

田 Joint probability:

- two independent event A and B

probability of A and B occurring, $P(A \cap B) = P(A) \times P(B)$
Both

- two dependent even A and B

probability of A and B occurring, $P(A \cap B) = P(A) \times P(B|A)$

田 Conditional probability:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

田 Total probability:

$$P(A) = \sum_{i=1}^n P(A|B_i) \times P(B_i)$$

田 Bayes Rule:

$$P(A|B) = \frac{P(B|A) P(A) \leftarrow \text{prior}}{P(B)}$$

田 Gaussian (Normal) Distribution:

$$N(x, \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

- modeling small random noise errors.

Maximum Likelihood Estimation:

distribution parameter

$$P(D|\theta) = \theta^3(1-\theta)^2$$

H, T, T, H, H

θ = Bernoulli

Distribution

$$L(\theta|D) = \log P(D|\theta)$$

$$\max_{\theta} = \theta^3(1-\theta)^2$$

$$\log(L) = 3\log\theta + 2\log(1-\theta)$$

$$\frac{d}{d\theta} \log L = \frac{3}{\theta} + \frac{2}{(1-\theta)} = 0$$

$$\Rightarrow \frac{3}{\theta} - \frac{2}{(1-\theta)} = 0$$

$$\theta^* = \frac{3}{5} \quad \text{or} \quad \hat{\theta}_{MLE} = \frac{m_H}{m_H + m_T}$$

parameter values that maximize the Likelihood data

MLE for Mean of Gaussian: $\hat{\mu}_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i$

MLE for variance of Gaussian: $\hat{\sigma}_{MLE}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2$

MAP:

$$\textcircled{1} \quad P(x|\theta) = \theta^m(1-\theta)^{n-m} \quad | \quad D = \{x_1, x_2, \dots, x_N\}$$

$$\textcircled{2} \quad \text{prior} = P(\theta) = \text{Beta}(\alpha, \beta) \quad | \quad B(\alpha, \beta) = \alpha-1, \beta-1$$

$$\textcircled{3} \quad \text{Posterior} \quad P(\theta|x) \propto P(x|\theta) \times P(\theta)$$

$$= \theta^m(1-\theta)^{n-m} \cdot \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

$$= \theta^{m+\alpha-1}(1-\theta)^{n-m+\beta-1}$$

For bigger sample, MAP \approx MLE,

$$\hat{\theta}_{MAP} = \frac{m_H + \alpha - 1}{m_H + \alpha - 1 + m_T + \beta - 1}$$

④ Linear Regression:

$$y(\vec{x}, \vec{w}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$
$$= w_0 + \sum_{j=1}^d w_j x_j$$

④ Notation Conventions:

$$y(\vec{x}, \vec{w}) = \sum_{j=0}^d w_j x_j$$
$$= \vec{w}^T \vec{x}$$

④ Squared Loss:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \left(\frac{x_n^T w - t_n}{\text{prediction}} \right)^2 =$$
$$\frac{1}{2} (\vec{x} \vec{w} - \vec{t})^T (\vec{x} \vec{w} - \vec{t})$$

④ Optimization:

minimize $E(w)$, find \vec{w} ,

$$\hat{w} = (\vec{x}^T \vec{x})^{-1} \vec{x}^T \vec{t}$$

④ Probability density function: (gaussian) Normal

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} (x-\mu)^2}$$

μ = mean

σ = variance

④ Likelihood of Linear Regression:

$$(2\pi\sigma^2)^{-\frac{N}{2}} e^{-\frac{1}{2\sigma^2} (\vec{t} - \vec{x} \vec{w})^T (\vec{t} - \vec{x} \vec{w})}$$

④ Ridge regression:

$$w = (I + X^T X)^{-1} X^T t$$

$$E(w) = (Xw - t)^T (Xw - t) + \lambda w^T w$$

if $\lambda \rightarrow 0$ overregularized

$\boxed{\lambda > 0}$ should be for regularization

if $\lambda \rightarrow \infty$ $\|w\| \geq 0$

underfitting

Polynomial Curve fitting:

$$x \in \mathbb{R}, y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_m x^m =$$

x = Non-linear

w = linear

$$\sum_{j=0}^m w_j x^j$$

General case: (Linear Basis Function)

$$y(x, w) = w_0 + w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_{n-1} \phi_{n-1}(x)$$

$$\phi(x) = \vec{x}$$

$$y(x, w) = \vec{w}^\top \phi(x)$$

Least squares:

$$w_{\text{ML}} = (\phi^\top \phi)^{-1} \phi^\top \vec{x}$$

B Generalization meaning in the context

of machine learning:

Good performance on any new,
unseen data.

