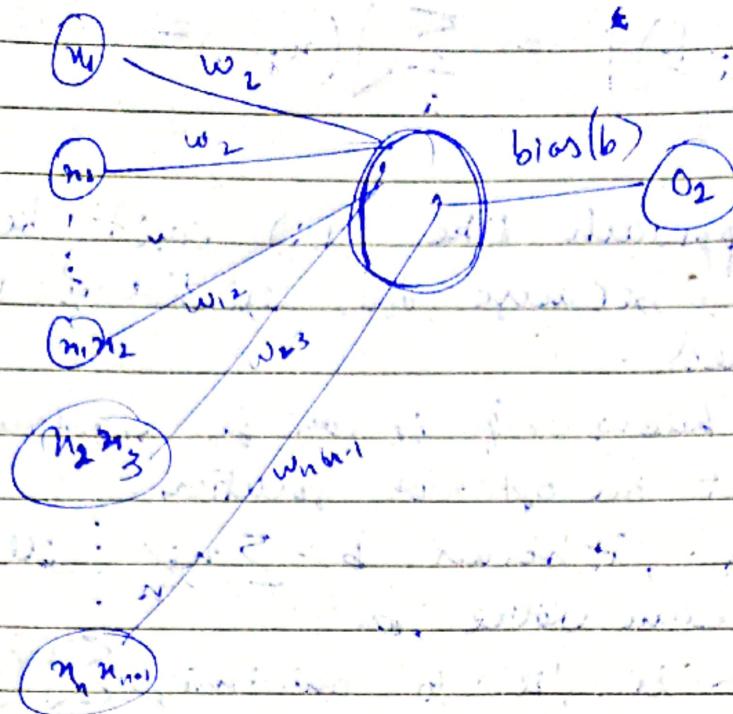
(2) Affine functions := linear function + constant  
 $\therefore$  (e.g.  $n^T w + b$ )

Linear functions:  $f(x) = m \cdot x + b$   
 on  $(n, 1)$   $\tilde{n} = \begin{pmatrix} n \\ 1 \end{pmatrix}$ ; higher dimensional space

$f(n) = n^T w + b$  with  $n \in \mathbb{R}^n$  (affine function)  
 $w \in (\underline{w_1, w_2, w_3, \dots})$

$f(n) = \tilde{n}^T w'$   $w' \in (\underline{w, b})^T$   
 (linear function)  $\tilde{n}^T \in (\underline{n, 1})$

(3)



(4)

Analytical solution, ... (min. Mean Squared Errors  $\rightarrow$ )  
Cost function

$$\partial_w \|y - Xw\|^2 = \frac{\partial}{\partial w} |y - X^T w - b|$$

$$\partial_w \|y - Xw\|^2 = 0 : z \vdash 2x^T(Xw - y)$$

$$X^T y = X^T X w$$

$$w^* = (X^T X)^{-1} X^T y$$

(optimized)  
solution

$\therefore$  unique  $w^*$  only for  $X^T X$  is invertible.

$\therefore X^T X$  is invertible  $\Rightarrow X$  has full rank,

$X \in \mathbb{R}^{n \times d}$  (Design Matrix)

Rank of a matrix : describes the dimension  
of the vector space generated by its rows or columns

In other words ; no. of independent columns.

basis vectors of the generated space.  
(each is  $n$ -dim)

$\therefore$  maximum dim. generated (Rank)  
 $\leq \min(n, m)$

(i) if  $X^T X$  is not invertible ( $\det(X^T X) = 0$ ), then no unique solution for  $w$  (as no. of solutions), which means, this algorithm would fail to provide the optimum solution for  $Wz(w, b)$ .

(ii) ① ~~May be~~, we can drop the dependent coordinate (as it is not helping in prediction of  $y$ ),

② If we add a small amount of coordinate-wise ( $n_1, n_2, \dots$ ) independent Gaussian Noise to all the entries of  $X$ . Then  $X^T X$  will be invertible with very high probability.