B Underfitting = high bias

B overfitting = high variance

④ Logistic function (Sigmoid Function):

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

④ Logistic regression probabilistic function:

$$P(Y=0|x) = \frac{1}{1 + \exp(w_0 + \sum_i w_i x_i)}$$

$$P(Y=1|x) = \frac{\exp(w_0 + \sum_i w_i x_i)}{1 + \exp(w_0 + \sum_i w_i x_i)}$$

④ odds ratio:

$$\frac{P(Y=1|x)}{P(Y=0|x)} = \exp(w_0 + \sum_i w_i x_i)$$

$$\log \left(\frac{P(Y=1|x)}{P(Y=0|x)} \right) = w_0 + \sum_i w_i x_i$$

Linear function

④ For more than two classes:

$k < K$ = total number of class

d = total feature

$$P(Y=Y_k|x) = \frac{\exp(w_{k0} + \sum_{i=1}^d w_{ki} x_i)}{1 + \sum_{j=1}^{k-1} \exp(w_{j0} + \sum_{i=1}^d w_{ji} x_i)}$$

$$k = K \quad P(Y=Y_K|x) = \frac{1}{1 + \sum_{j=1}^{k-1} \exp(w_{j0} + \sum_{i=1}^d w_{ji} x_i)}$$

④ Training Logistic Regression:

$$l(w) = \ln \prod P(Y^j|x^j, w)$$

$$= \sum_j [y^j (w_0 + \sum_i w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i w_i x_i^j))]$$

④ $\lambda(w)$ is a concave function.

For example:

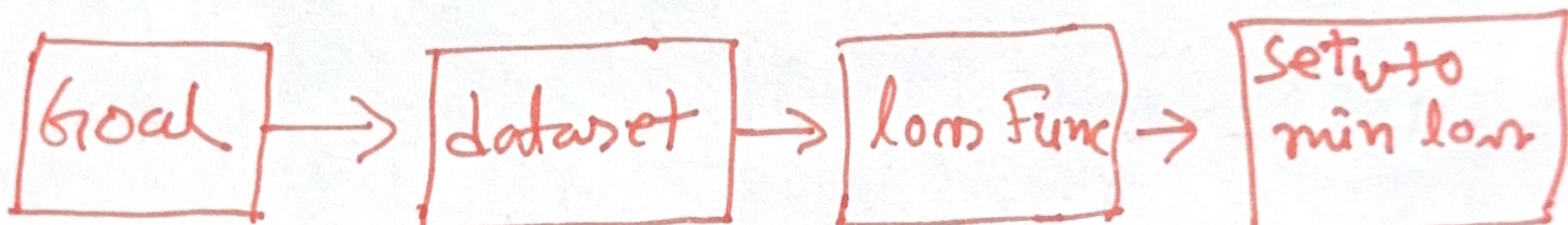
$$y = -9 - x^2 - y^2$$

coff of x^2 is negative.

④ No closed form solution.

④ easy to optimized (unique maximum)

④ Gradient descent:



learn w

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\} \quad \text{loss}_D(w) = \text{gradient descent}$$

$$\text{Argmin}_x f(x), \quad x_{t+1} = x_t - \eta f'(x_t)$$

or,

$$w \leftarrow w - \eta \nabla E_D(w)$$

④ Batch gradient:

$$\nabla E_D(w) = \left[\frac{d E_D(w)}{d w_0}, \dots, \frac{d E_D(w)}{d w_n} \right]$$

$$w \leftarrow w - \eta \nabla E_D(w)$$

④ Stochastic gradient: (Faster)

$$\nabla E_d(w) = \left[\frac{d E_d(w)}{d w_0}, \dots, \frac{d E_d(w)}{d w_n} \right]$$

$$w \leftarrow w - \eta \nabla E_d(w)$$

④ Gradient descent for logistic regression:

$$\begin{aligned} \frac{\partial \ell(w)}{\partial w_i} &= \sum_j x_i^j (y_j - p(y_j=1|x_i^j, w)) \\ &= \sum_j x_i^j \left(y_j - \frac{\exp(w_0 + \sum_{i=1}^j w_i x_i^j)}{1 + \exp(w_0 + \sum_{i=1}^j w_i x_i^j)} \right) \end{aligned}$$

↓ Ground truth ↓ prediction

④ Effect of Step Size:

$$w \leftarrow w - \eta \nabla E_p(w)$$

η = too large \Rightarrow oscillating

η = too small \Rightarrow slow convergence

④ Regularization:

$$w^* = \operatorname{argmax} \ln \sum_{j=1}^N p(y_j|x_i^j, w)$$

$$\text{non-regularized} \quad w_i^{t+1} \leftarrow w_i^t + \eta \sum_j x_i^j [y_j - p(y_j=1|x_i^j, w)]$$

regularized

$$w^* = \operatorname{argmax} \ln \sum_{j=1}^N p(y_j|x_i^j, w) - \frac{\lambda}{2} \frac{\|w\|^2}{\sum_{i=1}^N w_i^2}$$

$$w_i^{(t+1)} \leftarrow w_i^t + \eta \left[-\lambda w_i^t + \sum_j x_i^j [y_j - p(y_j=1|x_i^j, w)] \right]$$

1. Gradient ascent to maximize the cost function $l(w)$. [coz it's concave]
2. purpose of regularization in logistic Regression is

Soft margin

$$\text{Pb} \quad \frac{1}{2} \|w\|^2 + C \sum \xi_j$$

if margin ≥ 1 : Don't care
margin < 1 : pay linear penalty

$$\text{slack, } (w \cdot x + b)y_j \geq 1 - \xi_j$$

$$\rightarrow \xi_j = 1 - y_j(w \cdot x + b) \Rightarrow \text{Hinge loss}$$

$\rightarrow C = \text{Slack penalty } > 0$

Pb if $C = \infty$,
have to separate the data. **Overfitting**

if $C = 0$,
ignore data.

C = small, \rightarrow smooth and simple
decision boundary.

$\xi_j \geq 1$ if misclassified.

Pb multiple class.

$$(w_k, b_k)_{k=1,2,3}$$

$$y = \arg \max_k (w_k \cdot x + b_k)$$

C = chosen by cross-validation

■ $wx+b=0 \rightarrow$ be decision boundary

$$wx+b=+1$$

$$wx+b=-1$$

■ Goal: $\min w, b \rightarrow \max \gamma = \text{margin}$

for all data.

$$(wx+b)y_j \geq \gamma$$

■ projection on plane.

$$\underline{x_j} = \bar{x_j} + \underline{\lambda} \begin{bmatrix} w \\ \|w\| \end{bmatrix}$$

Projection \downarrow Distance

unit vector

for class:

$$x^+ = x^- + 2\lambda \frac{w}{\|w\|}$$

■ Unnormalized problem:

maximize γ, w, b

$$(wx_j+b)y_j \geq \gamma$$

■ Normalized problem:

$$\gamma = \frac{\|w\|}{w \cdot w} = \frac{1}{\sqrt{w \cdot w}} = \frac{1}{2\|w\|}$$

min \uparrow

$$(w \cdot x_j + b)y_j \geq 1 \quad \forall j \in \text{Dataset.}$$

min

④ Dual SVM - Linearly separable

$$L(w, b, \alpha) = \frac{1}{2} w \cdot w - \sum_j \alpha_j [(w \cdot x_j + b) y_j - 1]$$

$$\begin{array}{l} \max_{\alpha} \sum_j \alpha_j \\ \alpha \geq 0 \\ \sum_j \alpha_j y_j = 0 \end{array}$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w = \sum_j \alpha_j y_j x_j - 0$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_j \alpha_j y_j = 0 \quad \text{--- (1)}$$

Sub (1) and (1).

$$\sum \alpha_j - \frac{1}{2} \sum_{ij} \alpha_i \alpha_j y_i y_j \underline{x_i x_j}$$

$$\begin{array}{l} \sum_i \alpha_i y_i = 0 \\ \alpha_i \geq 0 \end{array} \quad \text{then Sol gives } \alpha_i$$

$$w =$$

④ Dual SVM - Sparsity.

Support vector if $\alpha_j \neq 0$.