(iii) It is necessary to add Gaussian Noise over all the coordinates rather than just to a single datapoint ( $y \rightarrow y + \epsilon$ ).

$$\mathbb{E}(X) = \begin{bmatrix} n_1 + \epsilon_1 & n_2 + \epsilon_2 & \dots \\ n_1 + \epsilon_1 & n_2 + \epsilon_2 & \dots \\ n_1 + \epsilon_1 & n_2 + \epsilon_2 & \dots \end{bmatrix}$$

$\epsilon$  is each element would be different

$$\therefore \mathbb{E}(X^T X) = \mathbb{E}[(n_i + \epsilon_i)(n_j + \epsilon_j)]$$

$$= \mathbb{E}[n_i n_j + \mathbb{E}[n_i] \mathbb{E}[n_j] + \mathbb{E}[n_i] \mathbb{E}[\epsilon_j] + \mathbb{E}[\epsilon_i] \mathbb{E}[n_j]]$$

$$\Rightarrow n_i n_j + 0 \quad 0 \quad 0;$$

$$\mathbb{E}(\epsilon_i) = 0$$

$$\mathbb{E}(n_i) = \underbrace{n_i}_{\text{deterministic}}$$

## (iv) minibatch Stochastic Descent

$$(\hat{w}, \hat{b}) \leftarrow (w, b) - \eta \frac{\sum_{i \in B_t} \hat{l}_i^{(i)}(w, b)}{|B_t|}$$

• Shows Trajectory of min loss function

The Stochastic Gradient Descent doesn't rely on the invertibility of  $X^T X$  but instead iteratively updates the regression coefficients based on randomly selected subset of the data.

The presence of linearly dependent columns (rank  $r < n$ ) can cause the loss surface to have multiple minima. (because now two/more columns are dependent)

Does it make sense to reach the saddle points? → payed - qualitatively characterizing neural networks optimization problems → Good Fellows.

$$(5i) P(E) = \frac{1}{2} \exp(-|E|), \quad y = n^T w + b \\ P(y/x) = \frac{1}{2} \exp(-|y - n^T w - b|). \\ -\log(P(y/x)) = -\log(\frac{1}{2}) + |y - n^T w - b|.$$

(ii) Maybe, it solvable for  $w, b$ ,  
but since the derivative of Mod is undefined;  
therefore a closed-form solution is not possible.

- (iii) Minibatch = SGD algorithm  
 → initialize  $(w, b)$  with random variable  
 → choose hyperparameters  $(n, \beta)$

$$\rightarrow (w, b) \leftarrow (w, b) - \frac{n}{|\beta|} \sum_{i \in \beta} \partial_{(w, b)} l^{(i)}(w, b),$$

$$l(w, b) = -\log(P(y/x))$$

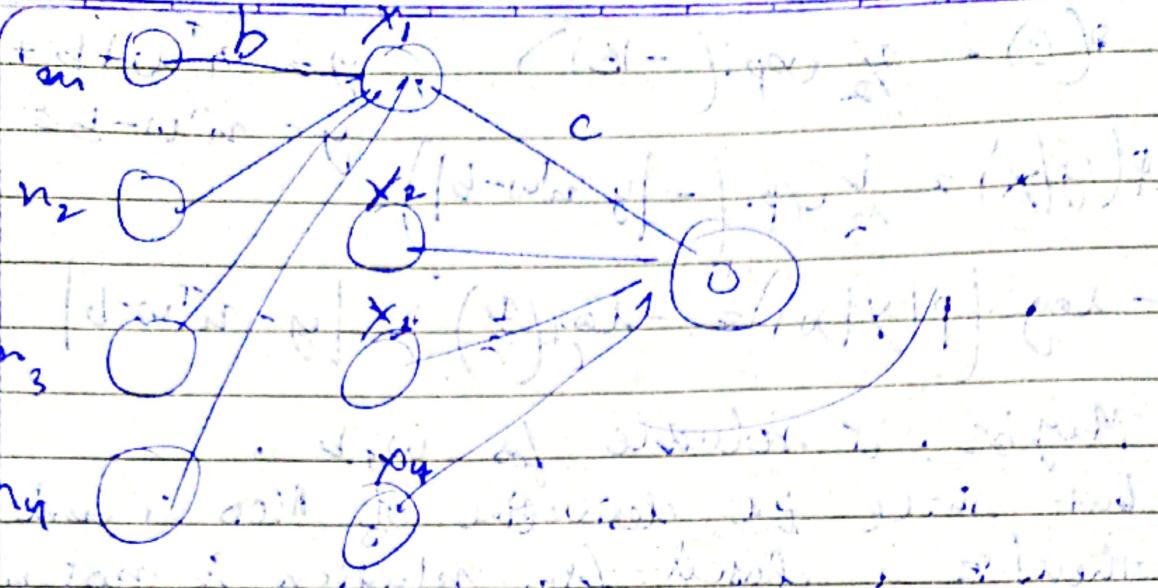
→ Stopping criterion (e.g.  $\Delta(w, b) \leq \epsilon$ )