④ kernel trick:

$$k(x_i, x_j) = \varphi(x_i)^\top \varphi(x_j)$$

④ Dual SVM - Non-separable

① Same or Linear.

$$c > \alpha_j > 0$$

④ Efficient dot product of poly.

$$k(u, v) = \phi(u) \cdot \phi(v) = (u \cdot v)^d$$

④ Derivation of kernel:

$$k(x, x') = \phi(x)^T \phi(x')$$

$k_{ii} = k(x_i, x_i)$ is positive semi-definite

④ SVM with kernel:

1. Choose feature and kernel

2. Solve dual SVM to get α_i

3. Compute w, b .

4. Classify wrt $(w \cdot \phi(x) + b)$

④ Avoid Overfitting by:

→ Setting C

→ Better kernel

→ Varying parameter of kernel

B) Boosting works well even if

$$\text{error} \leq 0.5$$

B) Highly sensitive to outliers. \rightarrow overfit

B) Weight update rule:

$$D_t(i) = \frac{1}{m}$$
$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha y_i h_t(x_i))}{Z_t}$$
$$Z_t = \sum_{i=1}^m D_t(i) \exp(-\alpha y_i h_t(x_i))$$

B) weight $\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right)$

B) training error,

$$\epsilon_t = \sum_{i=1}^m D_t \delta(h_t(x_i) \neq y_i)$$

$$+ \frac{1}{m} \sum_{i=1}^m \delta(h_t(x_i) \neq y_i) \leq \frac{1}{m} \sum_{i=1}^m \exp(-\gamma h_t(x_i))$$

$$h(x) = \sum_t \alpha_t h_t(x) \quad [\text{convex}]$$

B Boosting the margin

$$r(x_i) = y_i \frac{\alpha_1 h_1(x_i) + \dots + \alpha_t h_t(x_i)}{\alpha_1 + \alpha_2 + \dots + \alpha_t}$$

B Handler overfit By throwing all the classifiers have low variance, high bias.

② $k = \text{NN}:$

k too high \rightarrow overfitting

k too low \rightarrow underfitting

for a new data (x_1, x_2) ,

① Find the distance from all data points

② Define k , and find all the mindist. neighbours.

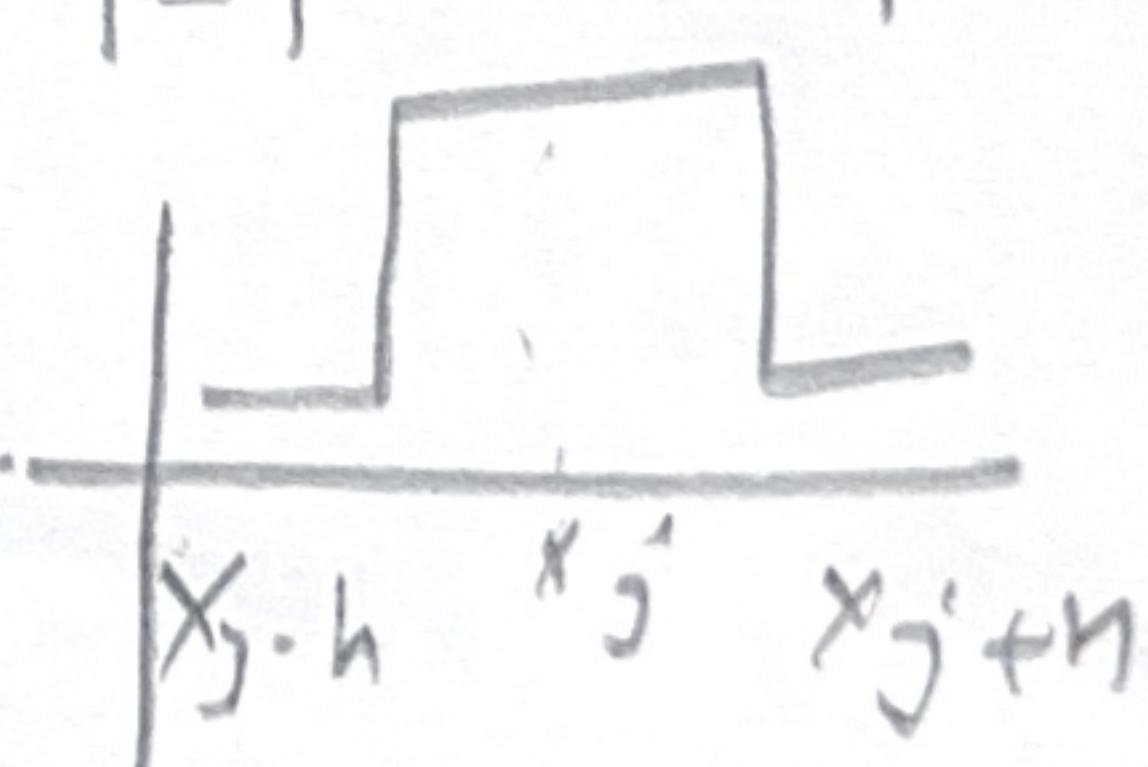
③ Select by vote.

③ kernel Regression:

$$\text{boxcar } k(x) = \frac{1}{2} I(x)$$



$$k\left(\frac{x_j - x}{h}\right)$$



h is low \rightarrow overfitting

h is high \rightarrow underfitting

Density Estimation:

$$\rho_i = \frac{n_i}{N\Delta'_i}$$

△ 0.04 overfit

△ 0.25 underfit

△ 0.08 perfect

more data = complex model

numbers of parameter \approx training data

Do k-means clustering:

- unsupervised,
- gets Density estimation
- groups on cluster
- sensitive to need choice. (poor convergence)

Algorithm:

① Random centroid.

② Classify each point,

$$c^{(t)}(j) \leftarrow \arg \min_{i=1 \dots k} \|u_i^{(t)} - x_j\|$$

③ Recenter centroid.

$$u_i^{(t+1)} \leftarrow \arg \min_u \sum_{j: c^{(t)}(j)=i} \|u - x_j\|$$

$u_i = \text{avg points (mean)}$

Cost function:

$$F(U, c) = \sum_{j=1}^m \|u_{c(j)} - x_j\|$$

So, min the distance

so we, fix U , optimize

$$F(U, c) = \sum_{i=1}^k \sum_{j: c(j)=i} \|u_i - x_j\|$$

fix c , optimize U

$i=1 \dots k$ is the number of centers.

j = datapoint
 x_j = data
 $c(j)$ = centroid of that data

④ Cent.

- Sensitive to outliers.
- only for vector data
- No method to choose number of clusters from cost function.

④ k-medoids:

- ① one center

$$c_i^{(t+1)} \leftarrow \underset{\substack{\text{argmin} \sum_{j=1}^n \|x_j - x_i\| \\ c(j) \neq k}}{}$$

④ Properties:

higher value of point

$$F(\mu, c) = \sum_{j=1}^n$$

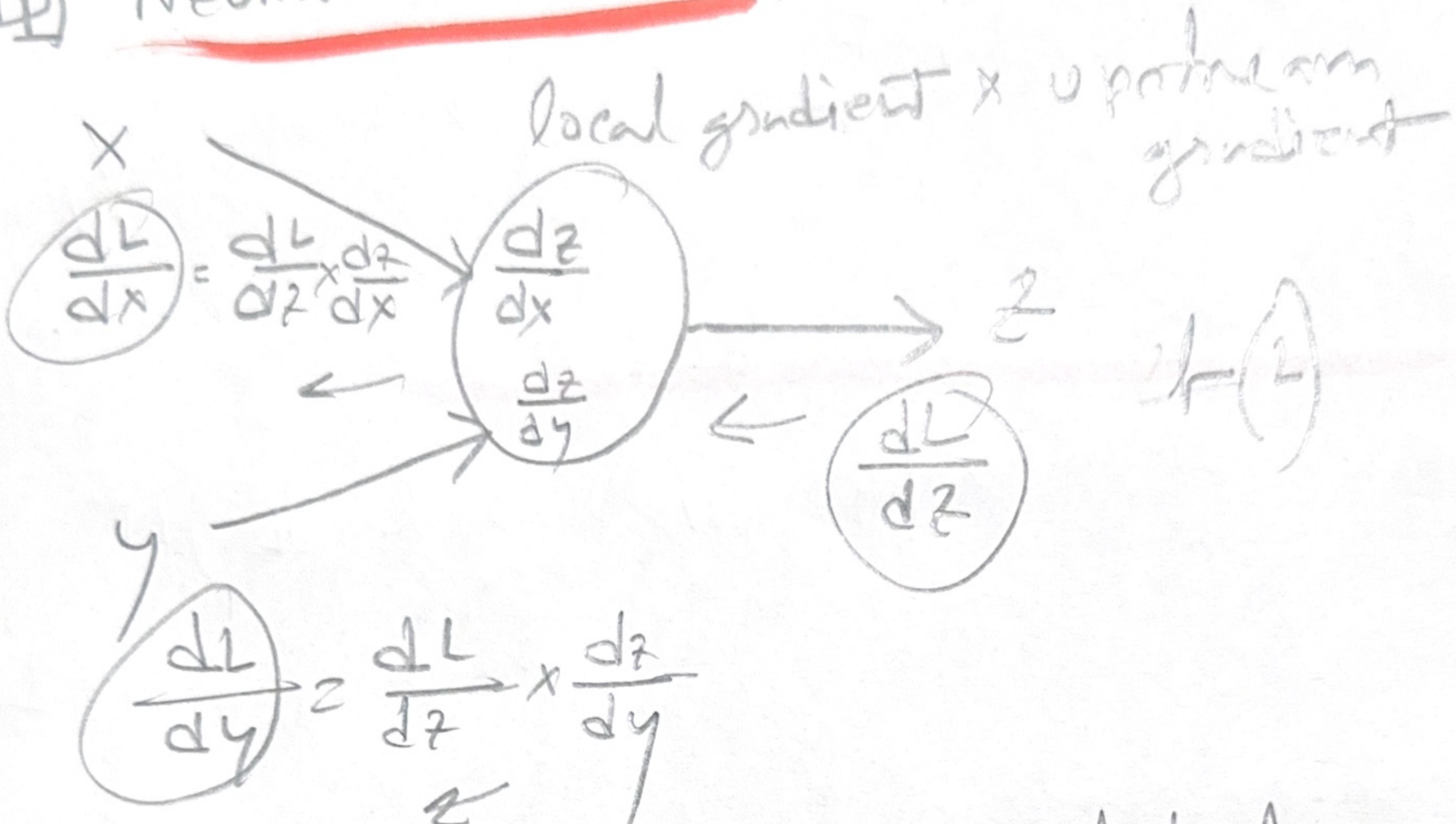
④ Soft k-means:

probability of point belongs to cluster

$$p(x, c_i) = \frac{\exp \{-d(x, c_i)\}}{\sum_j \exp \{-d(x, c_j)\}}$$

$$c_i = \frac{\sum_x x p(x, c_i)}{\sum_x p(x, c_i)} \quad \text{weighted center.}$$

④ Neural Network:



④ Activation in Layer Architecture:

$$t = w_1 x$$

$$t_2 = w_2 \max(0, w_1 x)$$

④ Activation function:

$$\text{Sigmoid } \sigma(x) = \frac{1}{1+e^{-x}}, \quad \frac{d\sigma(x)}{dx} = \sigma(x)(1-\sigma(x))$$

$$\text{ReLU } t(x) = \max(0, x), \quad \begin{cases} t'(x) = 0 & x \leq 0 \\ t'(x) = 1 & x > 0 \end{cases}$$

$$\text{LReLU } t(x) = \max(0, 0.01x, x), \quad \begin{cases} t'(x) = 0.01x & x \leq 0 \\ t'(x) = 1 & x > 0 \end{cases}$$

$$\text{ELU } t(x) = \begin{cases} x, & x > 0 \\ \alpha(\exp(x)-1), & x \leq 0 \end{cases} \quad t'(x) = 1 \quad x > 0$$

$$\text{Maxout } t(x) = \max(w_1^T b_1 + b_1 + w_2^T x_2 + b_2)$$

Q1 Output volume size: (conv)

$$w_2 = \frac{w_1 - F + 2P}{S} + 1 \times \text{Number of kernel/filter}$$

Q2 Number of parameter: (conv)

$$\text{filter} \times 3 \times (H \times W \times \text{channel} + 1) \times K$$

Common K = power of 2

torch.nn.conv2d

Common F=2, S=2

conv,

F=3, S=2

Input,

$$w_1 \times h_1 \times D_1$$

Output,

$$w_2 = \frac{w_1 - F + 2P}{S} + 1$$

$$h_2 = \frac{h_1 - F + 2P}{S} + 1$$

$D_2 = K$ = Number of filter

for, pooling

Input, $w_1 \times h_1 \times D_1$

$$w_2 = \frac{w_1 - F}{S} + 1$$

$$h_2 = \frac{h_1 - F}{S} + 1$$

$$D_2 = D_1$$