Let's assume for computation, we have  
smoothed  $|y - n^T w - b|$  near stationary points,  
we would result in the slope near stationary points  
tending to infinity  $\rightarrow$  This could cause the  
algorithm to diverge.

Solutions -

- ① Use a variant of gradient descent that includes a regularization term, such as gradient descent with momentum or RMSprop.)
- ② By skipping a chunk near the stationary point
- ③ Redefining to little accurate but compatible function

(8):



$$O_2 = (w_{11}C_1 + w_{21}C_2 + w_{31}C_3 + w_{41}C_4) n_1 + \\ (w_{12}C_1 + w_{22}C_2 + w_{32}C_3 + w_{42}C_4) n_2$$

$$= C^T \sum_i w_i n_i$$

- This will just like a single-layers only.
- This can't model non-linear relations, to model such relations more no. of linear variables are needed rather than layers.

→ If large no. of variable coefficients may become very small or large if the ranges of the variable are way apart.

(75) (i) Additive Gaussian Noise:  $\epsilon$  is added to the function

The Gaussian Noise varies from  $-\infty$  to  $\infty$ , which means some house prices might be negative.

But excluding those negative prices (extreme values), still the normal dist., maybe used for Prices.

The house prices in general increase with time over years but  $\epsilon$  would have an expectation of '0'. (Another  $\epsilon$  could be added which would change with time but its expectation value would be also around '0' which would mean that we need to modify it to have a bias to instate the fluctuations)  $\rightarrow$  which just means gaussian noise is not appropriate & needs lot of modification & tweaking to fit the "houses price problem" well.

(iii) Black-Scholes Model  $\rightarrow$  option Pricing

(ii) Log price may be better, because only few values are allowed. Houses to increase after certain values.

(8) i) Setting of Apple in a discrete quantity which Gaussian Noise can't model.

ii)  $p(k|\lambda) = \lambda^k e^{-\lambda} / k!$  {Poisson Distribution  
mean =  $E(k)$ . median =  $\lambda$  = rate function}

$\rightarrow \sum_{k=0}^{\infty} k p(k|\lambda)$  is total, avg no. of events  
in captures

$= \sum_{k=0}^{\infty} \left( \frac{\lambda^k e^{-\lambda}}{k!} \right) e^{\lambda} = \sum_{k=0}^{\infty} \lambda^k = \text{distribution over counts}$

$$= \lambda e^{-\lambda} \cdot \sum_{k=1}^{\infty} \lambda^{k-1} = \lambda e^{-\lambda} (\lambda - 1)!$$

$$= \lambda e^{-\lambda} e^{\lambda} = \lambda$$

iii)  $-\log(p(k|\lambda)) = -\sum_{i=1}^n \log(p(k_i|\lambda))$

like the maximum likelihood problem

$$= n\lambda - \left( \sum_{i=1}^n k_i \right) \log \lambda + \sum_{i=1}^n \log k_i!$$

$$\text{Loss function: } L(\lambda) = n\lambda - \left( \sum_{i=1}^n k_i \right) \log \lambda$$

minimizing this value maximizes likelihood over entire dataset. (Loss function for Poisson Distribution)

(iv)  $\lambda \rightarrow \log \lambda = t$ ,  $\lambda(\lambda) = L(t)$   
the function has  
to remain same.

$$\lambda(t) = n e^t - \left( \sum_{i=1}^n k_i \right) t \cancel{\log \lambda}$$